

CONTRIBUTORS: RAMAN ARORA, SANJEEV ARORA, JOAN BRUNA, NADAV COHEN, SIMON DU,

RONG GE, SURIYA GUNASEKAR, CHI JIN, JASON LEE, TENG YU MA, BEHNAM NEYSHABUR,

ZHAO SONG

# THEORY OF DEEP LEARNING



# Contents

1	<i>Basic Setup and some math notions</i>	15
1.1	<i>List of useful math facts</i>	16
1.1.1	<i>Probability tools</i>	16
1.1.2	<i>Singular Value Decomposition</i>	18
2	<i>Basics of Optimization</i>	19
2.1	<i>Gradient descent (GD)</i>	19
2.1.1	<i>Upperbound on the Taylor Expansion via Smoothness</i>	20
2.1.2	<i>Descent lemma for gradient descent</i>	20
2.2	<i>Stochastic gradient descent (SGD)</i>	21
2.3	<i>Accelerated Gradient Descent</i>	22
2.4	<i>Local Runtime Analysis of GD</i>	23
2.4.1	<i>Pre-conditioners</i>	24
2.5	<i>Correspondence of theory with practice</i>	25
3	<i>Note on overparametrized linear regression and kernel regression</i>	27
3.1	<i>Overparametrized least squares linear regression</i>	27
3.2	<i>Kernel least-squares regression</i>	28
4	<i>Note on Backpropagation and its Variants</i>	31
4.1	<i>Problem Setup</i>	31
4.1.1	<i>Multivariate Chain Rule</i>	33
4.1.2	<i>Naive feedforward algorithm (not efficient!)</i>	34

4.2	<i>Backpropagation (Linear Time)</i>	34
4.3	<i>Auto-differentiation</i>	35
4.4	<i>Notable Extensions</i>	36
4.4.1	<i>Hessian-vector product in linear time: Pearlmutter's trick</i>	37
5	<i>Advanced Optimization notions</i>	39
6	<i>Basics of generalization theory</i>	41
6.1	<i>Occam's razor formalized for ML</i>	41
6.1.1	<i>Motivation for generalization theory</i>	42
6.1.2	<i>Motivating example: Polynomial interpolation</i>	42
6.2	<i>Some simple upper bounds on generalization error</i>	43
6.3	<i>Data dependent complexity measures</i>	45
6.3.1	<i>Rademacher Complexity</i>	46
6.3.2	<i>Alternative Interpretation: Ability to correlate with random labels</i>	47
6.4	<i>Understanding limitations of the union-bound approach</i>	47
6.5	<i>A Compression-based framework</i>	49
6.5.1	<i>Example 1: Linear classifiers with margin</i>	51
6.5.2	<i>Example 2: Generalization bounds for deep nets using low rank approximations</i>	52
6.6	<i>PAC-Bayes bounds</i>	53
6.7	<i>Exercises</i>	55
7	<i>Tractable Landscapes for Nonconvex Optimization</i>	57
7.1	<i>Preliminaries and challenges in nonconvex landscapes</i>	58
7.2	<i>Cases with a unique global minimum</i>	59
7.2.1	<i>Generalized linear model</i>	60
7.2.2	<i>Alternative objective for generalized linear model</i>	61
7.3	<i>Symmetry, saddle points and locally optimizable functions</i>	62
7.4	<i>Case study: top eigenvector of a matrix</i>	64
7.4.1	<i>Characterizing all critical points</i>	64
7.4.2	<i>Finding directions of improvements</i>	66

8	<i>Escaping Saddle Points</i>	69
8.1	<i>Preliminaries</i>	69
8.2	<i>Perturbed Gradient Descent</i>	70
8.3	<i>Saddle Points Escaping Lemma</i>	72
8.3.1	<i>Improve or Localize</i>	73
8.3.2	<i>Bounding the Width of the Stuck Region</i>	73
9	<i>Algorithmic Regularization</i>	77
9.1	<i>Linear models in regression: squared loss</i>	78
9.1.1	<i>Geometry induced by updates of local search algorithms</i>	80
9.1.2	<i>Geometry induced by parameterization of model class</i>	83
9.2	<i>Matrix factorization</i>	83
9.3	<i>Linear Models in Classification</i>	83
9.3.1	<i>Gradient Descent</i>	84
9.3.2	<i>Steepest Descent</i>	85
9.4	<i>Homogeneous Models with Exponential Tailed Loss</i>	89
9.5	<i>Induced bias in function space</i>	91
10	<i>Ultra-wide Neural Networks and Neural Tangent Kernels</i>	93
10.1	<i>Evolution of the trained net</i>	94
10.1.1	<i>Behavior in the infinite limit</i>	95
10.2	<i>NTK: Simple 2-layer example</i>	96
10.3	<i>Explaining Optimization and Generalization of Ultra-wide Neural Networks via NTK</i>	99
10.3.1	<i>Understanding Generalization in 2-layer setting</i>	101
10.4	<i>NTK formula for Multilayer Fully-connected Neural Network</i>	102
10.5	<i>NTK in Practice</i>	104
10.6	<i>Exercises</i>	105
11	<i>Inductive Biases due to Algorithmic Regularization</i>	107
11.1	<i>Matrix Sensing</i>	108
11.1.1	<i>Gaussian Sensing Matrices</i>	110
11.1.2	<i>Matrix Completion</i>	113

11.2	<i>Deep neural networks</i>	115	
11.3	<i>Landscape of the Optimization Problem</i>	118	
11.3.1	<i>Implicit bias in local optima</i>	120	
11.3.2	<i>Landscape properties</i>	122	
11.4	<i>Role of Parametrization</i>	128	
11.4.1	<i>Related Work</i>	128	
12	<i>Unsupervised learning: Distribution Learning</i>	129	
12.1	<i>Possible goals of unsupervised learning</i>	129	
12.2	<i>Training Objective for Learning Distributions: Log Likelihood</i>	131	
12.2.1	<i>Notion of goodness for distribution learning</i>	131	
12.3	<i>Variational method</i>	133	
12.4	<i>Autoencoders and Variational Autoencoder (VAEs)</i>	134	
12.4.1	<i>Training VAEs</i>	135	
12.5	<i>Normalizing Flows</i>	136	
13	<i>Generative Adversarial Nets</i>	139	
13.1	<i>Distance between Distributions</i>	139	
13.2	<i>Introducing GANs</i>	140	
13.2.1	<i>Game-theoretic interpretation and implications for trainings</i>	142	
13.3	<i>Do GANs learn the distribution?</i>	142	
13.3.1	<i>Experimental verification: Birthday Paradox Test</i>	143	
14	<i>Self-supervised Learning</i>	145	
15	<i>Adversarial Examples</i>	147	
16	<i>Examples of Theorems, Proofs, Algorithms, Tables, Figures</i>	149	
16.1	<i>Example of Theorems and Lemmas</i>	149	
16.2	<i>Example of Long Equation Proofs</i>	149	

16.3 *Example of Algorithms* 151

16.4 *Example of Figures* 152

16.5 *Example of Tables* 153

16.6 *Exercise* 153

*Bibliography* 155





## List of Figures

- 2.1 Convex and Nonconvex Functions in two variables. For nonconvex functions GD will reach a stationary point, where gradient is zero. (Figure from kdnuggets.org) 21
- 2.2 How train and test loss behaved during SGD on PreResNet32 on CIFAR10 (50k datapoints). The test loss was estimated at various steps via a fixed held-out dataset. Initial learning rate was 1, and it was reduced by a factor of 10 at epoch 80 and epoch 300. (An epoch is a full pass over the training dataset, where the dataset has been randomly partitioned into batches for use in SGDs—in this case the batches were of size 128.) Note the complicated relationship between train and test loss, in particular, a slight rise in test loss can happen even when training loss is flat or going down. 26
- 2.3 **Edge of stability phenomenon.** When doing GD (as opposed to SGD) with a small learning rate  $\eta$ , the smoothness is observed to rise to  $2/\eta$ , whereupon during iterations one starts seeing loss increases as well as decreases, with a long term downward trend. No theoretical explanation is known as of now. The authors use “sharpness” instead of smoothness, which actually makes some sense because higher  $L$  corresponds to a more uneven landscape. 26
  
- 4.1 Why it suffices to compute derivatives with respect to nodes. 32
- 4.2 Multivariate chain rule: derivative with respect to node  $z$  can be computed using weighted sum of derivatives with respect to all nodes that  $z$  feeds into. 33
- 4.3 Vector version of above 36
  
- 7.1 Obstacles for nonconvex optimization. From left to right: local minimum, saddle point and flat region. 59
  
- 8.1 **Left:** Perturbation ball in 3D and “thin pancake” shape stuck region. **Right:** Perturbation ball in 2D and “narrow band” stuck region under gradient flow 74

- 9.1 **Steepest descent w.r.t  $\|\cdot\|_{4/3}$** : the global minimum to which steepest descent converges to depends on  $\eta$ . Here  $w_0 = [0, 0, 0]$ ,  $w_{\|\cdot\|}^* = \arg \min_{\psi \in G} \|\psi\|_{4/3}$  denotes the minimum norm global minimum, and  $w_{\eta \rightarrow 0}^\infty$  denotes the solution of infinitesimal SD with  $\eta \rightarrow 0$ . Note that even as  $\eta \rightarrow 0$ , the expected characterization does not hold, i.e.,  $w_{\eta \rightarrow 0}^\infty \neq w_{\|\cdot\|}^*$ . 82
- 10.1 Convergence rate vs. projections onto eigenvectors of the kernel matrix. 100
- 10.2 Generalization error vs. complexity measure. 101
- 11.1 Optimization landscape (top) and contour plot (bottom) for a single hidden-layer linear autoencoder network with one dimensional input and output and a hidden layer of width  $r = 2$  with dropout, for different values of the regularization parameter  $\lambda$ . Left: for  $\lambda = 0$  the problem reduces to squared loss minimization, which is rotation invariant as suggested by the level sets. Middle: for  $\lambda > 0$  the global optima shrink toward the origin. All local minima are global, and are equalized, i.e. the weights are parallel to the vector  $(\pm 1, \pm 1)$ . Right: as  $\lambda$  increases, global optima shrink further. 121
- 12.1 Visualization of Pearson's Crab Data as mixture of two Gaussians. (Credit: MIX homepage at McMaster University.) 130
- 12.2 Autoencoder defined using a density distribution  $p(h, x)$ , where  $h$  is the latent feature vector corresponding to visible vector  $x$ . The process of computing  $h$  given  $x$  is called "encoding" and the reverse is called "decoding." In general applying the encoder on  $x$  followed by the decoder would not give  $x$  again, since the composed transformation is a sample from a distribution. 130
- 12.3 Faces in the top row were produced by a VAE based method and those in the second row by RealNVP using normalizing flows. VAE is known for producing blurry images. RealNVP's output is much better, but still has visible artifacts. 137
- 12.4 Convex and Nonconvex Functions in two variables. For nonconvex functions GD will reach a stationary point, where gradient is zero. (Figure from kdnuggets.org) 138
- 16.1 A chasing sequence 152

## List of Tables

16.1 We ignore the  $O$  for simplicity. The  $\ell_\infty/\ell_2$  is the strongest possible guarantee, with  $\ell_2/\ell_2$  coming second,  $\ell_2/\ell_1$  third and exactly  $k$ -sparse being the weaker. We also note that all [RV08, CGV13, Bou14, HR16] obtain improved analyses of the Restricted Isometry property; the algorithm is suggested and analyzed (modulo the RIP property) in [BD08]. The work in [HIKP12] does not explicitly state the extension to the  $d$ -dimensional case, but can easily be inferred from the arguments. [HIKP12, IK14, Kap16, KVZ19] work when the universe size in each dimension are powers of 2.



# *Introduction*

This monograph discusses the emerging theory of deep learning. It is based upon a graduate seminar taught at Princeton University in Fall 2019 in conjunction with a Special Year on Optimization, Statistics, and Machine Learning at the Institute for Advanced Study.



# 1

## *Basic Setup and some math notions*

This Chapter introduces the basic nomenclature. Training/test error, generalization error etc. <<Tengyu notes: Todos: Illustrate with plots: a typical training curve and test curve

Mention some popular architectures (feed forward, convolutional, pooling, resnet, densenet) in a brief para each. >>

We review the basic notions in statistical learning theory.

- A space of possible data points  $\mathcal{X}$ .
- A space of possible labels  $\mathcal{Y}$ .
- A joint probability distribution  $\mathcal{D}$  on  $\mathcal{X} \times \mathcal{Y}$ . We assume that our training data consist of  $n$  data points

$$(x^{(1)}, y^{(1)}), \dots, (x^{(n)}, y^{(n)}) \stackrel{\text{i.i.d.}}{\sim} \mathcal{D},$$

each drawn independently from  $\mathcal{D}$ .

- Hypothesis space:  $\mathcal{H}$  is a family of hypotheses, or a family of predictors. E.g.,  $\mathcal{H}$  could be the set of all neural networks with a fixed architecture:  $\mathcal{H} = \{h_\theta\}$  where  $h_\theta$  is neural net that is parameterized by parameters  $\theta$ .
- Loss function:  $\ell : (\mathcal{X} \times \mathcal{Y}) \times \mathcal{H} \rightarrow \mathbb{R}$ .
  - E.g., in binary classification where  $\mathcal{Y} = \{-1, +1\}$ , and suppose we have a hypothesis  $h_\theta(x)$ , then the logistic loss function for the hypothesis  $h_\theta$  on data point  $(x, y)$  is

$$\ell((x, y), \theta) = \frac{1}{1 + \exp(-yh_\theta(x))}.$$

- Expected loss:

$$L(h) = \mathbb{E}_{(x, y) \sim \mathcal{D}} [\ell((x, y), h)].$$

Recall  $\mathcal{D}$  is the data distribution over  $\mathcal{X} \times \mathcal{Y}$ .

- Training loss (also known as empirical risk):

$$\widehat{L}(h) = \frac{1}{n} \sum_{i=1}^n \ell \left( (x^{(i)}, y^{(i)}), h \right),$$

where  $(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(n)}, y^{(n)})$  are  $n$  training examples drawn i.i.d. from  $\mathcal{D}$ .

- Empirical risk minimizer (ERM):  $\widehat{h} \in \arg \min_{h \in \mathcal{H}} \widehat{L}(h)$ .
- Regularization: Suppose we have a regularizer  $R(h)$ , then the regularized loss is

$$\widehat{L}_\lambda(h) = \widehat{L}(h) + \lambda R(h)$$

◀Suriya notes: Misc notations: gradient, hessian, norms>>

## 1.1 List of useful math facts

Now we list some useful math facts.

### 1.1.1 Probability tools

In this section we introduce the probability tools we use in the proof. Lemma 1.1.4, 1.1.5 and 1.1.6 are about tail bounds for random scalar variables. Lemma 1.1.7 is about cdf of Gaussian distributions. Finally, Lemma 1.1.8 is a concentration result on random matrices.

**Lemma 1.1.1** (Markov's inequality). *If  $x$  is a nonnegative random variable and  $t > 0$ , then the probability that  $x$  is at least  $t$  is at most the expectation of  $x$  divided by  $t$ :*

$$\Pr[x \geq t] \leq \mathbb{E}[x]/t.$$

**Lemma 1.1.2** (Chebyshev's inequality). *Let  $x$  denote a nonnegative random variable and  $t > 0$ , then*

$$\Pr[|x - \mathbb{E}[x]| \geq t] \leq \text{Var}[x]/t^2.$$

Next we present some concentration bounds regarding sum of independent random variables. The rule of thumb underlying concentration bounds is the *Central Limit Theorem*.

**Theorem 1.1.3** (Central Limit Thm, informal). *If  $X_1, X_2, \dots, X_n$  are independent random variables of mean  $\mu_1, \mu_2, \dots, \mu_n$  and variances  $\sigma_1^2, \sigma_2^2, \dots, \sigma_n^2$  then as  $n$  gets larger,  $\sum_i X_i$  behaves like the normal distribution  $\mathcal{N}(\sum_i \mu_i, \sum_i \sigma_i^2)$ .*



Concentration bounds are quantitative versions of this and work also in settings where  $p_i$ 's and  $\sigma_i$ 's could depend on  $n$ , the total number of variables. But in many setting the CLT is a good rule of thumb.

**Lemma 1.1.4** (Chernoff bound [Che52]). *Let  $X = \sum_{i=1}^n X_i$ , where  $X_i = 1$  with probability  $p_i$  and  $X_i = 0$  with probability  $1 - p_i$ , and all  $X_i$  are independent. Let  $\mu = \mathbb{E}[X] = \sum_{i=1}^n p_i$ . Then*

1.  $\Pr[X \geq (1 + \delta)\mu] \leq \exp(-\delta^2\mu/3), \forall \delta > 0;$
2.  $\Pr[X \leq (1 - \delta)\mu] \leq \exp(-\delta^2\mu/2), \forall 0 < \delta < 1.$

**Lemma 1.1.5** (Hoeffding bound [Hoe63]). *Let  $X_1, \dots, X_n$  denote  $n$  independent bounded variables in  $[a_i, b_i]$ . Let  $X = \sum_{i=1}^n X_i$ , then we have*

$$\Pr[|X - \mathbb{E}[X]| \geq t] \leq 2 \exp\left(-\frac{2t^2}{\sum_{i=1}^n (b_i - a_i)^2}\right).$$

**Lemma 1.1.6** (Bernstein inequality [Ber24]). *Let  $X_1, \dots, X_n$  be independent zero-mean random variables. Suppose that  $|X_i| \leq M$  almost surely, for all  $i$ . Then, for all positive  $t$ ,*

$$\Pr\left[\sum_{i=1}^n X_i > t\right] \leq \exp\left(-\frac{t^2/2}{\sum_{j=1}^n \mathbb{E}[X_j^2] + Mt/3}\right).$$

**Lemma 1.1.7** (Anti-concentration of Gaussian distribution). *Let  $X \sim N(0, \sigma^2)$ , that is, the probability density function of  $X$  is given by  $\phi(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{x^2}{2\sigma^2}}$ . Then*

$$\Pr[|X| \leq t] \in \left(\frac{2}{3} \frac{t}{\sigma}, \frac{4}{5} \frac{t}{\sigma}\right).$$

**Lemma 1.1.8** (Matrix Bernstein, Theorem 6.1.1 in [Tro15]). *Consider a finite sequence  $\{X_1, \dots, X_m\} \subset \mathbb{R}^{n_1 \times n_2}$  of independent, random matrices with common dimension  $n_1 \times n_2$ . Assume that*

$$\mathbb{E}[X_i] = 0, \forall i \in [m] \quad \text{and} \quad \|X_i\| \leq M, \forall i \in [m].$$

*Let  $Z = \sum_{i=1}^m X_i$ . Let  $\text{Var}[Z]$  be the matrix variance statistic of sum:*

$$\text{Var}[Z] = \max\left\{\left\|\sum_{i=1}^m \mathbb{E}[X_i X_i^\top]\right\|, \left\|\sum_{i=1}^m \mathbb{E}[X_i^\top X_i]\right\|\right\}.$$

*Then*

$$\mathbb{E}[\|Z\|] \leq (2\text{Var}[Z] \cdot \log(n_1 + n_2))^{1/2} + M \cdot \log(n_1 + n_2)/3.$$

*Furthermore, for all  $t \geq 0$ ,*

$$\Pr[\|Z\| \geq t] \leq (n_1 + n_2) \cdot \exp\left(-\frac{t^2/2}{\text{Var}[Z] + Mt/3}\right).$$

*explain these in a para*

A useful shorthand will be the following: If  $y_1, y_2, \dots, y_m$  are independent random variables each having mean 0 and taking values in  $[-1, 1]$ , then their average  $\frac{1}{m} \sum_i y_i$  behaves like a Gaussian variable with mean zero and variance at most  $1/m$ . In other words, the probability that this average is at least  $\epsilon$  in absolute value is at most  $\exp(-\epsilon^2 m)$ .

### 1.1.2 Singular Value Decomposition

TBD.

## 2

# Basics of Optimization

This chapter sets up the basic analysis framework for gradient-based optimization algorithms and discuss how it applies to deep learning. The algorithms work well in practice; the question for theory is to analyse them and give recommendations for practice. This has proved harder, and in recent years it has become clearer that classical ways of thinking about optimization may not match well with phenomena encountered in deep learning.

The basic conceptual framework in optimization builds upon simple Taylor approximation (Equation 2.1) of the loss function and thus relies upon derivatives (of various orders) of the loss function.

«Suriya notes: To ground optimization to our case, we can also mention that  $f$  is often of the either the ERM or stochastic optimization form  $L(w) = \sum l(w; x, y)$  - it might also be useful to mention that outside of this chapter, we typically use  $f$  as an alternative for  $h$  to denote a function computed»

### 2.1 Gradient descent (GD)

Suppose we wish to minimize a continuous function  $f(w)$  over  $\mathbb{R}^d$ .

$$\min_{w \in \mathbb{R}^d} f(w).$$

The gradient descent (GD) algorithm is

$$\begin{aligned} w_0 &= \text{initialization} \\ w_{t+1} &= w_t - \eta \nabla f(w_t) \end{aligned}$$

where  $\eta$  is called *step size* or *learning rate*. The choice of  $\eta$  is important and a main subject in the rest of the Chapter.

One motivation or justification of the GD is that the update direction  $-\nabla f(w_t)$  is the steepest descent direction locally. Consider the

Taylor expansion at a point  $w_t$

$$f(w) = f(w_t) + \underbrace{\langle \nabla f(w_t), w - w_t \rangle}_{\text{linear in } w} + \underbrace{\frac{1}{2}(w - w_t)^T \nabla^2 f(w_t)(w - w_t)}_{\text{quadratic in } w} + \dots \quad (2.1)$$

here  $\nabla^2(f)$  is the matrix of 2nd order derivatives called *Hessian*. Its  $(i, j)$  entry is  $\partial^2 f / \partial w_i \partial w_j$ . Note that it is a symmetric matrix.

Suppose we drop the higher-order term and only optimize the first order approximation within a neighborhood of  $w_t$

$$\begin{aligned} \arg \min_{w \in \mathbb{R}^d} & f(w_t) + \langle \nabla f(w_t), w - w_t \rangle \\ \text{s.t.} & \|w - w_t\|_2 \leq \epsilon \end{aligned}$$

**Problem 2.1.1.** Show that the optimizer of the program above is equal to  $w + \delta$  where  $\delta = -\alpha \nabla f(w_t)$  for some positive scalar  $\alpha$ .

In other words, to locally minimize the first order approximation of  $f(\cdot)$  around  $w_t$ , we should move towards the direction  $-\nabla f(w_t)$ .

The classic back-propagation algorithm (Chapter 4) is used to efficiently compute the gradient of the loss. Note that today's deep nets use nonlinear activations, notably ReLU, that make the function computed by the net non-differentiable. However, this differentiability is of the mild sort and does not appear to be an issue in practice.

### 2.1.1 Upperbound on the Taylor Expansion via Smoothness

The most basic analysis of training speed of GD involves the smoothness of the loss function.

**Definition 2.1.2** (*L-smooth*). A function  $f$  is *L-smooth* in a domain if for every  $w$  in the domain all eigenvalues of  $\nabla^2 f(w)$  lie in the interval  $[-L, L]$ .

**Problem 2.1.3.** Prove that if  $f$  is *L-smooth* then

$$f(w) \leq f(w_t) + \langle \nabla f(w_t), w - w_t \rangle + \frac{L}{2} \|w - w_t\|_2^2 \quad (2.2)$$

### 2.1.2 Descent lemma for gradient descent

The following says that with gradient descent and small enough learning rate, the function value always decreases unless the gradient at the iterate is zero. (Points where gradient is zero are called *stationary points*.)

**Lemma 2.1.4** (*Descent Lemma*). Suppose  $f$  is *L-smooth*. Then, if  $\eta < 1/L$ , we have

$$f(w_{t+1}) \leq f(w_t) - \frac{\eta}{2} \cdot \|\nabla f(w_t)\|_2^2$$

The proof uses the Taylor expansion. The main idea is that even using the upper provided by equation (2.2) suffices.

*Proof.* We have that

$$\begin{aligned} f(w_{t+1}) &= f(w_t - \eta \nabla f(w_t)) \\ &\leq f(w_t) + \langle \nabla f(w_t), -\eta \nabla f(w_t) \rangle + \frac{L}{2} \|\eta \nabla f(w_t)\|_2^2 \\ &= f(w_t) - (\eta - \eta^2 L/2) \|\nabla f(w_t)\|_2^2 \\ &\leq f(w_t) - \frac{\eta}{2} \cdot \|\nabla f(w_t)\|_2^2, \end{aligned}$$

where the second step follows from Eq. (2.2), and the last step follows from  $\eta \leq 1/L$ .  $\square$

We've shown GD stops making progress when the gradient  $\nabla$  becomes zero. Is this good enough?

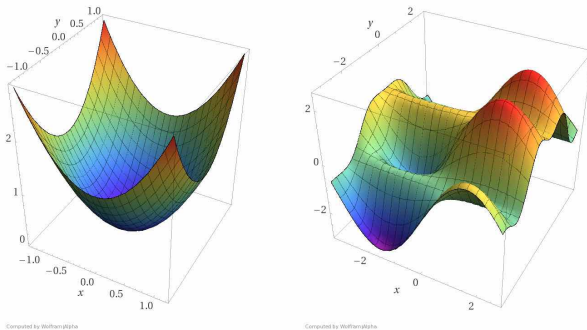


Figure 2.1: Convex and Nonconvex Functions in two variables. For nonconvex functions GD will reach a stationary point, where gradient is zero. (Figure from kdnuggets.org)

The loss function for deep learning is non-convex when the network has more than one layer. Thus GD is not guaranteed to produce a global optimum. Nevertheless, the solutions it finds in practice attain fairly low —near zero—value of the objective. (Recall that the loss function is usually nonnegative.) Elsewhere in the book this property of GD is explained in some concrete settings.

## 2.2 Stochastic gradient descent (SGD)

SGD is a very practical variant of gradient descent for large datasets. Recall that

$$\widehat{L}(h) = \frac{1}{n} \sum_{i=1}^n \ell((x^{(i)}, y^{(i)}), h).$$

Computing the gradient  $\nabla \widehat{L}(h)$  scales linearly in  $n$ , the size of the training dataset. Stochastic gradient descent (SGD) estimates the gradient by sampling a small number of training datapoints and

computing the average. By usual sampling theorems, the gradient estimate approaches true gradient as the sample size grows.

**The updates:** We simplify the notations a bit for ease of exposition. We consider optimizing the function

$$\frac{1}{n} \sum_{i=1}^n f_i(w)$$

So here  $f_i$  corresponds to  $\ell((x^i, y^{(i)}), h)$  in the statistical learning setting. At each iteration  $t$ , the SGD algorithm first samples  $i_1, \dots, i_B$  uniformly from  $[n]$ , and then computes the estimated gradient using the samples:

$$g_S(w) = \frac{1}{B} \sum_{k=1}^B \nabla f_{i_k}(w_t)$$

Here  $S$  is a shorthand for  $\{i_1, \dots, i_B\}$ . The SGD algorithm updates the iterate by

$$w_{t+1} = w_t - \eta \cdot g_S(w_t).$$

Note that if the learning rate  $\eta$  is very small, then the parameters will not change much over a sequence of updates, and so SGD would tend to be similar to (full) GD. However, when  $\eta$  is not too small and batch size is small, the gradient estimates from batches are highly noisy estimates of the true gradient. One would imagine this means SGD is worse than GD, but in practice SGD tends to do much better than GD. (Of course, being efficient, SGD allows far more iterations in the same computational budget.) Later chapters will cover theories that try to account for superiority of SGD over GD.

### 2.3 Accelerated Gradient Descent

The basic version of accelerated gradient descent algorithm is called heavy-ball algorithm. It has the following update rule:

$$w_{t+1} = w_t - \eta \nabla f(w_t) + \beta(w_{t+1} - w_t)$$

Here  $\beta(w_{t+1} - w_t)$  is the so-called momentum term. The motivation and the origin of the name of the algorithm comes from that it can be viewed as a discretization of the second order ODE:

$$\ddot{w} + a\dot{w} + b\nabla f(w) = 0$$

Another equivalent way to write the algorithm is

$$\begin{aligned} u_t &= -\nabla f(w_t) + \beta u_{t-1} \\ w_{t+1} &= w_t + \eta u_t \end{aligned}$$

Exercise: verify the two forms of the algorithm are indeed equivalent.

Another variant of the heavy-ball algorithm is due to Nesterov

$$\begin{aligned} u_t &= -\nabla f(w_t + \beta \cdot (u_t - u_{t-1})) + \beta \cdot u_{t-1}, \\ w_{t+1} &= w_t + \eta \cdot u_t. \end{aligned}$$

One can see that  $u_t$  stores a weighed sum of the all the historical gradient and the update of  $w_t$  uses all past gradient. This is another interpretation of the accelerate gradient descent algorithm

Nesterov gradient descent works similarly to the heavy ball algorithm empirically for training deep neural networks. It has the advantage of stronger worst case guarantees on convex functions. Both of the two algorithms can be used with stochastic gradient, but little is know about the theoretical guarantees about stochastic accelerate gradient descent.

## 2.4 Local Runtime Analysis of GD

When the iterations of GD are near a local minimum, the behavior of gradient descent is clearer because the function can be locally approximated by a quadratic function. In this section, we assume for simplicity that we are optimizing a convex quadratic function, and get some insight on how the curvature of the function influences the convergence of the algorithm.

We use gradient descent to optimize

$$\min_w \frac{1}{2} w^\top A w$$

where  $A \in \mathbb{R}^{d \times d}$  is a positive semidefinite matrix, and  $w \in \mathbb{R}^d$ .

Remark: w.l.o.g, we can assume that  $A$  is a diagonal matrix. **Diagonalization is a fundamental idea in linear algebra.** Suppose  $A$  has singular vector decomposition  $A = U \Sigma U^\top$  where  $\Sigma$  is a diagonal matrix. We can verify that  $w^\top A w = \hat{w}^\top \Sigma \hat{w}$  with  $\hat{w} = U^\top w$ . In other words, in a difference coordinate system defined by  $U$ , we are dealing with a quadratic form with a diagonal matrix  $\Sigma$  as the coefficient. Note the diagonalization technique here is only used for analysis.

Therefore, we assume that  $A = \text{diag}(\lambda_1, \dots, \lambda_d)$  with  $\lambda_1 \geq \dots \geq \lambda_d$ . The function can be simplified to

$$f(w) = \frac{1}{2} \sum_{i=1}^d \lambda_i w_i^2$$

The gradient descent update can be written as

$$x \leftarrow w - \eta \nabla f(w) = w - \eta \Sigma w$$

Here we omit the subscript  $t$  for the time step and use the subscript for coordinate. Equivalently, we can write the per-coordinate update rule

$$w_i \leftarrow w_i - \eta \lambda_i w_i = (1 - \lambda_i \eta) w_i$$

Now we see that if  $\eta > 2/\lambda_i$  for some  $i$ , then the absolute value of  $w_i$  will blow up exponentially and lead to an instable behavior. Thus, we need  $\eta \lesssim \frac{1}{\max \lambda_i}$ . Note that  $\max \lambda_i$  corresponds to the smoothness parameter of  $f$  because  $\lambda_1$  is the largest eigenvalue of  $\nabla^2 f = A$ . This is consistent with the condition in Lemma 2.1.4 that  $\eta$  needs to be small.

Suppose for simplicity we set  $\eta = 1/(2\lambda_1)$ , then we see that the convergence for the  $w_1$  coordinate is very fast — the coordinate  $w_1$  is halved every iteration. However, the convergence of the coordinate  $w_d$  is slower, because it's only reduced by a factor of  $(1 - \lambda_d/(2\lambda_1))$  every iteration. Therefore, it takes  $O(\lambda_d/\lambda_1 \cdot \log(1/\epsilon))$  iterations to converge to an error  $\epsilon$ . The analysis here can be extended to general convex function, which also reflects the principle that:

The condition number is defined as  $\kappa = \sigma_{\max}(A)/\sigma_{\min}(A) = \lambda_1/\lambda_d$ .

It governs the convergence rate of GD.

«Tengyu notes: add figure»

### 2.4.1 Pre-conditioners

From the toy quadratic example above, we can see that it would be more optimal if we can use a different learning rate for different coordinate. In other words, if we introduce a learning rate  $\eta_i = 1/\lambda_i$  for each coordinate, then we can achieve faster convergence. In the more general setting where  $A$  is not diagonal, we don't know the coordinate system in advance, and the algorithm corresponds to

$$w \leftarrow w - A^{-1} \nabla f(w)$$

In the even more general setting where  $f$  is not quadratic, this corresponds to the Newton's algorithm

$$w \leftarrow w - \nabla^2 f(w)^{-1} \nabla f(w)$$

Computing the hessian  $\nabla^2 f(w)$  can be computational difficult because it scales quadratically in  $d$  (which can be more than 1 million in practice). Therefore, approximation of the hessian and its inverse is used:

$$w \leftarrow w - \eta Q(w) \nabla f(w)$$



where  $Q(w)$  is supposed to be a good approximation of  $\nabla^2 f(w)$ , and sometimes is referred to as a pre-conditioner. In practice, often people first approximate  $\nabla^2 f(w)$  by a diagonal matrix and then take its inverse. E.g., one can use  $\text{diag}(\nabla f(w)\nabla f(w^\top))$  to approximate the Hessian, and then use the inverse of the diagonal matrix as the pre-conditioner.

«Tengyu notes: more on adagrad?»

## 2.5 Correspondence of theory with practice

Now we describe how the above theory compares with reality.

*Setting learning rate* The above theory shows how to set learning rate using smoothness parameter. Unfortunately, it is not easy to estimate the smoothness parameter for the loss function over the *entire* space of parameter vectors. For a particular parameter vector  $w$  however it is possible to estimate the maximum eigenvalue of the Hessian  $H = \nabla^2(f)$  at  $w$ . This uses the *power method*, which starts with a gaussian unit vector  $u$  and repeatedly computes  $u \leftarrow Hu/\|Hu\|_2$ . (Each iteration is efficiently implemented thanks to the Hessian-vector product computation described in Chapter 4.) This method is called the power method because it effectively amounts to computing  $H^t x/\|H^t x\|_2$ , which can be easily checked to converge to a vector that is a combination of the eigenvectors corresponding to the largest eigenvalue (in absolute value) of  $H$ . In particular, if  $v$  is the final vector,  $v^T H v$  would be a good approximation to the smoothness.

Of course, this only computes the smoothness parameter for a particular parameter vector  $w$ . As  $w$  changes, the smoothness can change too. Recomputing smoothness frequently would be computationally expensive. In practice the learning rate is set heuristically to some value. If GD does not lower the loss for a few iterations then learning rate is reduced by a small factor, like 2 or 5. In later chapters we will revisit the topic of learning rates.

*Edge of stability phenomenon.* The exposition of learning rates given above is canonical in classical optimization theory—it takes smoothness  $L$  as given and describes how learning rate must be set less than  $2/L$  to ensure consistent decrease in the loss. A recent paper <sup>1</sup> gives evidence that in deep nets the cart appears before the horse, so to speak. In other words, if we set the learning rate to some small  $\eta$ , the smoothness *adjusts* quickly to around  $2/\eta$ . This “edge of stability” phase appears to be important for good final performance.

<sup>1</sup> Jeremy Cohen, Simran Kaur, Yuanzhi Li, J Zico Kolter, and Ameet Talwalkar. Gradient descent on neural networks typically occurs at the edge of stability. *ICLR*, 2021

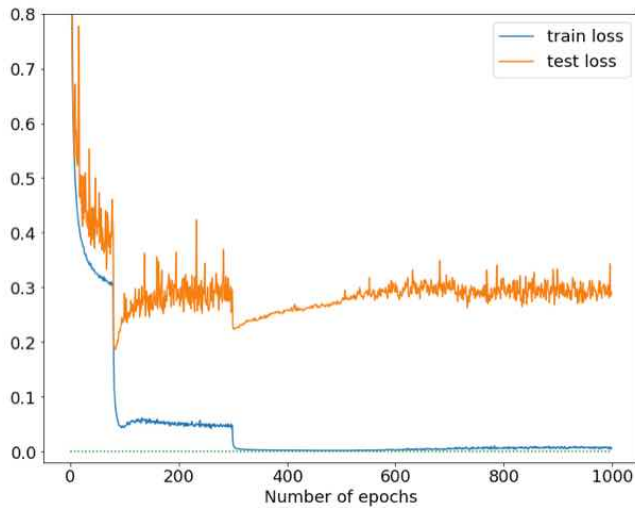


Figure 2.2: How train and test loss behaved during SGD on PreResNet32 on CIFAR10 (50k datapoints). The test loss was estimated at various steps via a fixed held-out dataset. Initial learning rate was 1, and it was reduced by a factor of 10 at epoch 80 and epoch 300. (An epoch is a full pass over the training dataset, where the dataset has been randomly partitioned into batches for use in SGDs—in this case the batches were of size 128.) Note the complicated relationship between train and test loss, in particular, a slight rise in test loss can happen even when training loss is flat or going down.

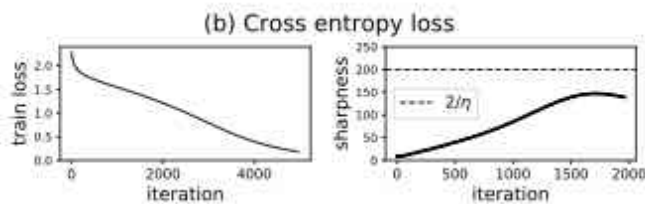


Figure 2.3: **Edge of stability phenomenon.** When doing GD (as opposed to SGD) with a small learning rate  $\eta$ , the smoothness is observed to rise to  $2/\eta$ , whereupon during iterations one starts seeing loss increases as well as decreases, with a long term downward trend. No theoretical explanation is known as of now. The authors use “sharpness” instead of smoothness, which actually makes some sense because higher  $L$  corresponds to a more uneven landscape.

### 3

## *Note on overparametrized linear regression and kernel regression*

This section analyzes gradient descent for a very classic model: least-squares linear regression. The problem is convex, and optimization does work well. We are interested primarily in the underdetermined version, where one has infinitely many zero-loss solutions and the interesting question is: what does gradient descent find? The analysis also extends to Kernel least squares regression.

Though classical, these analyses are the starting point for efforts (in later chapters) to understand over-parametrized deep nets, which also have an abundance of low cost solutions, and we wish to understand which ones are found by gradient descent and related algorithms.

### *3.1 Overparametrized least squares linear regression*

As in Chapter 1 we assume our training data consist of  $n$  data points, each drawn independently from a distribution  $\mathcal{D}$ ,

$$(x^{(1)}, y^{(1)}), \dots, (x^{(n)}, y^{(n)}) \stackrel{\text{i.i.d.}}{\sim} \mathcal{D}.$$

Here  $x^{(i)} \in \mathbb{R}^d$  and  $y^{(i)} \in \mathbb{R}$ . It will be convenient to write the  $\ell_2$  loss

$$\hat{L}(w) = \frac{1}{2} \sum_i (x^{(i)} \cdot w - y^{(i)})^2,$$

using matrix notation as  $\hat{L}(w) = \frac{1}{2} \|Xw - y\|^2$  where  $X$  is the matrix whose rows are the  $x^{(i)}$ 's,  $w$  is a column vector and  $y$  is the column vector whose  $i$ th entry is  $y^{(i)}$ . We're interested in the case where  $x^{(i)}$ 's are independent and  $d > n$ . Then the loss has infinitely many minimizers, all of which attain zero training loss. What does gradient descent find?

For simplicity, initialize gradient descent with starting point  $w_0 = 0$ . The gradient at any  $w$  is  $\nabla \widehat{L}(w) = X^T(Xw - y)$ . Thus gradient descent with learning rate  $\eta$  gives the following trajectory

$$w_{t+1} = w_t - \eta X^T(Xw_t - y) \quad (3.1)$$

$$= (I_{d \times d} - \eta X^T X)w_t + \eta X^T y \quad (3.2)$$

$$= \left( \sum_{j=0}^t (I_{d \times d} - \eta X^T X)^j \right) \eta X^T y \quad (3.3)$$

Assuming  $\eta$  was small enough by trial and error, specifically,  $\eta < 1/\lambda_{\max}(X^T X)$ , the infinite series as  $t \rightarrow \infty$  in the above equation converges to <sup>1</sup>

$$\lim_{t \rightarrow \infty} w_t = (X^T X)^\dagger X^T y \quad (3.4)$$

$$= X^T (X X^T)^{-1} y \quad (3.5)$$

<sup>1</sup> We're using that  $\sum_{i \geq 0} A^i = (I - A)^\dagger$  when the largest eigenvalue of positive semidefinite matrix  $A$  is less than 1 and  $Z^\dagger$  denotes the pseudo inverse of  $Z$ . Further, the second equality in eq. (3.5) can be verified by using the SVD of  $X$ .

which of course is the famous *pseudo-inverse* solution for overdetermined systems of linear equations. This solution is also the minimizing  $\ell_2$ -norm solution that fits the data:  $\arg \min_{w \in \mathbb{R}^d} \|w\|_2$  s.t.  $Xw = y$ .

### 3.2 Kernel least-squares regression

Kernel models involve a new representation for the data space  $\mathcal{X}$ , which represents point  $x \in \mathcal{X}$  as  $\phi(x)$  where  $\phi$  is a mapping from  $\mathbb{R}^d$  to a suitable Hilbert space known as the “Reproducing Kernel Hilbert Space” or RKHS. This means that the inner product  $\phi(x) \cdot \phi(y)$  is well-defined for every  $x_1, x_2 \in \mathcal{X}$ . It is customary to denote the inner product as  $K(x_1, x_2)$ , which is called the kernel function. The function class of interest are linear models over the transformed features  $h_w(x) = \phi(x) \cdot w$ . A plethora of useful kernels are known in mathematics and data science.

In data science the key property needed is that the inner product  $K(x_1, x_2)$  be efficiently computable. In fact in practice researchers design the kernel by starting with this property of efficient computability and never consider the underlying representation  $\phi(\cdot)$  because it plays no role in the training.<sup>2</sup> This property ensures one can solve the following kernel regression efficiently.

$$\ell(w) = \frac{1}{2} \sum_i (\phi(x)^{(i)} \cdot w - y^{(i)})^2.$$

While problematic at first sight—because  $\phi(x)^{(i)}$  is an infinite dimensional vector—but we quickly realize it is over-parametrized regression problem in disguise, where the data matrix  $XX^T$  now turns into

<sup>2</sup> For example, if the datapoints are unit vectors in  $\mathbb{R}^d$  then the Laplace kernel uses  $K(x_1, x_2) = \exp(-\|x_1 - x_2\|_2) = \exp(\sqrt{1 - 2x_1 \cdot x_2})$ . The reader should try to find the  $\phi(\cdot)$  function corresponding to this kernel. (Hint: look at the Taylor expansion of  $K(\cdot)$ .)

an  $n \times n$  gram matrix  $G$  where  $G_{ij} = \phi(x^{(i)}) \cdot \phi(x^{(j)}) = K(x^{(i)}, x^{(j)})$   
 <<Suriya notes: Edited to correct  $G = XX^T$  and not  $X^T X$  (3.4 equation is also equal to  $X^T(XX^T)^{-1}y$  which gives the gram-matrix computation for kernel-regression)>>. In fact computing  $G$  allows performing gradient descent without explicitly computing  $\phi(x^{(i)})$ . Expression (3.5) shows gradient descent ends with a classifier  $h(x) = \lim_{t \rightarrow \infty} w_t \cdot \phi(x)$  that maps an input point  $x \in \mathcal{X}$  to

$$h(x) = z^T G^{-1}y \quad (3.6)$$

where  $z$  is the column vector whose  $i$ th coordinate is  $K(x, x_i)$ . Alternatively, denoting  $\alpha = G^{-1}y \in \mathfrak{R}^n$ , the solution in eq. (3.6) is alternatively viewed as a weighted combinations of kernel evaluations at training points,

$$h(x) = \sum_i \alpha_i K(x, x_i). \quad (3.7)$$

Expression in eq. (3.7) is equivalent to minimizing the  $\ell_2$  norm of  $w$  while fitting the kernel regression objective, which viewed in the function space corresponds to the minimum RKHS norm solution wrt the kernel  $K$ .



## 4

# *Note on Backpropagation and its Variants*

Throughout the book we rely on computing the gradient of the loss with respect to model parameters. For deep nets, this computation is done with Backpropagation, a simple algorithm that uses the chain rule of calculus. For convenience we describe this more generally as a way to compute the sensitivity of the output of a neural network to all of its parameters, namely,  $\partial f / \partial w_i$ , where  $f$  is the output and  $w_i$  is the  $i$ th parameter. Here *parameters* can be edge weights or biases associated with nodes or edges of the network. Versions of this basic algorithm have been apparently independently rediscovered several times from 1960s to 1980s in several fields. This chapter introduces this algorithms as well as some advanced variants involving not just the gradient but also the Hessian.

In most of the book, the quantity of interest is the gradient of the training loss. But the above phrasing—computing gradient of the output with respect to the inputs—is fully general since one can simply add a new output node to the network that computes the training loss from the old output. Then the quantity of interest is indeed the gradient of this new output with respect to network parameters.

The importance of backpropagation derives from its efficiency. Assuming node operations take unit time, the running time is *linear*, specifically,  $O(\text{Network Size}) = O(V + E)$ , where  $V$  is the number of nodes in the network and  $E$  is the number of edges. As in many other settings in computer science—for example, sorting numbers—the naive algorithm would take quadratic time, and that would be hugely inefficient or even infeasible for today’s large networks.

### *4.1 Problem Setup*

Backpropagation applies only to acyclic networks with directed edges. (It can be heuristically applied to networks with cycles, as sketched later.) Without loss of generality, acyclic networks can be

visualized as being structured in numbered layers, with nodes in the  $t + 1$ th layer getting all their inputs from the outputs of nodes in layers  $t$  and earlier. We use  $f \in \mathbb{R}$  to denote the output of the network. In all our figures, the input of the network is at the bottom and the output on the top.

Our exposition uses the notion  $\partial f / \partial u$ , where  $f$  is the output and  $u$  is a node in the net. This means the following: suppose we cut off all the incoming edges of the node  $u$ , and fix/clamp the current values of all network parameters. Now imagine changing  $u$  from its current value. This change may affect values of nodes at higher levels that are connected to  $u$ , and the final output  $f$  is one such node. Then  $\partial f / \partial u$  denotes the rate at which  $f$  will change as we vary  $u$ . (Aside: Readers familiar with the usual exposition of back-propagation should note that there  $f$  is the training error and this  $\partial f / \partial u$  turns out to be exactly the "error" propagated back to on the node  $u$ .)

**Claim 4.1.1.** *To compute the desired gradient with respect to the parameters, it suffices to compute  $\partial f / \partial u$  for every node  $u$ .*

*Proof.* Follows from direct application of chain rule and we prove it by picture, namely Figure 4.1. Suppose node  $u$  is a weighted sum of the nodes  $z_1, \dots, z_m$  (which will be passed through a non-linear activation  $\sigma$  afterwards). That is, we have  $u = w_1 z_1 + \dots + w_m z_m$ . By Chain rule, we have

$$\frac{\partial f}{\partial w_1} = \frac{\partial f}{\partial u} \cdot \frac{\partial u}{\partial w_1} = \frac{\partial f}{\partial u} \cdot z_1.$$

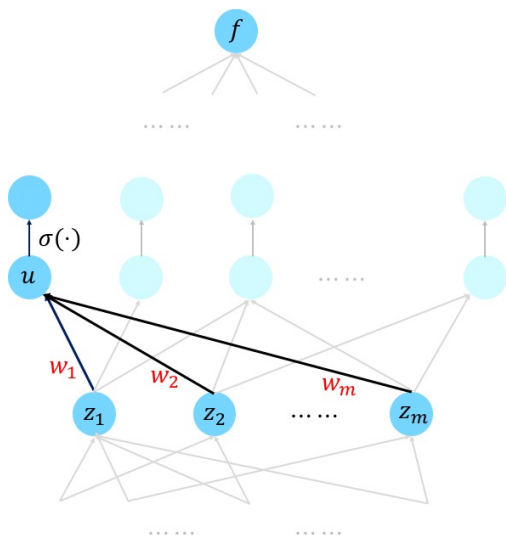


Figure 4.1: Why it suffices to compute derivatives with respect to nodes.

Hence, we see that having computed  $\partial f / \partial u$  we can compute  $\partial f / \partial w_1$ , and moreover this can be done locally by the endpoints of



the edge where  $w_1$  resides.  $\square$

#### 4.1.1 Multivariate Chain Rule

Towards computing the derivatives with respect to the nodes, we first recall the multivariate Chain rule, which handily describes the relationships between these partial derivatives (depending on the graph structure).

Suppose a variable  $f$  is a function of variables  $u_1, \dots, u_n$ , which in turn depend on the variable  $z$ . Then, multivariate Chain rule says that

$$\frac{\partial f}{\partial z} = \sum_{j=1}^n \frac{\partial f}{\partial u_j} \cdot \frac{\partial u_j}{\partial z}.$$

To illustrate, in Figure 4.2 we apply it to the same example as we used before but with a different focus and numbering of the nodes.

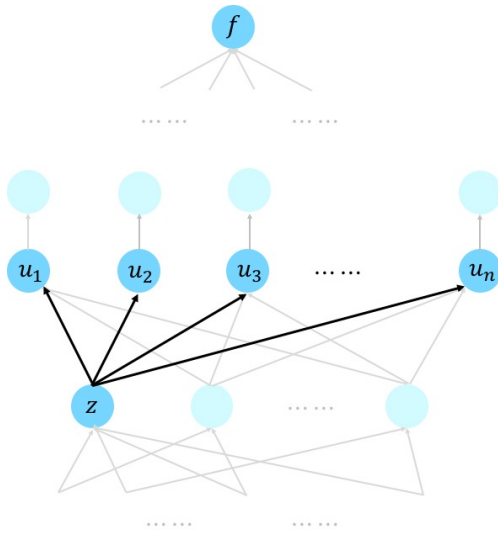


Figure 4.2: Multivariate chain rule: derivative with respect to node  $z$  can be computed using weighted sum of derivatives with respect to all nodes that  $z$  feeds into.

We see that given we've computed the derivatives with respect to all the nodes that is above the node  $z$ , we can compute the derivative with respect to the node  $z$  via a weighted sum, where the weights involve the local derivative  $\partial u_j / \partial z$  that is often easy to compute. This brings us to the question of how we measure running time. For book-keeping, we assume that

*Basic assumption:* If  $u$  is a node at level  $t + 1$  and  $z$  is any node at level  $\leq t$  whose output is an input to  $u$ , then computing  $\frac{\partial u}{\partial z}$  takes unit time on our computer.

### 4.1.2 Naive feedforward algorithm (not efficient!)

It is useful to first point out the naive quadratic time algorithm implied by the chain rule. Most authors skip this trivial version, which we think is analogous to teaching sorting using only quicksort, and skipping over the less efficient bubblesort.

The naive algorithm is to compute  $\partial u_i / \partial u_j$  for every pair of nodes where  $u_i$  is at a higher level than  $u_j$ . Of course, among these  $V^2$  values (where  $V$  is the number of nodes) are also the desired  $\partial f / \partial u_i$  for all  $i$  since  $f$  is itself the value of the output node.

This computation can be done in feedforward fashion. If such value has been obtained for every  $u_j$  on the level up to and including level  $t$ , then one can express (by inspecting the multivariate chain rule) the value  $\partial u_\ell / \partial u_j$  for some  $u_\ell$  at level  $t + 1$  as a weighted combination of values  $\partial u_i / \partial u_j$  for each  $u_i$  that is a direct input to  $u_\ell$ . This description shows that the amount of computation for a fixed  $j$  is proportional to the number of edges  $E$ . This amount of work happens for all  $j \in V$ , letting us conclude that the total work in the algorithm is  $O(VE)$ .

## 4.2 Backpropagation (Linear Time)

The more efficient backpropagation, as the name suggests, computes the partial derivatives in the reverse direction. Messages are passed in one wave backwards from higher number layers to lower number layers. (Some presentations of the algorithm describe it as dynamic programming.)

---

### Algorithm 1 Backpropagation

---

The node  $u$  receives a message along each outgoing edge from the node at the other end of that edge. It sums these messages to get a number  $S$  (if  $u$  is the output of the entire net, then define  $S = 1$ ) and then it sends the following message to any node  $z$  adjacent to it at a lower level:

$$S \cdot \frac{\partial u}{\partial z}$$


---

Clearly, the amount of work done by each node is proportional to its degree, and thus overall work is the sum of the node degrees. Summing all node degrees ends up double-counting each edge, and thus the overall work is  $O(\text{Network Size})$ .

To prove correctness, we prove the following:

**Lemma 4.2.1.** *At each node  $z$ , the value  $S$  is exactly  $\partial f / \partial z$ .*

*Proof.* Follows from simple induction on depth.

*Base Case:* At the output layer this is true, since  $\partial f / \partial f = 1$ .

*Inductive step:* Suppose the claim was true for layers  $t + 1$  and higher and  $u$  is at layer  $t$ , with outgoing edges go to some nodes  $u_1, u_2, \dots, u_m$  at levels  $t + 1$  or higher. By inductive hypothesis, node  $z$  indeed receives  $\frac{\partial f}{\partial u_j} \times \frac{\partial u_j}{\partial z}$  from each of  $u_j$ . Thus by Chain rule,

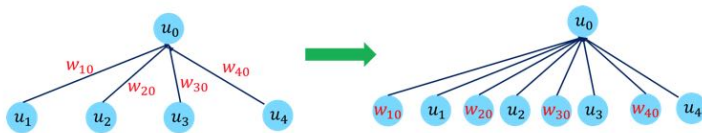
$$S = \sum_{i=1}^m \frac{\partial f}{\partial u_i} \frac{\partial u_i}{\partial z} = \frac{\partial f}{\partial z}.$$

This completes the induction and proves the Main Claim.  $\square$

### 4.3 Auto-differentiation

Since the exposition above used almost no details about the network and the operations that the node perform, it extends to every computation that can be organized as an acyclic graph whose each node computes a differentiable function of its incoming neighbors. This observation underlies many auto-differentiation packages found in deep learning environments: they allow computing the gradient of the output of such a computation with respect to the network parameters.

We first observe that Claim 4.1.1 continues to hold in this very general setting. This is without loss of generality because we can view the parameters associated to the edges as also sitting on the nodes (actually, leaf nodes). This can be done via a simple transformation to the network; for a single node it is shown in the picture below; and one would need to continue to do this transformation in the rest of the networks feeding into  $u_1, u_2, \dots$  etc from below.



Then, we can use the messaging protocol to compute the derivatives with respect to the nodes, as long as the local partial derivative can be computed efficiently. We note that the algorithm can be implemented in a fairly modular manner: For every node  $u$ , it suffices to specify (a) how it depends on the incoming nodes, say,  $z_1, \dots, z_n$  and (b) how to compute the partial derivative times  $S$ , that is,  $S \cdot \frac{\partial u}{\partial z_j}$ .

**Extension to vector messages** : In fact (b) can be done efficiently in more general settings where we allow the output of each node in the network to be a vector (or even matrix/tensor) instead of only a real number. Here we need to replace  $\frac{\partial u}{\partial z_j} \cdot S$  by  $\frac{\partial u}{\partial z_j}[S]$ , which denotes the result of applying the operator  $\frac{\partial u}{\partial z_j}$  on  $S$ . We note that to be consistent with the convention in the usual exposition of backpropagation, when  $y \in \mathbb{R}^p$  is a function of  $x \in \mathbb{R}^q$ , we use  $\frac{\partial y}{\partial x}$  to denote  $q \times p$  dimensional matrix with  $\partial y_j / \partial x_i$  as the  $(i, j)$ -th entry. Readers might notice that this is the transpose of the usual Jacobian matrix defined in mathematics. Thus  $\frac{\partial y}{\partial x}$  is an operator that maps  $\mathbb{R}^p$  to  $\mathbb{R}^q$  and we can verify  $S$  has the same dimension as  $u$  and  $\frac{\partial u}{\partial z_j}[S]$  has the same dimension as  $z_j$ .

For example, as illustrated below, suppose the node  $U \in \mathbb{R}^{d_1 \times d_3}$  is a product of two matrices  $W \in \mathbb{R}^{d_2 \times d_3}$  and  $Z \in \mathbb{R}^{d_1 \times d_2}$ . Then we have that  $\partial U / \partial Z$  is a linear operator that maps  $\mathbb{R}^{d_2 \times d_3}$  to  $\mathbb{R}^{d_1 \times d_3}$ , which naively requires a matrix representation of dimension  $d_2 d_3 \times d_1 d_3$ . However, the computation (b) can be done efficiently because

$$\frac{\partial U}{\partial Z}[S] = W^T S.$$

Such vector operations can also be implemented efficiently using today's GPUs.

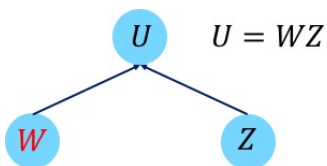


Figure 4.3: Vector version of above

#### 4.4 Notable Extensions

*Allowing weight tying:* In many neural architectures, the designer wants to force many network units such as edges or nodes to share the same parameter. For example, in including the ubiquitous convolutional net, the same filter has to be applied all over the image, which implies reusing the same parameter for a large set of edges between two layers of the net.

For simplicity, suppose two parameters  $a$  and  $b$  are supposed to share the same value. This is equivalent to adding a new node  $u$  and connecting  $u$  to both  $a$  and  $b$  with the operation  $a = u$  and  $b = u$ . Thus, by chain rule,

$$\frac{\partial f}{\partial u} = \frac{\partial f}{\partial a} \cdot \frac{\partial a}{\partial u} + \frac{\partial f}{\partial b} \cdot \frac{\partial b}{\partial u} = \frac{\partial f}{\partial a} + \frac{\partial f}{\partial b}.$$

Hence, equivalently, the gradient with respect to a shared parameter is the sum of the gradients with respect to individual occurrences.

*Backpropagation on networks with loops.* The above exposition assumed the network is acyclic. Many cutting-edge applications such as machine translation and language understanding use networks with directed loops (e.g., recurrent neural networks). These architectures—all examples of the "differentiable computing" paradigm below—can get complicated and may involve operations on a separate memory as well as mechanisms to shift attention to different parts of data and memory.

Networks with loops are trained using gradient descent as well, using *back-propagation through time* which consists of expanding the network through a finite number of time steps into an acyclic graph, with replicated copies of the same network. These replicas share the weights (weight tying!) so the gradient can be computed. In practice an issue may arise with *exploding or vanishing gradients*, which impact convergence. Such issues can be carefully addressed in practice by clipping the gradient or re-parameterization techniques such as *long short-term memory*. Recent work suggests that careful initialization of parameters can ameliorate some of the vanishing gradient problems.

The fact that the gradient can be computed efficiently for such general networks with loops has motivated neural net models with memory or even data structures (see for example *neural Turing machines* and *differentiable neural computer*). Using gradient descent, one can optimize over a family of parameterized networks with loops to find the best one that solves a certain computational task (on the training examples). The limits of these ideas are still being explored.

#### 4.4.1 Hessian-vector product in linear time: Pearlmutter's trick

It is possible to generalize backpropagation to work with 2nd order derivatives, specifically with the Hessian  $H$  which is the symmetric matrix whose  $(i, j)$  entry is  $\partial^2 f / \partial w_i \partial w_j$ . Sometimes  $H$  is also denoted  $\nabla^2 f$ . Just writing down this matrix takes quadratic time and memory, which is infeasible for today's deep nets. Surprisingly, using backpropagation it is possible to compute in linear time the matrix-vector product  $Hx$  for any vector  $x$ .

**Claim 4.4.1.** *Suppose an acyclic network with  $V$  nodes and  $E$  edges has output  $f$  and leaves  $z_1, \dots, z_m$ . Then there exists a network of size  $O(V + E)$  that has  $z_1, \dots, z_m$  as input nodes and  $\frac{\partial f}{\partial z_1}, \dots, \frac{\partial f}{\partial z_m}$  as output nodes.*

The proof of the Claim follows in straightforward fashion from implementing the message passing protocol as an acyclic circuit.

Next we show how to compute  $\nabla^2 f(z) \cdot v$  where  $v$  is a given fixed vector. Let  $g(z) = \langle \nabla f(z), v \rangle$  be a function from  $\mathbb{R}^d \rightarrow \mathbb{R}$ . Then by the Claim above,  $g(z)$  can be computed by a network of size  $O(V + E)$ . Now apply the Claim again on  $g(z)$ , we obtain that  $\nabla g(z)$  can also be computed by a network of size  $O(V + E)$ .

Note that by construction,

$$\nabla g(z) = \nabla^2 f(z) \cdot v.$$

Hence we have computed the Hessian vector product in network size time.

# 5

## *Advanced Optimization notions*

This chapter covers the basic 2nd order method (Newton's method) and then briefly discusses momentum, AdaGrad (as well as AdaDelta/RMSProp) and Adam. Some discussion of attempts to leverage Hessian-vector products and why they don't appear to help.

Using Hessian-vector products for DL dates back to J. Martens "Deep learning via Hessian-free optimization."





## 6

# *Basics of generalization theory*

Recall from Chapter 1 the language of Empirical Risk Minimization from Chapter 1. A datapoint  $x$  (for classification this is actually a pairing of a vector and a label) come from a distribution  $\mathcal{D}$  and  $S$  denotes the training sample. The loss of a hypothesis  $h$  on datapoint  $x$  is  $\ell(x, h)$ . (Since hypothesis in deep learning is given by a parameter vector  $w$  we may also represent this as  $\ell(x, w)$ .) In generalization theory we are interested in understanding the relationship between the test loss and the training loss (respectively):

$$L_{\mathcal{D}}(h) = \mathbb{E}_{x \in \mathcal{D}} [\ell(x, h)] \quad \text{and} \quad \hat{L}_S(h) = \mathbb{E}_{x \in S} [\ell(x, h)]. \quad (6.1)$$

(Here  $\hat{\cdot}$  refers to “empirical.” The training is considered a success if  $L_S(h)$  is small and the *generalization error*  $\Delta_S(h) = L_{\mathcal{D}}(h) - \hat{L}_S(h)$  is small too.

Generalization theory gives estimates of the number of training samples sufficient to guarantee low generalization error. The classic ideas described in this chapter give very loose (i.e., trivial) estimates for deep learning. A later chapter in this book describes recent attempts to come up with tighter estimates of sample complexity.

Generalization theory takes inspiration from an old philosophical principle called *Occam’s razor*: given a choice between a simpler theory of science and a more convoluted theory, both of which explain some empirical observations, we should trust the simpler one. For instance, Copernicus’s heliocentric theory of the solar system gained favor in science because it explained known facts more simply than the ancient Aristotelian theory. While this makes intuitive sense, Occam’s razor is a bit vague and hand-wavy. What makes a theory “simpler” or “better”?

### 6.1 *Occam’s razor formalized for ML*

The following is the mapping from the above intuitive notions to notions in ML. (For simplicity we focus only on supervised learning

here, and consider other settings in later chapters.)

<i>Observations/evidence</i>	$\leftrightarrow$	Training dataset $S$
<i>theory</i>	$\leftrightarrow$	hypothesis $h$
<i>All possible theories</i>	$\leftrightarrow$	hypothesis class $\mathcal{H}$
<i>Finding theory to fit observations</i>	$\leftrightarrow$	Minimize training loss to find $h \in \mathcal{H}$
<i>Theory is good (good predictions in new settings)</i>	$\leftrightarrow$	$h$ has low test loss
<i>Simpler theory</i>	$\leftrightarrow$	$h$ has shorter description

The notion of “shorter description” will be formalized in a variety of ways using a *complexity measure* for the class  $\mathcal{H}$ , denoted  $\mathcal{C}(\mathcal{H})$ , and use it to upper bound the generalization error.

Let  $S$  be a sample of  $m$  datapoints. Empirical risk minimization gives us  $\hat{h} = \arg \min \widehat{L}_S(h)$ . Of course, in deep learning we may not find the absolute optimum  $h$  but in practice the training loss becomes very small and near-zero. Intuitively, if generalization error is large then the hypothesis’s performance on training sample  $S$  does not accurately reflect the performance on the full distribution of examples, and we say it *overfitted* to the sample  $S$ .

The typical upper bound on generalization error <sup>1</sup> shows that with probability at least  $1 - \delta$  over the choice of training data, the following

$$\Delta_S(h) \leq \sqrt{\frac{\mathcal{C}(\mathcal{H}) + O(\log(1/\delta))}{m}} + \text{Sampling error term.} \quad (6.2)$$

Thus to drive the generalization error down it suffices to make  $m$  significantly larger than the “Complexity Measure”. Hence classes with lower complexity require fewer training samples, in line with Occam’s intuition.

### 6.1.1 Motivation for generalization theory

If the experimenter has already decided on the architecture, training algorithm etc. then generalization theory is of very limited use. They can use a held out dataset to estimate the generalization error.

The hope in developing generalization theory is that it provides insights into how to design architectures and algorithms that lead to good generalization. One uses generalization bounds to estimate the generalization error using the trained model  $h$  and the training dataset  $S$ . <sup>2</sup>

### 6.1.2 Motivating example: Polynomial interpolation

Suppose we are given  $n$  points  $(x^{(1)}, y^{(1)}), \dots, (x^{(n)}, y^{(n)})$  chosen according to the following probabilistic process:  $x^{(i)}$  is chosen from

<sup>1</sup> This is the format of typical generalization bound!

<sup>2</sup> Note: there are ways to cheat in how we’ve set up the problem of predicting generalization error using only the training data. For example, the bound could involve reserving a small heldout set from  $S$  and using it to compute an estimate of generalization error from it. This requires no particular insight into the trained classifier.

some distribution on  $[0, 1]$  and  $y^{(i)} = p(x^{(i)}) + \eta$  where  $p(\cdot)$  is an unknown degree  $d$  polynomial and  $\eta$  is a sample from the noise distribution  $\mathcal{N}(0, \sigma^2)$ . Since  $\mathbb{E}[\eta] = 0$  and  $\mathbb{E}[\eta^2] = \sigma^2$  the obvious way to find  $p$  is to minimize the least square fit to find the polynomial's coefficients  $\theta_0, \theta_1, \dots, \theta_d$ :

$$\ell(\vec{\theta}) = \sum_{i=1}^n (y^{(i)} - \sum_{j=0}^d \theta_j x^{(i)j})^2.$$

This is implicitly doing linear regression using a new data representation, whereby the point  $x \in \mathfrak{X}$  is represented using the vector  $(1, x, x^2, \dots, x^d)$ .

But what if we don't know  $d$  and try to fit a degree  $N$  polynomial where  $N \gg d$ ? Under what conditions would minimizing the above loss give us roughly the same polynomial as  $p(\cdot)$ ? A practical idea—noting the fact that the above loss is  $n\sigma^2$  even for the ground truth polynomial  $p(\cdot)$ —is to add for some large-ish  $\lambda > 0$  the regularizer  $\lambda \|\theta\|_2^2$  to the above loss, which signals to gradient descent that it is not important to reduce the least squares loss all the way to zero, and it is important to find solutions  $\theta$ 's of low norm. In this case  $\|\theta\|_2^2$  is the implicit complexity measure.

This example is intuitive and can be analysed more rigorously but requires the theory of orthonormal polynomials with respect to natural distributions on  $[0, 1]$ . See *citations??*

## 6.2 Some simple upper bounds on generalization error

Most generalization bounds involve a *union bound* over a set of possible hypotheses, namely, any events  $A_1, A_2, \dots$  satisfy  $\Pr[\cup_i A_i] \leq \sum_i \Pr[A_i]$ .

The first upper bound is an almost trivial example for this, but as we shall see is also at the heart of most other generalization bounds (albeit often hidden inside the proof). The bound shows that if a hypothesis class contains at most  $N$  distinct hypotheses, then  $\log N$  (i.e., the number of bits needed to describe a single hypothesis in this class) functions as a complexity measure.<sup>3</sup>

**Theorem 6.2.1** (Simple union bound). *If the loss function takes values in  $[0, 1]$  and hypothesis class  $\mathcal{H}$  contains  $N$  distinct hypotheses then with probability at least  $1 - \delta$ , every  $h \in \mathcal{H}$  satisfies*

$$\Delta_S(h) \leq 2\sqrt{(\log N + \log(1/\delta))/m}.$$

*Proof.* For any fixed hypothesis  $g$  imagine drawing a training sample of size  $m$ . Then  $\widehat{L}_S(g)$  is an average of i.i.d. variables and its expectation is  $L_D(g)$ . Concentration bounds imply that  $L_D(g) - \widehat{L}_S(g)$  has

<sup>3</sup> The union bound is also referred to as **uniform convergence framework** in many books. Often the hypothesis class is infinite but the proof discretizes it, as in Theorem 6.2.7.

a concentration property at least as strong as univariate Gaussian  $\mathcal{N}(0, 1/m)$ . The previous statement is true for all hypotheses  $g$  in the class, so the union bound implies that the probability is at most  $N \exp(-\epsilon^2 m/4)$  that this quantity exceeds  $\epsilon$  for *some* hypothesis in the class. Since  $h$  is the solution to ERM, we conclude that when  $\delta \leq N \exp(-\epsilon^2 m/4)$  then  $\Delta_S(h) \leq \epsilon$ . Simplifying and eliminating  $\epsilon$ , we obtain the theorem.  $\square$

At first sight the union bound appears useless for deep nets *per se* because the set of hypotheses—even after we have fixed the architecture—consists of all vectors in  $\mathbb{R}^k$ , where  $k$  is the number of real-valued parameters. This is an uncountable set!

**Example 6.2.2** (Possible way around?). *As we saw in Chapter 2, the end result of gradient descent on the loss is special: for instance it must be a stationary point (i.e.,  $\nabla = 0$ ) of the training loss. One can similarly imagine other conditions on Hessian  $\nabla^2(\cdot)$  and so forth. One could hope that the set of points in the loss landscape with such properties could be small and finite. This line of investigation hasn't yet worked out because current nets are so overparametrized (i.e., number of parameters far exceeding the number of training data points) that the set of such solution points in the landscape is also in general a continuous set (i.e., infinite). The next hope is to take into account the training algorithm, because not all solution points may be accessible via standard training algorithms. These are some ideas for restricting attention to a finite set of solutions, though they haven't yet worked out.*

There is a more obvious way to turn the set of solutions into a finite set: *discretization!* Suppose we assume that the  $\ell_2$  norm of the parameter vectors is at most 1, meaning the set of all deep nets has been identified with  $\text{Ball}(0, 1)$ . (Here  $\text{Ball}(w, r)$  refers to set of all points in  $\mathbb{R}^k$  within distance  $r$  of  $w$ .)

Suppose the loss is Lipschitz in the parameters: for every datapoint  $x$  and parameter vectors  $w_1, w_2$   $\|\ell(x, w_1) - \ell(x, w_2)\| \leq C\|w_1 - w_2\|_2$  for some explicit constant  $C$ . The arguments below only need local Lipschitz-ness, namely for the condition to hold for  $\|w_1 - w_2\|_2 \leq \rho$  for some explicit constant  $\rho$ . Furthermore it only requires Lipschitz-ness of the average loss on the training set, not loss on single data points.

Suppose  $\rho > 0$  is such that if  $w_1, w_2 \in \mathbb{R}^k$  satisfy  $\|w_1 - w_2\|_2 \leq \rho$  then the nets with these two parameter vectors have essentially the same loss on every input, meaning the losses differ by at most  $\gamma$  for some  $\gamma > 0$ . (NB: This means local Lipschitz constant at most  $\gamma/\rho$ .) It makes intuitive sense such a  $\rho$  must exist for every  $\gamma > 0$  since as we let  $\rho \rightarrow 0$  the two nets become equal.

**Problem 6.2.3.** Compute Lipschitz constant of the  $\ell_2$  regression loss: the loss on datapoint  $(x, y)$  is  $(w \cdot x - y)^2$ .

**Problem 6.2.4.** Compute Lipschitz constant of  $\ell_2$  loss for a two layer deep net with ReLU gates (zero bias) on the middle layer. Assume the two parameter vectors are infinitesimally close.

**Definition 6.2.5** ( $\rho$ -cover). A set of points  $w_1, w_2, \dots \in \mathbb{R}^k$  is a  $\rho$ -cover in  $\mathbb{R}^k$  if for every  $w \in \text{Ball}(0, 1)$  there is some  $w_i$  such that  $w \in \text{Ball}(w_i, \rho)$ .

**Lemma 6.2.6** (Existence of  $\rho$ -cover). There exists a  $\rho$ -cover of size at most  $(2/\rho)^k$ .

*Proof.* The proof is simple but clever. Let us pick  $w_1$  arbitrarily in  $\text{Ball}(0, 1)$ . For  $i = 2, 3, \dots$  do the following: arbitrarily pick any point in  $\text{Ball}(0, 1)$  outside  $\cup_{j < i} \text{Ball}(w_j, \rho)$  and designate it as  $w_{i+1}$ .

*A priori* it is unclear if this process will ever terminate. We now show it does after at most  $(2/\rho)^k$  steps. To see this, it suffices to note that  $\text{Ball}(w_i, \rho/2) \cap \text{Ball}(w_j, \rho/2) = \emptyset$  for all  $i < j$ . (Because if not, then  $w_j \in \text{Ball}(w_i, \rho)$ , which means that  $w_j$  could not have been picked during the above process.) Thus we conclude that the process must have stopped after at most

$$\text{volume}(\text{Ball}(0, 1)) / \text{volume}(\text{Ball}(0, \rho/2))$$

iterations, which is at most  $(2/\rho)^k$  since ball volume in  $\mathbb{R}^k$  scales as the  $k$ th power of the radius.

Finally, the sequence of  $w_i$ 's at the end must be a  $\rho$ -cover because the process stops only when no point can be found outside  $\cup_j \text{Ball}(w_j, \rho)$ .  $\square$

**Theorem 6.2.7** (Generalization bound for normed spaces). If (i) hypotheses are unit vectors in  $\mathbb{R}^k$  and (ii) every two hypotheses  $h_1, h_2$  with  $\|h_1 - h_2\|_2 \leq \rho$  differ in terms of loss on every datapoint by at most  $\gamma$  then <sup>4</sup>

$$\Delta_S(h) \leq \gamma + 2\sqrt{k \log(2/\rho) / m}.$$

*Proof.* Apply the union bound on the  $\rho$ -cover. Every other net can have loss at most  $\gamma$  higher than nets in the  $\rho$ -cover.  $\square$

### 6.3 Data dependent complexity measures

Thus far we considered complexity measures for hypothesis classes as a way to quantify their “complicatedness.” : the size of the hypothesis class (assuming it is finite) and the size of a  $\gamma$ -cover in it. Of course, the resulting bounds on sample complexity were still loose.

But these simple bounds hold for every data distribution  $\mathcal{D}$ . In practice, it seems clear that deep nets—or any learning method—works by being able to exploit properties of the input distribution

<sup>4</sup> Even ignoring the dependence on the Lipschitz constant, this bound requires the number of datapoints to grow linearly with the number of trainable parameters in the net. This does not begin to explain the surprising effectiveness of real-life deep learning, where the number of parameters greatly exceeds the number of training datapoints.

(e.g., convolutional structure exploits the fact that all subpatches of images can be processed very similarly). Thus one should try to prove some measure of complicatedness that depends on the data distribution.

### 6.3.1 Rademacher Complexity

Rademacher complexity is a complexity measure that depends on data distribution. For simplicity we will assume loss function takes values in  $[0, 1]$ .

The definition concerns the following thought experiment. Recall that the distribution  $\mathcal{D}$  is on labeled datapoints  $(x, y)$ . For simplicity we denote the labeled datapoint as  $z$ .

Now *Rademacher Complexity*<sup>5</sup> of hypothesis class  $\mathcal{H}$  on a distribution  $\mathcal{D}$  is defined as follows where  $l(z, h)$  is loss of hypothesis  $h$  on labeled datapoint  $z$ .

$$\mathcal{R}_{m,D}(\mathcal{H}) = \mathbb{E}_{S_1, S_2} \left[ \frac{1}{2m} \sup_{h \in \mathcal{H}} \left| \sum_{z \in S_1} l(z, h) - \sum_{z \in S_2} l(z, h) \right| \right], \quad (6.3)$$

where the expectation is over  $S_1, S_2$  are two iid samples (i.e., multisets) of size  $m$  each from the data distribution  $\mathcal{D}$ . The following theorem relates this to generalization error of the trained hypothesis.

**Theorem 6.3.1.** *If  $h$  is the hypothesis trained via ERM using a training set  $S_2$  of size  $m$ , then the probability (over  $S_2$ ) is  $> 1 - \delta$ , that*

$$\Delta_{S_2}(h) \leq 2\mathcal{R}_{m,D}(\mathcal{H}) + O((\log(1/\delta))/\sqrt{m}).$$

*Proof.* The generalization error  $\Delta_{S_2}(h) = L_{\mathcal{D}}(h) - \widehat{L}_{S_2}(h)$ , and ERM guarantees an  $h$  that maximizes this. Imagine we pick another  $m$  iid samples from  $\mathcal{D}$  to get another (multi)set  $S_1$  then with probability at least  $1 - \delta$  the loss on these closely approximates  $L_{\mathcal{D}}(h)$ :

$$\Delta_{S_2}(h) \leq \widehat{L}_{S_1}(h) - \widehat{L}_{S_2}(h) + O((\log(1/\delta))/\sqrt{m}).$$

Now we notice that  $S_1, S_2$  thus drawn are exactly like the sets drawn in the thought experiment<sup>6</sup> (6.3) and the maximizer  $h$  for this expression defined  $\mathcal{R}_{m,D}$ . So the right hand side is at most

$$2\mathcal{R}_{m,D}(\mathcal{H}) + O((\log(1/\delta))/\sqrt{m}).$$

□

**Example 1:** We can show that the Rademacher complexity of the set of linear classifiers (unit norm vectors  $U = \{w | w \in \mathbb{R}^d, \|w\|_2 = 1\}$ ), on a given sample  $S = \{x_1, x_2, \dots, x_m\}$  (each  $x_i \in \mathbb{R}^d$ ) is  $\leq \max_{i \in [m]} \|x_i\|_2 / \sqrt{m}$ .

<sup>5</sup> Standard accounts of this often confuse students, or falsely impress them with a complicated proof of Thm 6.3.1. In the standard definition, loss terms are weighted by iid  $\pm 1$  random variables. Its value is within  $\pm O(1/\sqrt{m})$  of the one in our definition.

<sup>6</sup> Here hypothesis  $h$  is allowed to depend on  $S_2$  but not  $S_1$ . In the thought experiment the supremum is over  $h$  that can depend on both. This discrepancy only helps the inequality, since the latter  $h$  can achieve a larger value. Note that the factor 2 is because of scaling of  $2m$  in (6.3).

**Example 2:** Consider the kernel classifier of the form  $h(x) = z^\top G^{-1}y$  we studied in Section 3.2 where  $G$  is the  $n \times n$  kernel matrix,  $y$  is the labels and  $z$  is the column vector whose  $i$ -th coordinate is  $K(x, x_i)$ . We can show the Rademacher complexity upper is  $\sqrt{2y^\top G y \cdot \text{Tr}(G)/n}$ . We will use this result in Chapter 10 to prove certain over-parameterized student nets can learn simple two-layer teacher nets.

### 6.3.2 Alternative Interpretation: Ability to correlate with random labels

Sometimes teachers explain Rademacher complexity more intuitively as *ability of classifiers in  $\mathcal{H}$  to correlate with random labelings of the data*. This is best understood for binary classification (i.e., labels are 0/1), and the loss function is also binary (loss 0 for correct label and 1 incorrect label). Now consider the following experiment: Pick  $S_1, S_2$  as in the definition of Rademacher Complexity, and imagine flipping the labels of  $S_1$ . Now average loss on  $S_2$  is  $1 - \widehat{L}_{S_2}(h)$ . Thus selecting  $h$  to maximise the right hand side of (6.3) is like finding an  $h$  that has low loss on  $S_1 \cup S_2$  where the labels have been flipped on  $S_1$ . In other words,  $h$  is able to achieve low loss on datasets where labels were flipped for some randomly chosen set of half of the training points.

When the loss is not binary a similar statement still holds qualitatively.

## 6.4 Understanding limitations of the union-bound approach

The phenomenon captured in the union bound approach and related approaches is also referred to as *uniform convergence*. If we have identified a finite set  $\mathcal{H}$  of hypotheses and sample  $S$  of datapoints is large enough then with probability is at least  $1 - \delta$  over choice of  $S$  that

$$\|L_{\mathcal{D}}(h) - \widehat{L}_S(h)\| \leq \epsilon \quad \forall h \in \mathcal{H}. \quad (6.4)$$

Here the important point to note is that a fixed sample set  $S$  can be used for good estimate of generalization error for *every* classifier  $h$  in the class. Of course, using  $\gamma$ -cover this kind of conclusion can be shown also for classes  $\mathcal{H}$  that are a continuous set, e.g. hypotheses with bounded  $\ell_2$  norm.

Now we describe a simple cautionary example from <sup>7</sup> pinpointing why this framework may be tricky to apply in modern settings, especially deep learning. It also illustrates how classifiers obtained at the end of training using gradient descent may not be easy to describe using norms—at least the obvious descriptions have too large a norm. The intuitive difficulty is that during training the classifier can pick up on a lot of information in the data that is “irrelevant” to

<sup>7</sup> V Nagarajan and Zico Kolter. Uniform convergence may be unable to explain generalization in deep learning. *NeurIPS*, 2019

the overall classification, and this may even account for most of the  $\ell_2$  norm. Thus a norm-based description gives a huge overestimate of the hypothesis class, much larger than the true set of hypotheses output-able by the training algorithm.

Suppose the points are in  $\mathfrak{R}^{D+K}$  and the labels are  $\pm 1$ . There is a fixed vector  $u \in \mathfrak{R}^k$  such that labeled datapoints  $(x, y)$  come from the following distribution: first label  $y$  is uniformly picked in  $\{\pm 1\}$  and then the first  $K$  coordinates of  $x$ —which we denote  $x_1$  for convenience—are set to the vector  $y \cdot u$ . The remaining  $D$  coordinates, denoted  $x_2$  consist of a random vector  $\mathfrak{R}^D$ , whose each coordinate is drawn independently from  $\mathcal{N}(0, 1/D)$ . Measure concentration implies  $x_2$  is distributed essentially like a random unit vector in  $\mathfrak{R}^D$ .

The classification can clearly be solved using the linear classifier  $x \rightarrow \text{sgn}(w^* \cdot x)$  where the first  $K$  coordinates contain  $w_1^* = u/\|u\|_2^2$ , and the last  $D$  coordinates contain  $w_2^* = 0$ .

Now we describe how to set the various parameters. As usual  $m$  denotes training set size. We set

$$m^{3/2} \sqrt{(\log 1/\delta)} \approx D \quad (6.5)$$

$$\|u\|_2^2 = \frac{4}{m} \quad (6.6)$$

Let's consider a simple training objective: find linear classifier  $h(x)$  that maximises  $y \cdot h(x)$ . Roughly speaking, this ignores the magnitude of  $h(x)$  and tries to align the sign of  $h(x)$  and  $y$ . Using learning rate  $\eta = 1$  and a sample  $S$  of  $m$  datapoints  $(x^i, y^i)$  for  $i = 1, \dots, m$  gradient descent produces a classifier with  $w = (w_1, w_2)$  where

$$w_1 = m \cdot u, \quad w_2 = \sum_i y^i x_2^i. \quad (6.7)$$

Notice that  $w_2$  is a sum of  $m$  random unit vectors, and its norm is fairly tightly concentrated around  $\sqrt{m}$ . In other words, unlike our ideal classifier  $w^*$  the learnt classifier has a lot of junk in the last  $D$  coordinates. The  $\ell_2$  norm of the learnt classifier is around  $\sqrt{m^2 \|u\|_2^2 + m}$ , and thus (6.6) implies this norm is  $\sqrt{5m}$ .

Let's check that the junk coordinates do not interfere with classification for most data points. Given a new data point  $x = (y \cdot u, x_2)$  where  $x_2$  is a random unit vector, the learnt classifier produces the answer  $w \cdot x = y \|u\|_2^2 + x_2 \cdot (\sum_i y^i x_2^i)$ . Since inner product between a fixed vector and a random gaussian vector  $\mathcal{N}(0, 1/D)$  is a univariate gaussian with standard deviation  $1/\sqrt{D}$  times the norm of the fixed vector, we see that the answer is correct with probability  $1 - \delta$  so long as

$$\|u\|_2^2 > \frac{2}{D} \sqrt{m(\log 1/\delta)},$$

which is the case. Thus the learnt classifier works fine on random test data points.



Now let's see a very natural set of  $m$  points where this same classifier gives the wrong answer. It consists of the same data points as above, except the  $x_2$  part gets a sign flip. That is,  $(y^i u, -x_2^i, y^i)$ . We denote this by  $S^{\text{flip}}$ . Note that viewed in isolation,  $S^{\text{flip}}$  also contains  $m$  iid samples from the data distribution just as  $S$  does (since in each data point the last  $D$  coordinates were picked from the distribution of random unit vectors, which is symmetric with respect to sign flip). Thus the following is surprising, but is easily verified.

**Problem 6.4.1.** *Show that for most  $S$ , the learnt classifier misclassifies every point in  $S^{\text{flip}}$ .*

Finally, we explain why this points to limitations of the uniform convergence framework (which we called *union bound* framework). Suppose we are trying to explain why gradient-based learning described above works. We let  $\mathcal{H}$  be the class of classifiers learnt using sample sets of size  $m$ . As argued above, these learnt classifiers do generalize well. However, what we have actually shown in Problem 6.4.1 is the following:

**Corollary 6.4.2.** *For almost all  $S$ , there is a classifier  $h \in \mathcal{H}$  that violates the property in (6.4).*

In other words, the union bound argument does not apply as is to show that the learnt classifier generalizes.

Note that the limitations shown above do not hold if we are allowed to modify/prune the classifier obtained at the end of training. One can imagine identifying non-influential coordinates in the learnt classifier via some simple test and realizing that the last  $D$  coordinates can be zero-ed out without greatly affecting accuracy. Then all learnt classifiers become scalar multiples of the ideal classifier  $w^*$ . In other words, the limitations shown here do not apply to the approach we describe in the next Section.

## 6.5 A Compression-based framework

Now we described a simple compression-based technique from <sup>8</sup> that formalizes a very simple idea. <sup>9</sup> Suppose the training dataset  $S$  contains  $m$  samples, and  $h$  is a classifier from a complicated class (e.g., deep nets with much more than  $m$  parameters) that incurs very low empirical loss. We are trying to understand from looking at  $h$  and  $S$  how well  $h$  will generalize. Now suppose we can compute a classifier  $g$  with discrete trainable parameters much less than  $m$  and which incurs similar loss on the training data as  $h$ . We call this an *approximator* for  $h$ . Then if  $g$  has sufficiently low description length, its generalization follows by simple union bound argument. <sup>10</sup>

<sup>8</sup> Sanjeev Arora, Rong Ge, Behnam Neyshabur, and Yi Zhang. Stronger generalization bounds for deep nets via a compression approach. In *Proc. ICML 2018*, pages 254–263, 2018

<sup>9</sup> Do not confuse this from another (unrelated) technique in generalization theory based upon *data compression*.

<sup>10</sup> This scenario is quite reminiscent of empirical work in network pruning, whereby trained deep nets are compressed using one of a long list of methods that prune away lots of parameters and retrain the rest. If network left after pruning is compact enough, one can conceivably prove generalization bounds for the pruned net. See

This framework has the advantage of staying with intuitive parameter counting and to avoid explicitly dealing with the hypothesis class that includes  $h$  (see note after Theorem 6.5.3). Notice, the mapping from  $f$  to  $g$  merely needs to *exist*, not to be efficiently computable. But in all our examples the mapping will be explicit and fairly efficient. Now we formalize the notions. The proofs are elementary via concentration bounds and appear in the appendix.

**Definition 6.5.1** ( $(\gamma, S)$ -compressible). *Let  $f$  be a classifier and  $G_{\mathcal{A}} = \{g_A | A \in \mathcal{A}\}$  be a class of classifiers. We say  $f$  is  $(\gamma, S)$ -compressible via  $G_{\mathcal{A}}$  if there exists  $A \in \mathcal{A}$  such that for any  $x \in S$ , we have for all  $y$*

$$|f(x)[y] - g_A(x)[y]| \leq \gamma.$$

We also consider a different setting where the compression algorithm is allowed a “helper string”  $s$ , which is arbitrary but fixed before looking at the training samples. Often  $s$  will contain random numbers.<sup>11</sup>

**Definition 6.5.2** ( $(\gamma, S)$ -compressible using helper string  $s$ ). *Suppose  $G_{\mathcal{A},s} = \{g_{A,s} | A \in \mathcal{A}\}$  is a class of classifiers indexed by trainable parameters  $A$  and fixed strings  $s$ . A classifier  $f$  is  $(\gamma, S)$ -compressible with respect to  $G_{\mathcal{A},s}$  using helper string  $s$  if there exists  $A \in \mathcal{A}$  such that for any  $x \in S$ , we have for all  $y$*

$$|f(x)[y] - g_{A,s}(x)[y]| \leq \gamma.$$

The following theorem is a simple application of the union bound method above.

**Theorem 6.5.3.** *Suppose  $G_{\mathcal{A},s} = \{g_{A,s} | A \in \mathcal{A}\}$  where  $A$  is a set of  $q$  parameters each of which can have at most  $r$  discrete values and  $s$  is a helper string. Let  $S$  be a training set with  $m$  samples. If the trained classifier  $f$  is  $(\gamma, S)$ -compressible via  $G_{\mathcal{A},s}$  with helper string  $s$ , then there exists  $A \in \mathcal{A}$  with high probability over the training set,*

$$L(g_A) \leq \widehat{L}_\gamma(f) + O\left(\sqrt{\frac{q \log r}{m}}\right),$$

where  $L(f) = \mathbb{E}_{(x,y) \in \mathcal{D}}[f(x)[y] \leq \max_{j \neq y} f(x)[j]]$  is the expected error and  $\widehat{L}_\gamma(f)$  is the proportion of data  $(x, y)$  satisfying  $f(x)[y] \leq \max_{j \neq y} f(x)[j]$  in the training set  $S$ .

*Remarks:* (1) The framework proves the generalization not of  $f$  but of its compression  $g_A$ . (An exception is if the two are shown to have similar loss at every point in the domain, not just the training set. This is the case in Theorem 6.5.6.)

(2) The previous item highlights the difference from what we called the union bound earlier (Theorem 6.2.1). There, one needs to

<sup>11</sup> A simple example is to let  $s$  be the random initialization used for training the deep net. Then one could compress the *difference* between the final weights and  $s$ ; this can give better generalization bounds.

fix a hypothesis class *independent* of the training set. By contrast we have no hypothesis class, only a *single* neural net that has some specific properties on a *single* finite training set. But if we can compress this specific neural net to a simpler neural nets with fewer parameters then we can use covering number argument on this simpler class to get the generalization of the compressed net.

(3) Issue (1) exists also in how researchers often apply the standard PAC-Bayes framework for deep nets (Section 6.6).

### 6.5.1 Example 1: Linear classifiers with margin

To illustrate the above compression method we use linear classifiers with high margins. Consider a simple family of linear classifiers, consisting of unit vectors  $c \in \mathbb{R}^d$  whose  $\pm 1$  output on input  $x$  is given by  $\text{sgn}(c \cdot x)$  (i.e., sign of the inner product with the datapoint). Assume that all data points are also unit vectors. Say  $c$  has *margin*  $\gamma$  if for all training pairs  $(x, y)$  we have  $y(c^\top x) \geq \gamma$ .

We show how to compress such a classifier with margin  $\gamma$  to one that has only  $O(1/\gamma^2)$  non-zero entries. First, assume all  $c_i$  have absolute value less than  $\gamma^2/8$ .

For each coordinate  $i$ , toss a coin with  $\Pr[\text{heads}] = p_i = 8c_i^2/\gamma^2$  and if it comes up heads set the coordinate to equal to  $c_i/p_i = \gamma^2/8c_i$ . This yields a vector  $\hat{c}$ . The expected number of non-zero entries in  $\hat{c}$  is  $\sum_{i=1}^d p_i = 8/\gamma^2$ . By Chernoff bound we know with high probability the number of non-zero entries is at most  $O(1/\gamma^2)$ .

Furthermore, variance of coordinate  $i$  of  $\hat{c}$  is  $2p_i(1-p_i)\frac{c_i^2}{p_i^2} \leq \frac{2c_i^2}{p_i} \leq \gamma^2/4$ . Therefore, for any unit vector  $u$  that is independent with the choice of  $\hat{c}$ , we have  $\mathbb{E}[\hat{c}^\top u] = c^\top u$ . Now we estimate variance of the random variable  $\hat{c}^\top u$ . It is  $\leq \gamma^2/4 \cdot \|u\|^2 \leq \gamma^2/4$ . By Chebyshev's inequality we know  $\Pr[|\hat{c}^\top u - c^\top u| \geq \gamma] \leq 1/4$ , so  $\hat{c}$  and  $c$  will make the same prediction for all  $u$  satisfying  $|c^\top u| \geq \gamma$ . We can then apply Theorem 6.5.3 on a discretized version of  $\hat{c}$  (via trivial rounding) to show that the sparsified classifier has good generalization with  $O(\log d/\gamma^2)$  samples.

**Problem 6.5.4.** Redo the proof above when some coordinates have absolute value more than  $\gamma^2/8$ .

This compressed classifier works correctly for a fixed input  $x$  with constant probability but not high probability. To fix this, one can recourse to the “compression with fixed string” model. The fixed string is a random linear transformation. When applied to unit vector  $c$ , it tends to equalize all coordinates and the guarantee  $|\hat{c}^\top u - c^\top u| < \gamma$  can hold with high probability. This random linear transformation can be fixed before seeing the training data.

**Problem 6.5.5.** *Prove the above property of random linear transformations. That is, let  $M$  be a random matrix of size  $O(1/\gamma^2) \times d$ , drawn from a suitable distribution you choose before seeing the unit vector  $c$  and the training data. Then, show that the following holds for fixed unit vectors  $c$  and  $u$  with high probability*

$$\|Mc\|_\infty = O(1), \quad |\langle Mc, Mu \rangle - \langle c, u \rangle| < \gamma.$$

This means we can compress a unit vector  $c$  to  $\hat{c} = M^\top Mc$ . Finally, Apply Theorem 6.5.3 on a discretized version of  $\hat{c}$  to show a good generalization bound with  $\tilde{O}(1/\gamma^2)$  samples, where  $\tilde{O}$  can hide polylog factors of  $d$  and  $1/\gamma$ .

### 6.5.2 Example 2: Generalization bounds for deep nets using low rank approximations

Some of the early generalization bounds for fully connected nets used the fact that layer matrices are often found to be low rank. (Or perhaps the final matrix minus the initialization.) We give a simple proof of such a result.

Realize that an  $h \times h$  matrix of rank  $r$  has effectively  $2hr$  parameters despite having  $h^2$  entries. We recall that for a square matrix  $A$  the spectral norm (i.e., largest singular value) is denoted  $\|A\|_2$  and sum of squares of singular values is denoted  $\|A\|_F^2$  where  $\|\cdot\|_F$  is also called *Frobenius norm*. The ratio  $\|A\|_F^2 / \|A\|_2^2$  is called *stable rank*, and it is clearly upper bounded by the rank. Often the layers of the trained net have low stable rank even though rank per se is high.

**Theorem 6.5.6.** <sup>(12)</sup> *For a depth- $d$  ReLU net with hidden layers of equal width  $h$  and single coordinate output, let  $A^1, A^2, \dots, A^d$  be weight matrices and  $\gamma$  be the output margin on a training set  $S$  of size  $m$ . Then the generalization error can be bounded by*

$$\tilde{O} \left( \sqrt{\frac{hd^2 \max_{x \in S} \|x\| \prod_{i=1}^d \|A^i\|_2^2 \sum_{i=1}^d \frac{\|A^i\|_F^2}{\|A^i\|_2^2}}{\gamma^2 m}} \right).$$

The second part of this expression ( $\sum_{i=1}^d \frac{\|A^i\|_F^2}{\|A^i\|_2^2}$ ) is sum of stable ranks of the layers, a natural measure of their true parameter count. The first part ( $\prod_{i=1}^d \|A^i\|_2^2$ ) is related to the Lipschitz constant of the network, namely, the maximum norm of the vector it can produce if the input is a unit vector. The Lipschitz constant of a matrix operator  $B$  is just its spectral norm  $\|B\|_2$ . Since the network applies a sequence of matrix operations interspersed with ReLU, and ReLU is 1-Lipschitz we conclude that the Lipschitz constant of the full network is at most  $\prod_{i=1}^d \|A^i\|_2$ .

<sup>12</sup> Behnam Neyshabur, Srinadh Bhojanapalli, David McAllester, and Nathan Srebro. A pac-bayesian approach to spectrally-normalized margin bounds for neural networks. *ICLR*, 2018

To prove Theorem 6.5.6 we use the following lemma to compress the matrix at each layer to a matrix of smaller rank. Since a matrix of rank  $r$  can be expressed as the product of two matrices of inner dimension  $r$ , it has  $2hr$  parameters (instead of the trivial  $h^2$ ). (Furthermore, the parameters can be discretized via trivial rounding to get a compression with discrete parameters as needed by Definition 6.5.1.)

**Lemma 6.5.7.** *For any matrix  $A \in \mathbb{R}^{m \times n}$ , let  $\hat{A}$  be the truncated version of  $A$  where singular values that are smaller than  $\delta \|A\|_2$  are removed. Then  $\|\hat{A} - A\|_2 \leq \delta \|A\|_2$  and  $\hat{A}$  has rank at most  $\|A\|_F^2 / (\delta^2 \|A\|_2^2)$ .*

*Proof.* Let  $r$  be the rank of  $\hat{A}$ . By construction, the maximum singular value of  $\hat{A} - A$  is at most  $\delta \|A\|_2$ . Since the remaining singular values are at least  $\delta \|A\|_2$ , we have  $\|A\|_F \geq \|\hat{A}\|_F \geq \sqrt{r} \delta \|A\|_2$ .  $\square$

For each  $i$  replace layer  $i$  by its compression using the above lemma, with  $\delta = \gamma(3d\|x\| \prod_{i=1}^d \|A^i\|_2)^{-1}$ . How much error does this introduce at each layer and how much does it affect the output after passing through the intermediate layers (and getting magnified by their Lipschitz constants)? Since  $A - \hat{A}^i$  has spectral norm (i.e., Lipschitz constant) at most  $\delta \|A^i\|_2$ , the error at the output due to changing layer  $i$  in isolation is at most  $\prod_{j=i+1}^d \|A^j\|_2 \cdot \delta \|A^i\|_2 \cdot \prod_{j=1}^{i-1} \|A^j\|_2 \cdot \|x\| \leq \gamma/3d$ . Rest of the proof is left to the reader and generalization bound follows immediately from Theorem 6.5.3.

**Problem 6.5.8.** *Complete the above proof using a simple induction (see <sup>13</sup> if needed) to show the total error incurred in all layers is strictly bounded by  $\gamma$ . That is, for an input  $x$ , the change in the deep net output is smaller than  $\gamma$  after replacing every weight matrix  $A^i$  with its truncated version  $\hat{A}^i$ .*

<sup>13</sup> Behnam Neyshabur, Srinadh Bhojanapalli, David McAllester, and Nathan Srebro. A pac-bayesian approach to spectrally-normalized margin bounds for neural networks. *ICLR*, 2018

## 6.6 PAC-Bayes bounds

These bounds due to McAllester (1999) [McA99] are in principle the tightest, meaning previous bounds in this chapter are its subcases. They are descended from an old philosophical tradition of considering the logical foundations for belief systems, which often uses Bayes' Theorem. For example, in the 18th century, Laplace sought to give meaning to questions like "What is the probability that the sun will rise tomorrow?" The answer to this question depends upon the person's prior beliefs (e.g., degree of scientific knowledge) as well as their empirical observation that the sun has risen every day in their lifetime.

Coming back to ML, PAC-Bayes bounds assume that experimenter (i.e. machine learning expert) has some prior distribution  $P$  over the hypothesis  $\mathcal{H}$ . If asked to classify without seeing any concrete training data, the experimenter would pick a hypothesis  $h$  according

to  $P$  (denoted  $h \sim P$ ) and classify using it  $h$ . After seeing the training data and running computations, the experimenter's distribution changes to the posterior  $Q$ , meaning now if asked to classify they would pick  $h \sim Q$  and use that. Thus the expected training loss is

$$\mathbb{E}_{h \sim Q} [L_{\mathcal{D}}(h)].$$

We say  $Q$  is a *valid* posterior with respect to  $P$  if any hypothesis  $h$  that gets zero probability under  $P$  also gets zero probability under  $Q$ .

**Example 6.6.1.**  $P$  could be the standard normal distribution, which assigns nonzero probability to every vector. For any sample set  $S$ , we could let  $Q$  be the distribution on parameter vectors obtained by vanilla deep learning using  $S$ : that is, initialize parameters using random Gaussian, and train with SGD with a predetermined learning rate schedule. Since SGD is a stochastic process (due to randomness of batches) it leads to a natural distribution  $Q$  on trained classifiers at the end of training. Notice,  $P, Q$  are a valid pair because  $P$  assigns nonzero probability to every  $h$ . As this example emphasizes, choice of  $P$  and  $Q$  can be arbitrary and the bound will hold for every fixed choice.

The following form of PAC-Bayes bound is from <sup>14</sup>.

**Theorem 6.6.2** (PAC-Bayes bound). *Let  $\mathcal{D}$  be the data distribution and  $P$  be a prior distribution over hypothesis class  $\mathcal{H}$  and  $\delta > 0$ . If  $S$  is a set of i.i.d. samples of size  $m$  from  $\mathcal{D}$  and  $Q$  is any valid posterior (possibly depending arbitrarily on  $S$ ) then  $\Delta_S(Q) = \mathbb{E}_{h \sim Q} [L_{\mathcal{D}}(h) - \hat{L}_S(h)]$  satisfies the following bound with probability  $1 - \delta$ ,*

$$\Delta_S(Q) \leq \sqrt{\frac{D(Q||P) + \ln(2m/\delta)}{2(m-1)}},$$

where  $D(Q||P) = \mathbb{E}_{h \sim Q} [\ln(Q(h)/P(h))]$  is the so-called KL-divergence<sup>15</sup>.

In other words, generalization error can be upper bounded using the (square root of) KL-divergence of the distributions, plus some terms that arise from concentration bounds.

**Example 6.6.3.** *Suppose  $h$  is any classifier and  $P, Q$  are the distribution that assigns probability 1 to  $h$  and zero to all other hypotheses. Then  $D(Q||P) = 0$ , and by Hoeffding bound we have  $\Delta_S(Q) = \Delta_S(h) \leq \sqrt{\frac{\log(1/\delta)}{2m}}$ . The inequality in PAC-Bayes bound is satisfied.*

**Problem 6.6.4.** *Derive the union bound Theorem 6.2.1 using PAC-Bayes.*

Now we're ready to prove Theorem 6.6.2. For simplicity, we only prove a slightly weaker version:

$$\Delta_S(Q) \leq \sqrt{\frac{2(D(Q||P) + \ln(2/\delta))}{m}}.$$

<sup>14</sup> John Langford. *Quantitatively tight sample complexity bounds*. PhD Thesis CMU, 2002

<sup>15</sup> This is a measure of distance between distributions, meaningful when  $P$  dominates  $Q$ , in the sense that every  $h$  with nonzero probability in  $Q$  also has nonzero probability in  $P$ .

We'll also make a simplifying assumption:  $z = 2\sqrt{m}\Delta_S(h)$  which is a sum of  $m$  iid random variables, behaves *exactly* like a normal distribution  $\mathcal{N}(0, 1)$ . Of course, in truth this is true only in the limit  $m \rightarrow \infty$ . This assumption can be removed by using a more quantitative argument with Hoeffding bound. We will also use a mathematical fact from calculus: the expected value of  $e^{x^2/4}$  is  $\sqrt{2}$  when  $x$  is drawn from  $\mathcal{N}(0, 1)$ .

*Proof.* (Theorem 6.6.2, weaker version) Rearranging the expression in the theorem statement, we see that it gives an upper bound of  $\ln(2/\delta)$  on  $(m/2) \mathbb{E}_{h \sim Q} [\Delta(h)^2] - D(Q||P)$ . By Jensen's inequality<sup>16</sup> applied to the square function  $f(x) = x^2$ , this expression is at most  $(m/2) \mathbb{E}_{h \sim Q} [\Delta(h)^2] - D(Q||P)$ . We show this is upper bounded by  $\ln(2/\delta)$ . The steps are:

$$\begin{aligned} &= \mathbb{E}_{h \sim Q} \left[ (m/2) \Delta(h)^2 - \ln(Q(h)/P(h)) \right] \\ &= \mathbb{E}_{h \sim Q} \left[ \ln \left( \exp((m/2) \Delta(h)^2) \cdot P(h)/Q(h) \right) \right] \\ &\leq \ln \left( \mathbb{E}_{h \sim Q} \left[ \exp((m/2) \Delta(h)^2) \cdot P(h)/Q(h) \right] \right), \end{aligned}$$

where the last inequality uses Jensen's inequality along with the concavity of  $\ln$ . Also, since taking expectation over  $h \sim Q$  is effectively like summing with a weighting by  $Q(h)$ , we have<sup>17</sup>

$$\ln \mathbb{E}_{h \sim Q} \left[ \exp((m/2) \Delta(h)^2) \cdot P(h)/Q(h) \right] = \ln \mathbb{E}_{h \sim P} \left[ \exp((m/2) \Delta(h)^2) \right].$$

Recapping, we thus have the following for a fixed dataset  $S$ :

$$(m/2) \mathbb{E}_{h \sim Q} [\Delta(h)^2] - D(Q||P) \leq \ln \left( \mathbb{E}_{h \sim P} \left[ e^{(m/2) \Delta(h)^2} \right] \right) \quad (6.8)$$

Now using the fact that belief  $P$  was fixed before seeing  $S$  (i.e., is independent of  $S$ ):

$$\mathbb{E}_S \left[ \mathbb{E}_{h \sim P} \left[ e^{(m/2) \Delta(h)^2} \right] \right] = \mathbb{E}_{h \sim P} \left[ \mathbb{E}_S \left[ e^{(m/2) \Delta(h)^2} \right] \right] = \sqrt{2} \leq 2.$$

Thus, (1) implies that with probability  $1 - \delta$  over  $S$ ,

$$\mathbb{E}_{h \sim P} \left[ e^{(m/2) \Delta(h)^2} \right] \leq 2/\delta \quad (6.9)$$

and now by taking logarithm of both sides the proof is completed.  $\square$

## 6.7 Exercises

1. Assume the loss function  $\ell$  is 1-Lipschitz. Consider the kernel classifier of the form  $h(x) = z^\top G^{-1}y$  we studied in Section 3.2

<sup>16</sup> Jensen's Inequality: For a concave function  $f$  and random variable  $X$ ,  $\mathbb{E}[f(X)] \leq f(\mathbb{E}[X])$ . For convex function the inequality is reversed.

<sup>17</sup> Often when you see KL-divergence in machine learning, you will see this trick being used to switch the distribution over which expectation is taken!

where  $G$  is the  $n \times n$  kernel matrix,  $y$  is the labels and  $z$  is the column vector whose  $i$ -th coordinate is  $K(x, x_i)$ . Prove that its Rademacher complexity is upper bounded  $\sqrt{2y^\top G y \cdot \text{Tr}(G)}/n$ . (Hint: view kernel classifier as a linear classifier in RKHS)



## Tractable Landscapes for Nonconvex Optimization

Deep learning relies on optimizing complicated, nonconvex loss functions. Finding the global minimum of a nonconvex objective is NP-hard in the worst case. However in deep learning simple algorithms such as stochastic gradient descent often the objective value to zero or near-zero at the end. This chapter focuses on the *optimization landscape* defined by a nonconvex objective and identifies properties of these landscapes that allow simple optimization algorithms to find global minima (or near-minima). These properties thus far apply to simpler nonconvex problems than deep learning, and it is open how to analyse deep learning with such landscape analysis.

*Warm-up: Convex Optimization* To understand optimization landscape, one can first look at optimizing a convex function. If a function  $f(w)$  is convex, then it satisfies many nice properties, including

$$\forall \alpha \in [0, 1], w, w', f(\alpha w + (1 - \alpha)w') \leq \alpha f(w) + (1 - \alpha)f(w'). \quad (7.1)$$

$$\forall w, w', f(w') \geq f(w) + \langle \nabla f(w), w' - w \rangle. \quad (7.2)$$

These equations characterize important geometric properties of the objective function  $f(w)$ . In particular, Equation (7.1) shows that all the global minima of  $f(w)$  must be connected, because if  $w, w'$  are both globally optimal, anything on the segment  $\alpha w + (1 - \alpha)w'$  must also be optimal. Such properties are important because it gives a characterization of all the global minima. Equation (7.2) shows that every point with  $\nabla f(w) = 0$  must be a global minimum, because for every  $w'$  we have  $f(w') \geq f(w) + \langle \nabla f(w), w' - w \rangle \geq f(w)$ . Such properties are important because it connects a local property (gradient being 0) to global optimality.

In general, optimization landscape looks for properties of the objective function that characterizes its local/global optimal points (such as Equation (7.1)) or connects local properties with global optimality (such as Equation (7.2)).

## 7.1 Preliminaries and challenges in nonconvex landscapes

We have been discussing global/local minimum informally, here we first give a precise definition:

**Definition 7.1.1** (Global/Local minimum). *For an objective function  $f(w) : \mathbb{R}^d \rightarrow \mathbb{R}$ , a point  $w^*$  is a global minimum if for every  $w$  we have  $f(w^*) \leq f(w)$ . A point  $w$  is a local minimum/maximum if there exists a radius  $\epsilon > 0$  such that for every  $\|w' - w\|_2 \leq \epsilon$ , we have  $f(w) \leq f(w')$  ( $f(w) \geq f(w')$  for local maximum). A point  $w$  with  $\nabla f(w) = 0$  is called a critical point, for smooth functions all local minimum/maximum are critical points.*

Throughout the chapter, we will always work with functions whose global minimum exists, and use  $f(w^*)$  to denote the optimal value of the function<sup>1</sup>. For simplicity we focus on optimization problems that do not have any constraints ( $w \in \mathbb{R}^d$ ). It is possible to extend everything in this chapter to optimization with nondegenerate equality constraints, which would require definitions of gradient and Hessians with respect to a manifold and is out of the scope for this book.

<sup>1</sup> Even though there might be multiple global minima  $w^*$ , the value  $f(w^*)$  is unique by definition.

*Spurious local minimum* The first obstacle in nonconvex optimization is a *spurious local minimum*.

**Definition 7.1.2** (Spurious local minimum). *For an objective function  $f(w) : \mathbb{R}^d \rightarrow \mathbb{R}$ , a point  $w$  is a spurious local minimum if it is a local minimum, but  $f(w) > f(w^*)$ .*

Many of the simple optimization algorithms are based on the idea of local search, thus are not able to escape from a spurious local minimum. As we will later see, many nonconvex objectives do not have spurious local minima.

*Saddle points* The second obstacle in nonconvex optimization is a *saddle point*. The simplest example of a saddle point is  $f(w) = w_1^2 - w_2^2$  at the point  $w = (0, 0)$ . In this case, if  $w$  moves along direction  $(\pm 1, 0)$ , the function value increases; if  $w$  moves along direction  $(0, \pm 1)$ , the function value decreases.

**Definition 7.1.3** (Saddle point). *For an objective function  $f(w) : \mathbb{R}^d \rightarrow \mathbb{R}$ , a point  $w$  is a saddle point if  $\nabla f(w) = 0$ , and for every radius  $\epsilon > 0$ , there exists  $w^+, w^-$  within distance  $\epsilon$  of  $w$  such that  $f(w^-) < f(w) < f(w^+)$ .*

This definition covers all cases but makes it very hard to verify whether a point is a saddle point. In most cases, it is possible to tell whether a point is a saddle point, local minimum or local maximum based on its Hessian.

**Claim 7.1.4.** For an objective function  $f(w) : \mathbb{R}^d \rightarrow \mathbb{R}$  and a critical point  $w$  ( $\nabla f(w) = 0$ ), we know

- If  $\nabla^2 f(w) \succ 0$ ,  $w$  is a local minimum.
- If  $\nabla^2 f(w) \prec 0$ ,  $w$  is a local maximum.
- If  $\nabla^2 f(w)$  has both a positive and a negative eigenvalue,  $w$  is a saddle point.

These criteria are known as second order sufficient conditions in optimization. Intuitively, one can prove this claim by looking at the second-order Taylor expansion. The three cases in the claim do not cover all the possible Hessian matrices. The remaining cases are considered to be degenerate, and can either be a local minimum, local maximum or a saddle point<sup>2</sup>.

*Flat regions* Even if a function does not have any spurious local minima or saddle point, it can still be nonconvex, see Figure 7.1. In high dimensions such functions can still be very hard to optimize. The main difficulty here is that even if the norm  $\|\nabla f(w)\|_2$  is small, unlike convex functions one cannot conclude that  $f(w)$  is close to  $f(w^*)$ . However, often in such cases one can hope the function  $f(w)$  to satisfy some relaxed notion of convexity, and design efficient algorithms accordingly. We discuss one of such cases in Section 7.2.

<sup>2</sup> One can consider the  $w = 0$  point of functions  $w^4$ ,  $-w^4$ ,  $w^3$ , and it is a local minimum, maximum and saddle point respectively.

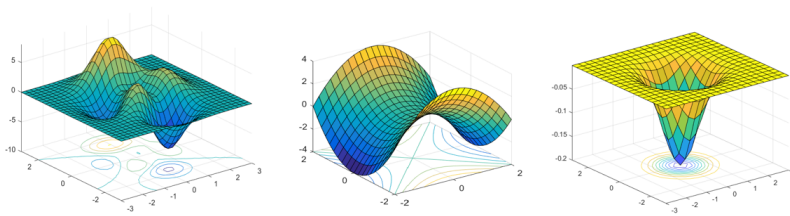


Figure 7.1: Obstacles for non-convex optimization. From left to right: local minimum, saddle point and flat region.

## 7.2 Cases with a unique global minimum

We first consider the case that is most similar to convex objectives. In this section, the objective functions we look at have no spurious local minima or saddle points. In fact, in our example the objective is only going to have a unique global minimum. The only obstacle in optimizing these functions is that points with small gradients may not be near-optimal.

The main idea here is to identify properties of the objective and also a *potential function*, such that the potential function keeps de-

creasing as we run simple optimization algorithms such as gradient descent. Many properties were used in previous literature, including

**Definition 7.2.1.** *Let  $f(w)$  be an objective function with a unique global minimum  $w^*$ , then*

*Polyak-Lojasiewicz  $f$  satisfies Polyak-Lojasiewicz if there exists a value  $\mu > 0$  such that for every  $w$ ,  $\|\nabla f(x)\|_2^2 \geq \mu(f(w) - f(w^*))$ .*

*weakly-quasi-convex  $f$  is weakly-quasi-convex if there exists a value  $\tau > 0$  such that for every  $w$ ,  $\langle \nabla f(w), w - w^* \rangle \geq \mu(f(w) - f(w^*))$ .*

*Restricted Secant Inequality (RSI)  $f$  satisfies RSI if there exists a value  $\tau$  such that for every  $w$ ,  $\langle \nabla f(w), w - w^* \rangle \geq \mu \|w - w^*\|_2^2$ .*

Any one of these three properties can imply fast convergence together with some smoothness of  $f$ .

**Claim 7.2.2.** *If an objective function  $f$  satisfies one of Polyak-Lojasiewicz, weakly-quasi-convex or RSI, and  $f$  is smooth<sup>3</sup>, then gradient descent converges to global minimum with a geometric rate<sup>4</sup>.*

Intuitively, Polyak-Lojasiewicz condition requires that the gradient to be nonzero for any point that is not a global minimum, therefore one can always follow the gradient and further decrease the function value. This condition can also work in some settings when the global minimum is not unique. Weakly-quasi-convex and RSI are similar in the sense that they both require the (negative) gradient to be correlated with the correct direction - direction from the current point  $w$  to the global minimum  $w^*$ .

In this section we are going to use generalized linear model as an example to show how some of these properties can be used.

### 7.2.1 Generalized linear model

In generalized linear model (also known as isotonic regression) [KS09, KSK11], the input consists of samples  $\{x^{(i)}, y^{(i)}\}$  that are drawn from distribution  $\mathcal{D}$ , where  $(x, y) \sim \mathcal{D}$  satisfies

$$y = \sigma(w_*^\top x) + \epsilon. \quad (7.3)$$

Here  $\sigma : \mathbb{R} \rightarrow \mathbb{R}$  is a known monotone function,  $\epsilon$  is a noise that satisfies  $\mathbb{E}[\epsilon|x] = 0$ , and  $w_*$  is the unknown parameter that we are trying to learn.

In this case, it is natural to consider the following expected loss

$$L(w) = \frac{1}{2} \mathbb{E}_{(x,y) \sim \mathcal{D}} \left[ (y - \sigma(w^\top x))^2 \right]. \quad (7.4)$$

<sup>3</sup> Polyak-Lojasiewicz and RSI requires standard smoothness definition as in Equation (2.2), weakly-quasi-convex requires a special smoothness property detailed in [HMR18].

<sup>4</sup> The potential functions for Polyak-Lojasiewicz and weakly-quasi-convex are function value  $f$ ; potential function for RSI is the squared distance  $\|w - w_*\|_2^2$ .

Of course, in practice one can only access the training loss which is an average over the observed  $\{x^{(i)}, y^{(i)}\}$  pairs. For simplicity we work with the expected loss here. The difference between the two losses can be bounded using techniques in Chapter ??.

Generalized linear model can be viewed as learning a single neuron where  $\sigma$  is its nonlinearity.

We will give high level ideas on how to prove properties such as weakly-quasi-convex or RSI for generalized linear model. First we rewrite the objective as:

$$\begin{aligned} L(w) &= \frac{1}{2} \mathbb{E}_{(x,y) \sim \mathcal{D}} \left[ (y - \sigma(w^\top x))^2 \right] \\ &= \frac{1}{2} \mathbb{E}_{(x,\epsilon)} \left[ (\epsilon + \sigma(w_*^\top x) - \sigma(w^\top x))^2 \right]. \\ &= \frac{1}{2} \mathbb{E}_\epsilon [\epsilon^2] + \frac{1}{2} \mathbb{E}_x \left[ (\sigma(w_*^\top x) - \sigma(w^\top x))^2 \right]. \end{aligned}$$

Here the second equality uses Definition of the model (Equation (7.3)), and the third equality uses the fact that  $\mathbb{E}[\epsilon|x] = 0$  (so there are no cross terms). This decomposition is helpful as the first term  $\frac{1}{2} \mathbb{E}_\epsilon [\epsilon^2]$  is now just a constant.

Now we can take the derivative of the objective:

$$\nabla L(w) = \mathbb{E}_x \left[ (\sigma(w^\top x) - \sigma(w_*^\top x)) \sigma'(w^\top x) x \right].$$

Notice that both weakly-quasi-convex and RSI requires that the objective to be correlated with  $w - w_*$ , so we compute

$$\langle \nabla L(w), w - w_* \rangle = \mathbb{E}_x \left[ (\sigma(w^\top x) - \sigma(w_*^\top x)) \sigma'(w^\top x) (w^\top x - w_*^\top x) \right].$$

The goal here is to show that the RHS is bigger than 0. A simple way to see that is to use the intermediate value theorem:  $\sigma(w^\top x) - \sigma(w_*^\top x) = \sigma'(\xi)(w^\top x - w_*^\top x)$ , where  $\xi$  is a value between  $w^\top x$  and  $w_*^\top x$ . Then we have

$$\langle \nabla L(w), w - w_* \rangle = \mathbb{E}_x \left[ \sigma'(\xi) \sigma'(w^\top x) (w^\top x - w_*^\top x)^2 \right].$$

In the expectation in the RHS, both derivatives ( $\sigma'(\xi), \sigma'(w^\top x)$ ) are positive as  $\sigma$  is monotone, and  $(w^\top x - w_*^\top x)^2$  is clearly nonnegative. By making more assumptions on  $\sigma$  and the distribution of  $x$ , it is possible to lowerbound the RHS in the form required by either weakly-quasi-convex or RSI. We leave this as an exercise.

### 7.2.2 Alternative objective for generalized linear model

There is another way to find  $w_*$  for generalized linear model that is more specific to this setting. In this method, one estimate a different

“gradient” for generalized linear model:

$$\nabla g(w) = \mathbb{E}_{x,y} [(\sigma(w^\top x) - y)x] = \mathbb{E}_x [(\sigma(w^\top x) - \sigma(w_*^\top x))x]. \quad (7.5)$$

The first equation gives a way to estimate this “gradient”. The main difference here is that in the RHS we no longer have a factor  $\sigma'(w^\top x)$  as in  $\nabla L(w)$ . Of course, it is unclear why this formula is the gradient of some function  $g$ , but we can construct the function  $g$  in the following way:

Let  $\tau(x)$  be the integral of  $\sigma(x)$ :  $\tau(x) = \int_0^x \sigma(x') dx'$ . Define  $g(w) := \mathbb{E}_x [\tau(w^\top x) - \sigma(w_*^\top x)w^\top x]$ . One can check  $\nabla g(w)$  is indeed the function in Equation (7.5). What’s very surprising is that  $g(w)$  is actually a convex function with  $\nabla g(w_*) = 0$ ! This means that  $w_*$  is a global minimum of  $g$  and we only need to follow  $\nabla g(w)$  to find it. Nonconvex optimization is unnecessary here.

Of course, this technique is quite special and uses a lot of structure in generalized linear model. However similar ideas were also used in <sup>5</sup> to learn a single neuron. In general, when one objective is hard to analyze, it might be easier to look for an alternative objective that has the same global minimum but easier to optimize.

5

### 7.3 Symmetry, saddle points and locally optimizable functions

In the previous section, we saw some conditions that allow nonconvex objectives to be optimized efficiently. However, such conditions often do not apply to neural networks, or more generally any function that has some symmetry properties.

More concretely, consider a two-layer neural network  $h_\theta(x) : \mathbb{R}^d \rightarrow \mathbb{R}$ . The parameters  $\theta$  is  $(w_1, w_2, \dots, w_k)$  where  $w_i \in \mathbb{R}^d$  represents the weight vector of the  $i$ -th neuron. The function can be evaluated as  $h_\theta(x) = \sum_{i=1}^k \sigma(\langle w_i, x \rangle)$ , where  $\sigma$  is a nonlinear activation function. Given a dataset  $(x^{(1)}, y^{(1)}), \dots, (x^{(n)}, y^{(n)}) \stackrel{\text{i.i.d.}}{\sim} \mathcal{D}$ , one can define the training loss and expected loss as in Chapter 1. Now the objective for this neural network  $f(\theta) = L(h_\theta) = \mathbb{E}_{(x,y) \sim \mathcal{D}} [\ell((x,y), h_\theta)]$  has *permutation symmetry*. That is, for any permutation  $\pi(\theta)$  that permutes the weights of the neurons, we know  $f(\theta) = f(\pi(\theta))$ .

The symmetry has many implications. First, if the global minimum  $\theta^*$  is a point where not all neurons have the same weight vector (which is very likely to be true), then there must be equivalent global minimum  $f(\pi(\theta^*))$  for every permutation  $\pi$ . An objective with this symmetry must also be nonconvex, because if it were convex, the point  $\bar{\theta} = \frac{1}{k!} \sum_{\pi} \pi(\theta^*)$  (where  $\pi$  sums over all the permutations) is a convex combination of global minima, so it must also be a global minimum. However, for  $\bar{\theta}$  the weight vectors of the neurons are all

equal to  $\frac{1}{k} \sum_{i=1}^k w_i$  (where  $w_i$  is the weight of  $i$ -th neuron in  $\theta^*$ ), so  $h_{\bar{\theta}}(x) = k\sigma(\langle \frac{1}{k} \sum_{i=1}^k w_i, x \rangle)$  is equivalent to a neural network with a single neuron. In most cases a single-neuron network should not achieve the global minimum, so by proof of contradiction we know  $f$  should not be convex.

It's also possible to show that functions with symmetry must have saddle points<sup>6</sup>. Therefore to optimize such a function, the algorithm needs to be able to either avoid or escape from saddle points. More concretely, one would like to find a *second order stationary point*.

<sup>6</sup> Except some degenerate cases such as constant functions.

**Definition 7.3.1** (Second order stationary point (SOSP)). *For an objective function  $f(w) : \mathbb{R}^d \rightarrow \mathbb{R}$ , a point  $w$  is a second order stationary point if  $\nabla f(w) = 0$  and  $\nabla^2 f(w) \succeq 0$ .*

The conditions for second order stationary point are known as the second order necessary conditions for a local minimum. Of course, generally an optimization algorithm will not be able to find an exact second order stationary point (just like in Section ?? we only show gradient descent finds a point with small gradient, but not 0 gradient). The optimization algorithms can be used to find an approximate second order stationary point:

**Definition 7.3.2** (Approximate second order stationary point). *For an objective function  $f(w) : \mathbb{R}^d \rightarrow \mathbb{R}$ , a point  $w$  is a  $(\epsilon, \gamma)$ -second order stationary point (later abbreviated as  $(\epsilon, \gamma)$ -SOSP) if  $\|\nabla f(w)\|_2 \leq \epsilon$  and  $\lambda_{\min}(\nabla^2 f(w)) \geq -\gamma$ .*

Later in Chapter ?? we will show that simple variants of gradient descent can in fact find  $(\epsilon, \gamma)$ -SOSPs efficiently.

Now we are ready to define a class of functions that can be optimized efficiently and allow symmetry and saddle points.

**Definition 7.3.3** (Locally optimizable functions). *An objective function  $f(w)$  is locally optimizable, if for every  $\tau > 0$ , there exists  $\epsilon, \gamma = \text{poly}(\tau)$  such that every  $(\epsilon, \gamma)$ -SOSP  $w$  of  $f$  satisfies  $f(w) \leq f(w_*) + \tau$ .*

Roughly speaking, an objective function is locally optimizable if every local minimum of the function is also a global minimum, and the Hessian of every saddle point has a negative eigenvalue. Similar class of functions were called “strict saddle” or “ridable” in some previous results. Many nonconvex objectives, including matrix sensing [BNS16a, PKCS17, GJZ17a], matrix completion [GLM16a, GJZ17a], dictionary learning [SQW16a], phase retrieval [SQW18], tensor decomposition [GHJY15a], synchronization problems [BBV16] and certain objective for two-layer neural network [GLM18] are known to be locally optimizable.

## 7.4 Case study: top eigenvector of a matrix

In this section we look at a simple example of a locally optimizable function. Given a symmetric PSD matrix  $M \in \mathbb{R}^{d \times d}$ , our goal is to find its top eigenvector (eigenvector that corresponds to the largest eigenvalue). More precisely, using SVD we can write  $M$  as

$$M = \sum_{i=1}^d \lambda_i v_i v_i^\top.$$

Here  $v_i$ 's are orthonormal vectors that are eigenvectors of  $M$ , and  $\lambda_i$ 's are the eigenvalues. For simplicity we assume  $\lambda_1 > \lambda_2 \geq \lambda_3 \geq \dots \geq \lambda_d \geq 0$ <sup>7</sup>.

There are many objective functions whose global optima gives the top eigenvector. For example, using basic definition of spectral norm, we know for PSD matrix  $M$  the global optima of

$$\max_{\|x\|_2=1} x^\top M x$$

is the top eigenvector of  $M$ . However, this formulation requires a constraint. We instead work with an unconstrained version whose correctness follows from Eckhart-Young Theorem:

$$\min_{x \in \mathbb{R}^d} f(x) := \frac{1}{4} \|M - x x^\top\|_F^2. \quad (7.6)$$

Note that this function does have a symmetry in the sense that  $f(x) = f(-x)$ . Under our assumptions, the only global minima of this function are  $x = \pm \sqrt{\lambda_1} v_1$ . We are going to show that these are also the only second order stationary points. We will give two proof strategies that are commonly used to prove the locally optimizable property.

### 7.4.1 Characterizing all critical points

The first idea is simple – we will just try to solve the Equation  $\nabla f(x) = 0$  to get the position of all critical points; then for the critical points that are not the desired global minimum, try to prove that they are local maximum or saddle points.

*Computing gradient and Hessian* Before we solve the equation  $\nabla f(x) = 0$  for the objective function  $f(x)$  defined in Equation (7.6), we first give a simple way of computing the gradient and Hessian. We will first expand  $f(x + \delta)$  (where  $\delta$  should be thought of as a small

<sup>7</sup> Note that the only real assumption here is  $\lambda_1 > \lambda_2$ , so the top eigenvector is unique. Other inequalities are without loss of generality.



perturbation):

$$\begin{aligned}
f(x + \delta) &= \frac{1}{4} \|M - (x + \delta)(x + \delta)^\top\|_F^2 \\
&= \frac{1}{4} \|M - xx^\top - (x\delta^\top + \delta x^\top) - \delta\delta^\top\|_F^2 \\
&= \frac{1}{4} \|M - xx^\top\|_F^2 - \frac{1}{2} \langle M - xx^\top, x\delta^\top + \delta x^\top \rangle \\
&\quad + \left[ \frac{1}{4} \|x\delta^\top + \delta x^\top\|_F^2 - \frac{1}{2} \langle M - xx^\top, \delta\delta^\top \rangle \right] + o(\|\delta\|_2^2).
\end{aligned}$$

Note that in the last step, we have collected the terms based on the degree of  $\delta$ , and ignored all the terms that are smaller than  $o(\|\delta\|_2^2)$ . We can now compare this expression with the Taylor's expansion of  $f(x + \delta)$ :

$$f(x + \delta) = f(x) + \langle \nabla f(x), \delta \rangle + \frac{1}{2} \delta^\top [\nabla^2 f(x)] \delta + o(\|\delta\|_2^2).$$

By matching terms, we immediately have

$$\begin{aligned}
\langle \nabla f(x), \delta \rangle &= -\frac{1}{2} \langle M - xx^\top, x\delta^\top + \delta x^\top \rangle, \\
\delta^\top [\nabla^2 f(x)] \delta &= \frac{1}{2} \|x\delta^\top + \delta x^\top\|_F^2 - \langle M - xx^\top, \delta\delta^\top \rangle.
\end{aligned}$$

These can be simplified to give the actual gradient and Hessian<sup>8</sup>

$$\nabla f(x) = (xx^\top - M)x, \quad \nabla^2 f(x) = \|x\|_2^2 I + 2xx^\top - M. \quad (7.7)$$

<sup>8</sup> In fact in the next subsection we will see it is often good enough to know how to compute  $\langle \nabla f(x), \delta \rangle$  and  $\delta^\top [\nabla^2 f(x)] \delta$ .

*Characterizing critical points* Now we can execute the original plan. First set  $\nabla f(x) = 0$ , we have

$$Mx = xx^\top x = \|x\|_2^2 x.$$

Luckily, this is a well studied equation because we know the only solutions to  $Mx = \lambda x$  are if  $\lambda$  is an eigenvalue and  $x$  is (a scaled version) of the corresponding eigenvector. Therefore we know  $x = \pm \sqrt{\lambda_i} v_i$  or  $x = 0$ . These are the only critical points of the objective function  $f(x)$ .

Among these critical points,  $x = \pm \sqrt{\lambda_1} v_1$  are our intended solutions. Next we need to show for every other critical point, its Hessian has a negative eigendirection. We will do this for  $x = \pm \sqrt{\lambda_i} v_i (i > 1)$ . By definition, it suffices to show there exists a  $\delta$  such that  $\delta^\top [\nabla^2 f(x)] \delta < 0$ . The main step of the proof involves guessing what is this direction  $\delta$ . In this case we will choose  $\delta = v_1$  (we will give more intuitions about how to choose such a direction in the next subsection).

When  $x = \pm\sqrt{\lambda_i}v_i$ , and  $\delta = v_1$ , we have

$$\delta^\top [\nabla^2 f(x)]\delta = v_1^\top [ \|\sqrt{\lambda_i}v_i\|_2^2 I + 2\lambda_i v_i v_i^\top - M ] v_1 = \lambda_i - \lambda_1 < 0.$$

Here the last step uses the fact that  $v_i$ 's are orthonormal vectors and  $v_1^\top M v_1 = \lambda_1$ . The proof for  $x = 0$  is very similar. Combining all the steps above, we proved the following claim:

**Claim 7.4.1** (Properties of critical points). *The only critical points of  $f(x)$  are of the form  $x = \pm\sqrt{\lambda_i}v_i$  or  $x = 0$ . For all critical points except  $x = \pm\sqrt{\lambda_1}v_1$ ,  $\nabla^2 f(x)$  has a negative eigenvalue.*

This claim directly implies that the only second order stationary points are  $x = \pm\sqrt{\lambda_1}v_1$ , so all second order stationary points are also global minima.

#### 7.4.2 Finding directions of improvements

The approach in Section 7.4.1 is straight-forward. However, in more complicated problems it is often infeasible to enumerate all the solutions for  $\nabla f(x) = 0$ . What we proved in Section 7.4.1 is also not strong enough for showing  $f(x)$  is locally optimizable, because we only proved every exact SOSP is a global minimum, and a locally optimizable function requires every approximate SOSP to be close to a global minimum. We will now give an alternative approach that is often more flexible and robust.

For every point  $x$  that is not a global minimum, we define its direction of improvements as below:

**Definition 7.4.2** (Direction of improvement). *For an objective function  $f$  and a point  $x$ , we say  $\delta$  is a direction of improvement (of  $f$  at  $x$ ) if  $|\langle \nabla f(x), \delta \rangle| > 0$  or  $\delta^\top [\nabla^2 f(x)]\delta < 0$ . We say  $\delta$  is an  $(\epsilon, \gamma)$  direction of improvement (of  $f$  at  $x$ ) if  $|\langle \nabla f(x), \delta \rangle| > \epsilon \|\delta\|_2$  or  $\delta^\top [\nabla^2 f(x)]\delta < -\gamma \|\delta\|_2^2$ .*

Intuitively, if  $\delta$  is a direction of improvement for  $f$  at  $x$ , then moving along one of  $\delta$  or  $-\delta$  for a small enough step can decrease the objective function. In fact, if a point  $x$  has a direction of improvement, it cannot be a second order stationary point; if a point  $x$  has an  $(\epsilon, \gamma)$  direction of improvement, then it cannot be an  $(\epsilon, \gamma)$ -SOSP.

Now we can look at the contrapositive of what we were trying to prove in the definition of locally optimizable functions: if every point  $x$  with  $f(x) > f(x^*) + \tau$  has an  $(\epsilon, \gamma)$  direction of improvement, then every  $(\epsilon, \gamma)$ -second order stationary point must satisfy  $f(x) \leq f(x^*) + \delta$ . Therefore, our goal in this part is to find a direction of improvement for every point that is not globally optimal.

For simplicity, we will look at an even simpler version of the top eigenvector problem. In particular, we consider the case where  $M = zz^\top$  is a rank-1 matrix, and  $z$  is a unit vector. In this case, the objective function we defined in Equation (7.6) becomes

$$\min_x f(x) = \frac{1}{4} \|zz^\top - xx^\top\|_F^2. \quad (7.8)$$

The intended global optimal solutions are  $x = \pm z$ . This problem is often called the matrix factorization problem as we are given a matrix  $M = zz^\top$  and the goal is to find a decomposition  $M = xx^\top$ .

<sup>9</sup> Note that we only observe  $M$ , not  $z$ .

Which direction should we move to decrease the objective function? In this problem we only have the optimal direction  $z$  and the current direction  $x$ , so the natural guesses would be  $z, x$  or  $z - x$ . Indeed, these directions are enough:

**Lemma 7.4.3.** *For objective function (7.8), there exists a universal constant  $c > 0$  such that for any  $\tau < 1$ , if neither  $x$  or  $z$  is an  $(c\tau, 1/4)$ -direction of improvement for the point  $x$ , then  $f(x) \leq \tau$ .*

The proof of this lemma involves some detailed calculation. To get some intuition, we can first think about what happens if neither  $x$  or  $z$  is a direction of improvement.

**Lemma 7.4.4.** *For objective function (7.8), if neither  $x$  or  $z$  is a direction of improvement of  $f$  at  $x$ , then  $f(x) = 0$ .*

*Proof.* We will use the same calculation for gradient and Hessian as in Equation (7.7), except that  $M$  is now  $zz^\top$ . First, since  $x$  is not a direction of improvement, we must have

$$\langle \nabla f(x), x \rangle = 0 \implies \|x\|_2^4 = \langle x, z \rangle^2. \quad (7.9)$$

If  $z$  is not a direction of improvement, we know  $z^\top [\nabla^2 f(x)] z \geq 0$ , which means

$$\|x\|^2 + 2\langle x, z \rangle^2 - 1 \geq 0 \implies \|x\|^2 \geq 1/3.$$

Here we used the fact that  $\langle x, z \rangle^2 \leq \|x\|_2^2 \|z\|_2^2 = \|x\|_2^2$ . Together with Equation (7.9) we know  $\langle x, z \rangle^2 = \|x\|_2^4 \geq 1/9$ .

Finally, since  $z$  is not a direction of improvement, we know  $\langle \nabla f(x), z \rangle = 0$ , which implies  $\langle x, z \rangle (\|x\|_2^2 - 1) = 0$ . We have already proved  $\langle x, z \rangle^2 \geq 1/9 > 0$ , thus  $\|x\|_2^2 = 1$ . Again combining with Equation (7.9) we know  $\langle x, z \rangle^2 = \|x\|_2^4 = 1$ . The only two vectors with  $\langle x, z \rangle^2 = 1$  and  $\|x\|_2^2 = 1$  are  $x = \pm z$ .  $\square$

The proof of Lemma 7.4.3 is very similar to Lemma 7.4.4, except we need to allow slacks in every equation and inequality we use. The additional benefit of having the more robust Lemma 7.4.3 is that the

proof is also robust if we don't have access to the exact objective - in settings where only a subset of coordinates of  $zz^\top$ <sup>10</sup>, one can still prove that the objective function is locally optimizable, and hence find  $z$  by nonconvex optimization.

Lemma 7.4.4 and Lemma 7.4.3 both use directions  $x$  and  $z$ . It is also possible to use the direction  $x - z$  when  $\langle x, z \rangle \geq 0$  (and  $x + z$  when  $\langle x, z \rangle < 0$ ). Both ideas can be generalized to handle the case when  $M = ZZ^\top$  where  $Z \in \mathbb{R}^{d \times r}$ , so  $M$  is a rank- $r$  matrix.

<sup>10</sup>This setting is known as *matrix completion* and has been widely applied to recommendation systems.

# 8

## *Escaping Saddle Points*

Gradient descent (GD) and stochastic gradient descent (SGD) are the workhorses of large-scale machine learning. While classical theory focused on analyzing the performance of these methods in *convex* optimization problems, the most notable successes in machine learning have involved *nonconvex* optimization, and a gap has arisen between theory and practice.

Indeed, traditional analyses of GD and SGD show that both algorithms converge to stationary points efficiently. But these analyses do not take into account the possibility of converging to saddle points. Motivated by the geometric characterizations in the last chapter, the central difficulty in solving many nonconvex machine learning problems becomes escaping saddle points.

In this chapter, we will discuss a simple perturbed form of gradient descent, which is capable of escaping saddle points very efficiently. Particularly, in terms of convergence rate and dimension dependence, it is almost as if the saddle points are not there!

### 8.1 Preliminaries

«Chi notes: Many definitions have appeared in the earlier chapters. Coordination may be required.»

In this chapter, we are interested in solving general unconstrained optimization problems of the form:

$$\min_{x \in \mathbb{R}^d} f(x),$$

where  $f$  is a smooth function that can be nonconvex. In particular we assume that  $f$  has Lipschitz gradients and Lipschitz Hessians, which ensures that the gradient and Hessian can not change too rapidly.

**Definition 8.1.1.** A differentiable function  $f$  is  $\ell$ -**gradient Lipschitz** if:

$$\|\nabla f(x_1) - \nabla f(x_2)\| \leq \ell \|x_1 - x_2\| \quad \forall x_1, x_2.$$

**Definition 8.1.2.** A twice-differentiable function  $f$  is  $\rho$ -Hessian Lipschitz if:

$$\left\| \nabla^2 f(x_1) - \nabla^2 f(x_2) \right\| \leq \rho \|x_1 - x_2\| \quad \forall x_1, x_2.$$

For minimization problems, both saddle points and local maxima are clearly undesirable. Our focus will be “saddle points,” although our results also apply directly to local maxima as well. Unfortunately, distinguishing saddle points from local minima for smooth functions is still NP-hard in general [Nesoo]. To avoid these hardness results, we focus on a subclass of saddle points.

**Definition 8.1.3** (strict saddle point). For a twice-differentiable function  $f$ ,  $x$  is a **strict saddle point** if  $x$  is a stationary point and  $\lambda_{\min}(\nabla^2 f(x)) < 0$ .

A generic saddle point must satisfy that  $\lambda_{\min}(\nabla^2 f(x)) \leq 0$ . Being “strict” simply rules out the case where  $\lambda_{\min}(\nabla^2 f(x)) = 0$ . We reformulate our goal as that of finding stationary points that are not strict saddle points.

**Definition 8.1.4** (SOSP). For twice-differentiable function  $f(\cdot)$ ,  $x$  is a **second-order stationary point** if

$$\nabla f(x) = 0, \quad \text{and} \quad \nabla^2 f(x) \succeq 0.$$

**Definition 8.1.5** ( $\epsilon$ -SOSP). For a  $\rho$ -Hessian Lipschitz function  $f(\cdot)$ ,  $x$  is an  $\epsilon$ -**second-order stationary point** if:

$$\|\nabla f(x)\| \leq \epsilon \quad \text{and} \quad \nabla^2 f(x) \succeq -\sqrt{\rho\epsilon} \cdot I.$$

Definition 8.1.5 characterizes an  $\epsilon$ -approximate version of SOSP so that we can discuss rates. The condition on the Hessian in Definition 8.1.5 uses the Hessian Lipschitz parameter  $\rho$  to retain a single accuracy parameter and to match the units of the gradient and Hessian<sup>1</sup>, following the convention of [NP06].

## 8.2 Perturbed Gradient Descent

According to the update equations, Gradient Descent (GD) makes a non-zero step only when the gradient is non-zero, and thus in the nonconvex setting it will be stuck at saddle points if initialized there. We thus consider a simple variant of GD which adds perturbations to the iterates at each step.

$$x_{t+1} \leftarrow x_t - \eta(\nabla f(x_t) + \zeta_t), \quad \zeta_t \sim \mathcal{N}(0, (r^2/d)I)$$

At each iteration, Perturbed Gradient Descent (PGD) is almost the same as gradient descent, except it adds a small isotropic random

<sup>1</sup> By matching the “units”, we can make the optimization results invariant to simple rescaling  $g(x) = af(bx)$  for scalar  $a, b > 0$ . Here,  $\rho$  has the scaling of third-order derivative,  $\epsilon$  has the scaling of the first-order derivative, so  $\sqrt{\rho\epsilon}$  has the scaling of the second-order derivative.

Gaussian perturbation to the gradient. The perturbation  $\xi_t$  is sampled from a zero-mean Gaussian with covariance  $(r^2/d)\mathbf{I}$  so that  $\mathbb{E} \|\xi_t\|^2 = r^2$ . Parameter  $r$  control the effective radius of the perturbation, which is often chosen to be very small. [JNG<sup>+</sup>19] proves that this simple form of PGD is capable of escaping strict saddle points and finding SOSP efficiently.

In this chapter, to show the insights behind PGD, we turn to an alternative variant of the algorithm, which has a slightly more complicated form, but a easier analysis. The variant we consider here performs the following two steps at each iteration:

1. If  $\|\nabla f(x_t)\| \leq \epsilon$  and no perturbation has been added in the last  $\mathcal{T}$  steps, then add small perturbation  $x_t \leftarrow x_t - \eta \tilde{\xi}_t$  where  $\tilde{\xi}_t \sim \text{Uniform}(\mathbb{B}_0(r))$ .
2.  $x_{t+1} \leftarrow x_t - \eta \nabla f(x_t)$ .

where  $\mathbb{B}_0(r)$  is the Euclidean ball centered at 0 with radius  $r$ . This variant of PGD only adds perturbations when gradient is small and no perturbation has been added in the last  $\mathcal{T}$  steps, thus adds less stochasticity to the algorithm compared to the original form of PGD. In the following theorem, we provide theoretical guarantees for this variant of PGD as follows:

**Theorem 8.2.1.** *Assume  $f$  is  $\ell$ -gradient Lipschitz, and  $\rho$ -Hessian Lipschitz. For any  $\epsilon, \delta > 0$ , if we choose  $\eta = 1/\ell$ ,  $r = \tilde{\Theta}(\epsilon)$ ,  $\mathcal{T} = \tilde{\Theta}(\ell/\sqrt{\rho\epsilon})$ , and run PGD for more than  $\tilde{O}(\ell(f(x_0) - f^*)/\epsilon^2)$  iterations, then with probability at least  $1 - \delta$ , at least one of the iterates will be  $\epsilon$ -SOSP.*

Here  $\tilde{O}(\cdot), \tilde{\Theta}(\cdot)$  hides absolute constant and poly-logarithmic dependence in  $d, \ell, \rho, \epsilon, \delta$  and  $f(x_0) - f^*$ . Our choice of  $\mathcal{T}$  is, up to logarithmic factors, the ratio of the gradient Lipschitz parameter  $\ell$  and the Hessian accuracy tolerance in  $\epsilon$ -SOSP —  $\sqrt{\rho\epsilon}$ .

We remark that, in classic optimization literature, it is known that GD finds an  $\epsilon$ -first-order stationary point (a point  $x$  satisfying  $\|\nabla f(x)\| \leq \epsilon$ ) in  $O(\ell(f(x_0) - f^*)/\epsilon^2)$  iterations [Nes98]. Theorem 8.2.1 shows that PGD finds second-order stationary points in almost the same time as GD finds first-order stationary points, up to only logarithmic factors. In particular, despite there might be only one escaping direction within the  $d$ -dimensional space at saddle points, the dimension dependency of PGD is only polylogarithmic. This is extremely important in high-dimensional settings, such as training deep neural networks. This provides a compelling explanation why strict saddle points are computationally benign for first-order gradient methods.

The overall proof strategy is as follows. According to Definition 8.1.5, if an iterate  $x_t$  is not an  $\epsilon$ -second-order stationary point, then  $x_t$

must either have large gradient or be an approximate saddle point. We prove the following two claims:

1. Large gradient ( $\|\nabla f(x_t)\| > \epsilon$ ), then function value decreases significantly in one step:  $f(x_{t+1}) - f(x_t) \leq -\Omega(\epsilon^2/\ell)$ .
2. Approximate saddle point ( $\|\nabla f(x_t)\| \leq \epsilon$  and  $\lambda_{\min}(\nabla^2 f(x_t)) < -\sqrt{\rho\epsilon}$ ), then, with high probability, function value decreases significantly in  $\mathcal{T}$  steps:  $f(x_{t+\mathcal{T}}) - f(x_t) \leq -\tilde{\Omega}(\mathcal{T} \cdot \epsilon^2/\ell)$ .

That is, in either case, the function value will decrease by  $\tilde{\Omega}(\epsilon^2/\ell)$  on average per step. Since the function value can be no less than the optimal value  $f^*$ , we know after  $\tilde{O}(\ell(f(x_0) - f^*)/\epsilon^2)$  steps, at least one of the iterates must to  $\epsilon$ -second-order stationary point. The first claim immediately follows from the descent lemma (Lemma 2.1.4). In the next section, we will show how to prove the second claim.

### 8.3 Saddle Points Escaping Lemma

In this section, we formally prove that if the starting point has a strictly negative eigenvalue of the Hessian, then adding a perturbation and following by gradient descent will yield a significant decrease in function value in  $\mathcal{T}$  iterations.

**Lemma 8.3.1** (Saddle Points Escaping Lemma). *Under the setting of Theorem 8.2.1, if  $\tilde{x}$  satisfies  $\|\nabla f(\tilde{x})\| \leq \epsilon$ , and  $\lambda_{\min}(\nabla^2 f(\tilde{x})) \leq -\sqrt{\rho\epsilon}$ , then let  $x_0 = \tilde{x} + \eta\xi$  ( $\xi \sim \text{Uniform}(B_0(r))$ ), and run gradient descent starting from  $x_0$ . With probability at least  $1 - \delta$ , we have*

$$f(x_{\mathcal{T}}) - f(\tilde{x}) \leq -\tilde{\Omega}(\mathcal{T} \cdot \epsilon^2/\ell)$$

where  $x_{\mathcal{T}}$  is the  $\mathcal{T}^{\text{th}}$  gradient descent iterate starting from  $x_0$ .

Recall that  $\mathcal{T} = \tilde{\Theta}(\ell/\sqrt{\rho\epsilon})$ . This implies that both saddle point escaping time, and the amount of function decrease depend on dimension  $d$  only polylogarithmically. To prove this lemma, we will show the followings:

- (*Improve or Localize*) If gradient descent keeps making little progress for a certain number of iterations, then all the iterates within those iterations must be stuck in a small Euclidean ball.
- (*Stuck probability is small around saddle points*) If the Hessian has a significant negative eigenvalue, then after a random perturbation, with high probability, gradient descent will not be stuck in a small Euclidean ball for a long time.

Combining two statements above, we conclude that GD in the second statement must make significant progress after a certain number of iterations, which proves Lemma 8.3.1.



### 8.3.1 Improve or Localize

We first prove the following lemma which says that if the function value does not decrease too much over  $t$  iterations, then all iterates  $\{x_\tau\}_{\tau=0}^t$  will remain in a small neighborhood of  $x_0$ .

**Lemma 8.3.2** (Improve or Localize). *Assume function  $f$  is  $\ell$ -gradient Lipschitz, and run GD with  $\eta \leq 1/\ell$ , then for any  $t \geq \tau > 0$ , we have:*

$$\|x_\tau - x_0\| \leq \sqrt{2\eta t(f(x_0) - f(x_t))}.$$

*Proof.* Given the gradient update,  $x_{t+1} = x_t - \eta \nabla f(x_t)$ , we have that for any  $\tau \leq t$ :

$$\begin{aligned} \|x_\tau - x_0\| &\leq \sum_{\tau=1}^t \|x_\tau - x_{\tau-1}\| \stackrel{(1)}{\leq} \left[ t \sum_{\tau=1}^t \|x_\tau - x_{\tau-1}\|^2 \right]^{\frac{1}{2}} \\ &= \left[ \eta^2 t \sum_{\tau=1}^t \|\nabla f(x_{\tau-1})\|^2 \right]^{\frac{1}{2}} \stackrel{(2)}{\leq} \sqrt{2\eta t(f(x_0) - f(x_t))}, \end{aligned}$$

where step (1) uses Cauchy-Schwarz inequality, and step (2) is due to the descent lemma (Lemma 2.1.4).  $\square$

Lemma 8.3.2 immediately implies that if  $f(x_\mathcal{T}) - f(x_0) \geq -\tilde{O}(\mathcal{T} \cdot \epsilon^2/\ell)$ , i.e. GD does not make enough progress in  $\mathcal{T}$  steps after perturbation, then we immediately have that  $\|x_t - x_0\| \leq \tilde{O}(\epsilon\mathcal{T}/\ell)$  for all  $t \in [\mathcal{T}]$ .

### 8.3.2 Bounding the Width of the Stuck Region

Second, we show that the probability for GD sequence to get stuck is small if initialized with a point around saddle point with random perturbation. Recall in Lemma 8.3.1 that  $x_0 \sim \text{Uniform}(\mathbb{B}_{\tilde{x}}(\eta r))$ . We refer to  $\mathbb{B}_{\tilde{x}}(\eta r)$  as the *perturbation ball*, and define the *stuck region* within the perturbation ball to be the set of points starting from which GD makes little progress in  $\mathcal{T}$  steps:

$$\begin{aligned} \mathcal{X}_{\text{stuck}} &:= \{x \in \mathbb{B}_{\tilde{x}}(\eta r) \mid \{x_t\}_{t=0}^{\mathcal{T}} \text{ is a GD sequence with} \\ &\quad x_0 = x, \text{ and } \forall t \in [\mathcal{T}], \|x_t - x_0\| \leq \tilde{O}(\epsilon\mathcal{T}/\ell)\}. \end{aligned}$$

See Figure 8.1 for illustrations. Since  $x_0$  sampled uniformly from this perturbation ball, the probability GD got stuck after perturbation is equal to the ratio of the volume of the stuck region and the volume of the perturbation ball. Therefore, we want to show that the stuck region has small volume.

In general, the shape of the stuck region can be very complicated, so it is very difficult to directly compute its volume. A crucial observation here is that, despite we do not know the shape of the stuck

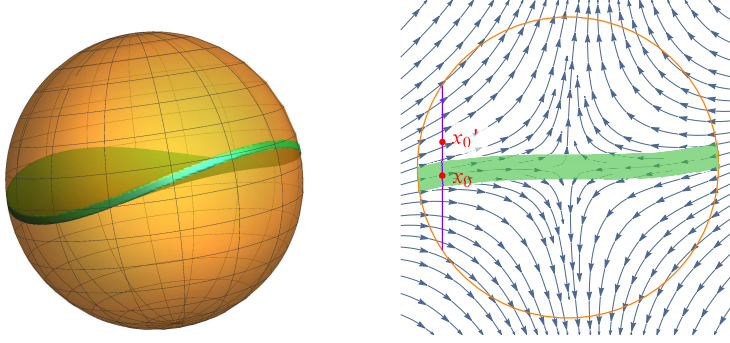


Figure 8.1: **Left:** Perturbation ball in 3D and “thin pancake” shape stuck region. **Right:** Perturbation ball in 2D and “narrow band” stuck region under gradient flow

region, we can prove the width of  $\mathcal{X}_{\text{stuck}}$  along the minimum eigenvalue direction of  $\nabla^2 f(\tilde{x})$  is small. In fact, if the width is at most  $\eta\omega$ , then we have  $\text{Vol}(\mathcal{X}_{\text{stuck}}) \leq \text{Vol}(\mathbb{B}_0^{d-1}(\eta r))\eta\omega$ , and thus,

$$\begin{aligned} \Pr(x_0 \in \mathcal{X}_{\text{stuck}}) &= \frac{\text{Vol}(\mathcal{X}_{\text{stuck}})}{\text{Vol}(\mathbb{B}_{\tilde{x}}^d(\eta r))} \leq \frac{\eta\omega \times \text{Vol}(\mathbb{B}_0^{d-1}(\eta r))}{\text{Vol}(\mathbb{B}_0^d(\eta r))} \\ &= \frac{\omega}{r\sqrt{\pi}} \frac{\Gamma(\frac{d}{2} + 1)}{\Gamma(\frac{d}{2} + \frac{1}{2})} \leq \frac{\omega}{r} \cdot \sqrt{\frac{d}{\pi}} \end{aligned}$$

To achieve failure probability at most  $\delta$ , we hope  $\omega \leq O(\delta r / \sqrt{d})$ . We bound the width of the stuck region  $\mathcal{X}_{\text{stuck}}$  by the novel techniques of **coupling sequences**—consider two GD sequences  $\{x_t\}_{t=0}^{\mathcal{T}}, \{x'_t\}_{t=0}^{\mathcal{T}}$  which satisfy: (1)  $\max\{\|x_0 - \tilde{x}\|, \|x'_0 - \tilde{x}\|\} \leq \eta r$ ; and (2)  $x_0 - x'_0 = \eta\omega e_1$ , where  $e_1$  is the minimum eigenvector of  $\nabla^2 f(\tilde{x})$ , and  $\omega \geq \omega_0$  for some threshold  $\omega_0$ .

**Lemma 8.3.3.** *For any  $\omega_0 \in (0, \epsilon]$ , under the setting of Lemma 8.3.1, if  $\{x_t\}_{t=0}^{\mathcal{T}}, \{x'_t\}_{t=0}^{\mathcal{T}}$  are coupling sequences as specified above, then for  $\mathcal{T} \geq \Omega(\kappa \cdot \log(\epsilon\kappa/\omega_0))$  where  $\kappa := \ell/\sqrt{\rho\epsilon}$ , we have,*

$$\exists t \in [\mathcal{T}], \quad \max\{\|x_t - x_0\|, \|x'_t - x'_0\|\} \geq \tilde{\Omega}(\epsilon\mathcal{T}/\ell)$$

Lemma 8.3.3 claims that for any pair of  $x_0, x'_0$  whose difference is on the  $e_1$  direction, with length greater or equal to  $\eta\omega_0$ , at least one of  $x_0, x'_0$  is outside  $\mathcal{X}_{\text{stuck}}$ . This directly implies the width of  $\mathcal{X}_{\text{stuck}}$  in  $e_1$  direction is  $\eta\omega_0$ . Lemma 8.3.3 further claims that the width  $\eta\omega_0$  can be made arbitrarily small by paying only logarithmic factors in the choice of  $\mathcal{T}$ .

*Proof.* We prove by contradiction. Assume the contrary of Lemma 8.3.3 is true— $\max\{\|x_t - x_0\|, \|x'_t - x'_0\|\} \leq \tilde{O}(\epsilon\mathcal{T}/\ell)$  for all  $t \in [\mathcal{T}]$ , i.e. both GD sequences stuck in a small Euclidean ball for  $\mathcal{T}$  steps.

We can write out the update equation for the difference of the

couple sequences  $\widehat{x}_t := x_t - x'_t$  as:

$$\begin{aligned}\widehat{x}_{t+1} &= \widehat{x}_t - \eta[\nabla f(x_t) - \nabla f(x'_t)] = (\mathbf{I} - \eta\mathcal{H})\widehat{x}_t - \eta\Delta_t\widehat{x}_t \\ &= \underbrace{(\mathbf{I} - \eta\mathcal{H})^{t+1}\widehat{x}_0}_{\mathfrak{p}(t+1)} - \underbrace{\eta \sum_{\tau=0}^t (\mathbf{I} - \eta\mathcal{H})^{t-\tau} \Delta_\tau \widehat{x}_\tau}_{\mathfrak{q}(t+1)},\end{aligned}$$

where  $\mathcal{H} = \nabla^2 f(\widetilde{x})$  and  $\Delta_t = \int_0^1 [\nabla^2 f(x'_t + \theta(x_t - x'_t)) - \mathcal{H}] d\theta$ . We note that term  $\mathfrak{p}(t)$  is the formula of  $\widehat{x}_t$  if function  $f$  is quadratic around  $\widetilde{x}$ , and  $\mathfrak{q}(t)$  is the approximation error term caused by function  $f$  being non-quadratic.

We will show later that the quadratic term is the dominating term, in the sense that  $\|\mathfrak{q}(t)\| \leq \|\mathfrak{p}(t)\|/2$  for all  $t \in [\mathcal{T}]$ . Given this is true, since  $\widehat{x}_0 = \eta\omega e_1$ , we have

$$\|\mathfrak{p}(t)\| \geq (1 + \sqrt{\epsilon\rho}/\ell)^t \cdot \eta\omega_0$$

This term grows exponentially. Therefore, by choosing  $\mathcal{T} \geq \Omega(\kappa \cdot \log(\epsilon\kappa/\omega_0))$ , we have  $\|\widehat{x}_t\| \geq \|\mathfrak{p}(t)\|/2 \geq \widetilde{\Omega}(\epsilon\mathcal{T}/\ell)$ , which contradicts the assumption that both GD sequences stuck in a small Euclidean ball with radius  $\widetilde{O}(\epsilon\mathcal{T}/\ell)$  for  $\mathcal{T}$  steps (we note  $\|x_0 - x'_0\| \leq 2\eta r \ll \widetilde{O}(\epsilon\mathcal{T}/\ell)$ ). This proves Lemma 8.3.3.

For the remaining of the proof, we only need to verify by induction that  $\|\mathfrak{q}(t)\| \leq \|\mathfrak{p}(t)\|/2$  for all  $t \in [\mathcal{T}]$ . The claim is true for the base case  $t = 0$  as  $\|\mathfrak{q}(0)\| = 0 \leq \|\widehat{x}_0\|/2 = \|\mathfrak{p}(0)\|/2$ . Now suppose the induction claim is true up to  $t$ . Denote  $\lambda_{\min}(\mathcal{H}) = -\gamma$ . Note that  $\widehat{x}_0$  lies in the direction of the minimum eigenvector of  $\mathcal{H}$ . Thus for any  $\tau \leq t$ , we have:

$$\|\widehat{x}_\tau\| \leq \|\mathfrak{p}(\tau)\| + \|\mathfrak{q}(\tau)\| \leq 2\|\mathfrak{p}(\tau)\| = 2(1 + \eta\gamma)^\tau \eta\omega.$$

By the Hessian Lipschitz property, we further have

$$\|\Delta_t\| \leq \rho \max\{\|x_t - \widetilde{x}\|, \|x'_t - \widetilde{x}\|\} \leq \widetilde{O}(\rho\epsilon\mathcal{T}/\ell) = \widetilde{O}(\sqrt{\rho\epsilon})$$

therefore:

$$\begin{aligned}\|\mathfrak{q}(t+1)\| &= \left\| \eta \sum_{\tau=0}^t (\mathbf{I} - \eta\mathcal{H})^{t-\tau} \Delta_\tau \widehat{x}_\tau \right\| \\ &\leq \eta \sum_{\tau=0}^t \|\Delta_\tau\| \|(\mathbf{I} - \eta\mathcal{H})^{t-\tau}\| \|\widehat{x}_\tau\| \leq \widetilde{O}(\eta\sqrt{\rho\epsilon}) \sum_{\tau=0}^t (1 + \eta\gamma)^\tau \eta\omega \\ &\leq \widetilde{O}(1)(1 + \eta\gamma)^t \eta\omega \leq \widetilde{O}(1) \|\mathfrak{p}(t+1)\|,\end{aligned}$$

where the second-to-last inequality uses  $t+1 \leq \mathcal{T}$ , and  $\widetilde{O}(\eta\mathcal{T}\sqrt{\rho\epsilon}) = \widetilde{O}(1)$ . Finally, with a careful treatment of constant and logarithmic factors, we can in fact make this  $\widetilde{O}(1)$  term less or equal to  $1/2$  (we omit the detail here). This finishes the inductive proof.  $\square$



## *Algorithmic Regularization*

Large scale neural networks used in practice are highly over-parameterized with far more trainable model parameters compared to the number of training examples. Consequently, the optimization objectives for learning such high capacity models have many global minima that fit training data perfectly. However, minimizing the training loss using specific optimization algorithms take us to not just any global minima, but some special global minima – in this sense the choice of optimization algorithms introduce a implicit form of inductive bias in learning which can aid generalization.

In over-parameterized models, specially deep neural networks, much, if not most, of the inductive bias of the learned model comes from this implicit regularization from the optimization algorithm. For example, early empirical work on this topic (ref. [NTS15a, NSS15, HS97, KMN<sup>+</sup>16, ZBH<sup>+</sup>16a, CCS<sup>+</sup>16, DPBB17, ADG<sup>+</sup>16, Ney17, WRS<sup>+</sup>17, HHS17, Smi18]) show that deep models often generalize well even when trained purely by minimizing the training error without any explicit regularization, and even when the networks are highly overparameterized to the extent of being able to fit random labels. Consequently, there are many zero training error solutions, all global minima of the training objective, most of which generalize horribly. Nevertheless, our choice of optimization algorithm, typically a variant of gradient descent, seems to prefer solutions that do generalize well. This generalization ability cannot be explained by the capacity of the explicitly specified model class (namely, the functions representable in the chosen architecture). Instead, the optimization algorithm biasing toward a “simple” model, minimizing some implicit “regularization measure”, say  $R(w)$ , is key for generalization. Understanding the implicit inductive bias, *e.g.* via characterizing  $R(w)$ , is thus essential for understanding how and what the model learns. For example, in linear regression it can be shown that minimizing an under-determined model (with more parameters than samples) using gradient descent yields the minimum  $\ell_2$  norm solution (see

Proposition 9.1.1), and for linear logistic regression trained on linearly separable data, gradient descent converges in the direction of the hard margin support vector machine solution (Theorem 9.3.2), even though the norm or margin is not explicitly specified in the optimization problem. In fact, such analysis showing implicit inductive bias from optimization algorithm leading to generalization is not new. In the context of boosting algorithms, (author?) [EHJT04] and (author?) [Tel13] established connections of gradient boosting algorithm (coordinate descent) to  $\ell_1$  norm minimization, and  $\ell_1$  margin maximization, respectively. minimization was observed. Such minimum norm or maximum margin solutions are of course very special among all solutions or separators that fit the training data, and in particular can ensure generalization [BM03, KST09].

In this chapter, we largely present results on algorithmic regularization of vanilla gradient descent when minimizing unregularized training loss in regression and classification problem over various simple and complex model classes. We briefly discuss general algorithmic families like steepest descent and mirror descent.

*Meanings of “implicit regularization due to training algorithm.”*

Results in the current chapter tend to show that the solution obtained by applying training algorithm  $A$  on *Objective 1* essentially to convergence (e.g. to stationary point of gradient descent), also satisfies KKT local optimality conditions for some other *Objective 2*. In many results *Objective 2* is simply *Objective 1* with a regularizer term, typically involving some norm of the solution. Hence we can think of the training algorithm as *implicitly regularizing* the objective.

While these results give good insight into the effect of the training a few caveats are in order, especially if we seek takeaways for deep learning. First, even though the solution found happens to be a KKT point of *Objective 2*, it may be never (or almost never) observed if we actually do standard training on *Objective 2*.<sup>1</sup> Second, the results in this chapter are often stated for training carried out to infinite time, which may also limit their applicability to real life.

In later chapters we will see a different type of analysis, which analyses the trajectory followed by the solution as it evolves during training. This *dynamic* view of training quickly gets complicated (as opposed to the more static view taken in understanding stationary points) and has not been achieved for realistic deep nets yet.

## 9.1 Linear models in regression: squared loss

SURIYA: PLS SEE CHAPTER 3 AND MODIFY THE WRITEUP AS NEEDED.

We first demonstrate the algorithmic regularization in a simple

<sup>1</sup> Recall that in a nonconvex landscape the solution obtained at the end of training is greatly affected by the initialization, and in deep learning the initialization is very special.

linear regression setting where the prediction function is specified by a linear function of inputs:  $f_w(x) = w^\top x$  and we have the following empirical risk minimization objective.

$$L(w) = \sum_{i=1}^n \left( w^\top x^{(i)} - y^{(i)} \right)^2. \quad (9.1)$$

Such simple models are natural starting points to build analytical tools for extending to complex models, and such results provide intuitions for understanding and improving upon the empirical practices in neural networks. Although the results in this section are specified for squared loss, the results and proof technique extend for any smooth loss a unique finite root: where  $\ell(\hat{y}, y)$  between a prediction  $\hat{y}$  and label  $y$  is minimized at a unique and finite value of  $\hat{y}$  [GLSS18a].

We are particularly interested in the case where  $n < d$  and the observations are realizable, i.e.,  $\min_w L(w) = 0$ . Under these conditions, the optimization problem in eq. (9.1) is underdetermined and has multiple global minima denoted by  $\mathcal{G} = \{w : \forall i, w^\top x^{(i)} = y^{(i)}\}$ . In this and all the following problems we consider, the goal is to answer: *Which specific global minima do different optimization algorithms reach when minimizing  $L(w)$ ?*

The following proposition is the simplest illustration of the algorithmic regularization phenomenon.

**Proposition 9.1.1.** *Consider gradient descent updates  $w_t$  for the loss in eq. (9.1) starting with initialization  $w_0$ . For any step size schedule that minimizes the loss  $L(w)$ , the algorithm returns a special global minimizer that implicitly also minimizes the Euclidean distance to the initialization:  $w_t \rightarrow \operatorname{argmin}_{w \in \mathcal{G}} \|w - w_0\|_2$ .*

*Proof.* The key idea is in noting that that the gradients of the loss function have a special structure. For the linear regression loss in eq. (9.1)  $\forall w, \nabla L(w) = \sum_i (w^\top x^{(i)} - y^{(i)}) x^{(i)} \in \operatorname{span}(\{x^{(i)}\})$  - that is the gradients are restricted to a  $n$  dimensional subspace that is independent of  $w$ . Thus, the gradient descent updates from initialization  $w_t - w_0 = \sum_{t' < t} \eta w_{t'}$ , which linearly accumulate gradients, are again constrained to the  $n$  dimensional subspace. It is now easy to check that there is a unique global minimizer that both fits the data ( $w \in \mathcal{G}$ ) as well as is reachable by gradient descent ( $w \in w_0 + \operatorname{span}(\{x^{(i)}\})$ ). By checking the KKT conditions, it can be verified that this unique minimizer is given by  $\operatorname{argmin}_{w \in \mathcal{G}} \|w - w_0\|_2^2$ .  $\square$

In general overparameterized optimization problems, the characterization of the implicit bias or algorithmic regularization is often not this elegant or easy. For the same model class, changing the algorithm, or changing associated hyperparameter (like step size

and initialization), or even changing the specific parameterization of the model class can change the implicit bias. For example, (author?) [WRS<sup>+</sup>17] showed that for some standard deep learning architectures, variants of SGD algorithm with different choices of momentum and adaptive gradient updates (AdaGrad and Adam) exhibit different biases and thus have different generalization performance; (author?) [KMN<sup>+</sup>16], (author?) [HHS17] and (author?) [Smi18] study how the size of the mini-batches used in SGD influences generalization; and (author?) [NSS15] compare the bias of path-SGD (steepest descent with respect to a scale invariant path-norm) to standard SGD.

A comprehensive understanding of how all the algorithmic choices affect the implicit bias is beyond the scope of this chapter (and also the current state of research). However, in the context of this chapter, we specifically want to highlight the role of *geometry* induced by optimization algorithm and specific parameterization, which are discussed briefly below.

### 9.1.1 Geometry induced by updates of local search algorithms

The relation of gradient descent to implicit bias towards minimizing Euclidean distance to initialization is suggestive of the connection between algorithmic regularization to the geometry of updates in local search methods. In particular, gradient descent iterations can be alternatively specified by the following equation where the  $t + 1$ th iterate is derived by minimizing the a local (first order Taylor) approximation of the loss while constraining the step length in Euclidean norm.

$$w_{t+1} = \operatorname{argmin}_w \langle w, \nabla L(w_t) \rangle + \frac{1}{2\eta} \|w - w_t\|_2^2. \quad (9.2)$$

Motivated by the above connection, we can study other families of algorithms that work under different and non-Euclidean geometries. Two convenient families are mirror descent w.r.t. potential  $\psi$  [BT03, NY83] and steepest descent w.r.t. general norms [BV04].

**Mirror descent w.r.t. potential  $\psi$**  Mirror descent updates are defined for any strongly convex and differentiable potential  $\psi$  as

$$\begin{aligned} w_{t+1} &= \operatorname{argmin}_w \eta \langle w, \nabla L(w_t) \rangle + D_\psi(w, w_t), \\ \implies \nabla \psi(w_{t+1}) &= \nabla \psi(w_t) - \eta \nabla L(w_t) \end{aligned} \quad (9.3)$$

where  $D_\psi(w, w') = \psi(w) - \psi(w') - \langle \nabla \psi(w'), w - w' \rangle$  is the *Bregman divergence* [Bre67] w.r.t.  $\psi$ . This family captures updates where the geometry is specified by the Bregman divergence  $D_\psi$ . Examples of



potentials  $\psi$  for mirror descent include the squared  $\ell_2$  norm  $\psi(w) = 1/2\|w\|_2^2$ , which leads to gradient descent; the entropy potential  $\psi(w) = \sum_i w[i] \log w[i] - w[i]$ ; the spectral entropy for matrix valued  $w$ , where  $\psi(w)$  is the entropy potential on the singular values of  $w$ ; general quadratic potentials  $\psi(w) = 1/2\|w\|_D^2 = 1/2 w^\top D w$  for any positive definite matrix  $D$ ; and the squared  $\ell_p$  norms for  $p \in (1, 2]$ .

From eq. (9.3), we see that rather than  $w_t$  (called primal iterates), it is the  $\nabla\psi(w_t)$  (called dual iterates) that are constrained to the low dimensional data manifold  $\nabla\psi(w_0) + \text{span}(\{x^{(i)}\})$ . The arguments for gradient descent can now be generalized to get the following result.

**Theorem 9.1.2.** *For any realizable dataset  $\{x^{(i)}, y^{(i)}\}_{n=1}^N$ , and any strongly convex potential  $\psi$ , consider the mirror descent iterates  $w_t$  from eq. (9.3) for minimizing the empirical loss  $L(w)$  in eq. (9.1). For all initializations  $w_0$ , if the step-size schedule minimizes  $L(w)$ , i.e.,  $L(w_t) \rightarrow 0$ , then the asymptotic solution of the algorithm is given by*

$$w_t \rightarrow \arg \min_{w: \forall i, w^\top x^{(i)} = y^{(i)}} D_\psi(w, w_0). \tag{9.4}$$

In particular, if we start at  $w_0 = \arg \min_w \psi(w)$  (so that  $\nabla\psi(w_0) = 0$ ), then we get to  $\arg \min_{w \in \mathcal{G}} \psi(w)$ .<sup>2</sup>

<sup>2</sup> The analysis of Theorem 9.1.2 and Proposition 9.1.1 also hold when instancewise stochastic gradients are used in place of  $\nabla L(w_t)$ .

**Steepest descent w.r.t. general norms** Gradient descent is also a special case of steepest descent (SD) w.r.t a generic norm  $\|\cdot\|$  [BV04] with updates given by,

$$w_{t+1} = w_t + \eta_t \Delta w_t, \text{ where } \Delta w_t = \arg \min_v \langle \nabla L(w_t), v \rangle + \frac{1}{2} \|v\|^2. \tag{9.5}$$

Examples of steepest descent include gradient descent, which is steepest descent w.r.t  $\ell_2$  norm and coordinate descent, which is steepest descent w.r.t  $\ell_1$  norm. In general, the update  $\Delta w_t$  in eq. (9.5) is not uniquely defined and there could be multiple direction  $\Delta w_t$  that minimize eq. (9.5). In such cases, any minimizer of eq. (9.5) is a valid steepest descent update.

Generalizing gradient descent and mirror descent, we might expect the steepest descent iterates to converge to the solution closest to initialization in corresponding norm,  $\arg \min_{w \in \mathcal{G}} \|w - w_0\|$ . This is indeed the case for quadratic norms  $\|v\|_D = \sqrt{v^\top D v}$  when eq. 9.5 is equivalent to mirror descent with  $\psi(w) = 1/2\|w\|_D^2$ . Unfortunately, this does not hold for general norms as shown by the following results.

**Example 1.** In the case of coordinate descent, which is a special case of steepest descent w.r.t. the  $\ell_1$  norm, (author?) [EHJT04] studied this phenomenon in the context of gradient boosting: observing that sometimes but *not always* the optimization path of coordinate descent

given by  $\Delta w_{t+1} \in \text{conv} \left\{ -\eta_t \frac{\partial L(w_t)}{\partial w_{[j_t]}} e_{j_t} : j_t = \text{argmax}_j \left| \frac{\partial L(w_t)}{\partial w_{[j]}} \right| \right\}$ , coincides with the  $\ell_1$  regularization path given by,  $\hat{w}(\lambda) = \text{arg min}_w L(w) + \lambda \|w\|_1$ . The specific coordinate descent path where updates average all the optimal coordinates and the step-sizes are infinitesimal is equivalent to forward stage-wise selection, a.k.a.  $\epsilon$ -boosting [Frio1]. When the  $\ell_1$  regularization path  $\hat{w}(\lambda)$  is monotone in each of the coordinates, it is identical to this stage-wise selection path, i.e., to a coordinate descent optimization path (and also to the related LARS path) [EHJTo4]. In this case, at the limit of  $\lambda \rightarrow 0$  and  $t \rightarrow \infty$ , the optimization and regularization paths, both converge to the minimum  $\ell_1$  norm solution. However, when the regularization path  $\hat{w}(\lambda)$  is not monotone, which can and does happen, the optimization and regularization paths diverge, and forward stage-wise selection can converge to solutions with sub-optimal  $\ell_1$  norm.

Example 2. The following example shows that even for  $\ell_p$  norms where the  $\|\cdot\|_p^2$  is smooth and strongly convex, the global minimum returned by the steepest descent depends on the step-size. Consider minimizing  $L(w)$  with dataset  $\{(x^{(1)} = [1, 1, 1], y^{(1)} = 1), (x^{(2)} = [1, 2, 0], y^{(2)} = 10)\}$  using steepest descent updates w.r.t. the  $\ell_{4/3}$  norm. The empirical results for this problem in Figure 9.1 clearly show that steepest descent converges to a global minimum that depends on the step-size and even in the continuous step-size limit of  $\eta \rightarrow 0$ ,  $w_t$  does not converge to the expected solution of  $\text{arg min}_{w \in G} \|w - w_0\|$ .

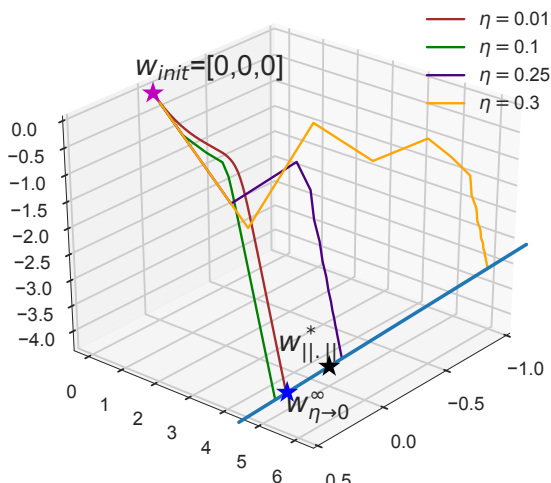


Figure 9.1: Steepest descent w.r.t  $\|\cdot\|_{4/3}$ : the global minimum to which steepest descent converges to depends on  $\eta$ . Here  $w_0 = [0, 0, 0]$ ,  $w_{\|\cdot\|}^* = \text{arg min}_{\psi \in G} \|\psi\|_{4/3}$  denotes the minimum norm global minimum, and  $w_{\eta \rightarrow 0}^\infty$  denotes the solution of infinitesimal SD with  $\eta \rightarrow 0$ . Note that even as  $\eta \rightarrow 0$ , the expected characterization does not hold, i.e.,  $w_{\eta \rightarrow 0}^\infty \neq w_{\|\cdot\|}^*$ .

In summary, for squared loss, we characterized the implicit bias of generic mirror descent algorithm in terms of the potential function and initialization. However, even in simple linear regression, for steepest descent with general norms, we were unable to get a useful

characterization. In contrast, in Section 9.3.2, we study logistic like strictly monotonic losses used in classification, where we *can* get a characterization for steepest descent.

### 9.1.2 Geometry induced by parameterization of model class

In many learning problems, the same model class can be parameterized in multiple ways. For example, the set of linear functions in  $\mathbb{R}^d$  can be parameterized in a canonical way as  $w \in \mathbb{R}^d$  with  $f_w(x) = w^\top x$ , but also equivalently by  $u, v \in \mathbb{R}^d$  with  $f_{u,v}(x) = (u \cdot v)^\top x$  or  $f_{u,v}(x) = (u^2 - v^2)^\top x$ . All such equivalent parameterizations lead to equivalent training objectives, however, in overparameterized models, using gradient descent on different parameterizations lead to different induced biases in the function space. For example, (author?) [GWB<sup>+</sup>17, GLSS18b] demonstrated this phenomenon in matrix factorization and linear convolutional networks, where these parameterizations were shown to introduce interesting and unusual biases towards minimizing nuclear norm, and  $\ell_p$  (for  $p = 2/\text{depth}$ ) norm in Fourier domain, respectively. In general, these results are suggestive of role of architecture choice in different neural network models, and shows how even while using the same gradient descent algorithm, different geometries in the function space can be induced by the different parameterizations.

## 9.2 Matrix factorization

«Suriya notes: I would like to include this section here but can also move to a separate chapter. Ideally, summarize our 2017 paper, Tengyu's 2018 paper and Nadav's 2019 paper. May be we can discuss this after Nadav's lecture?»

## 9.3 Linear Models in Classification

We now turn to studying classification problems with logistic or cross-entropy type losses. We focus on binary classification problems where  $y^{(i)} \in \{-1, 1\}$ . Many continuous surrogate of the 0-1 loss including logistic, cross-entropy, and exponential loss are examples of strictly monotone loss functions  $\ell$  where the behavior of the implicit bias is fundamentally different, and as are the situations when the implicit bias can be characterized.

We look at classification models that fit the training data  $\{x^{(i)}, y^{(i)}\}_i$  with linear decision boundaries  $f(x) = w^\top x$  with decision rule given by  $\hat{y}(x) = \text{sign}(f(x))$ . In many instances of the proofs, we also assume without loss of generality that  $y^{(i)} = 1$  for all  $i$ , since for linear models, the sign of  $y^{(i)}$  can equivalently be absorbed into  $x^{(i)}$ . We

again look at unregularized empirical risk minimization objective of the form in eq. (9.1), but now with strictly monotone losses. When the training data  $\{x^{(i)}, y^{(i)}\}_n$  is not linearly separable, the empirical objective  $L(w)$  can have a finite global minimum. However, if the dataset is linearly separable, i.e.,  $\exists w : \forall i, y^{(i)} w^\top x^{(i)} > 0$ , the empirical loss  $L(w)$  is again ill-posed, and moreover  $L(w)$  does not have any finite minimizer, i.e.,  $L(w) \rightarrow 0$  only as  $\|w\| \rightarrow \infty$ . Thus, for any sequence  $\{w_t\}_{t=0}^\infty$ , if  $L(w_t) \rightarrow 0$ , then  $w_t$  necessarily diverges to infinity rather than converge, and hence we cannot talk about  $\lim_{t \rightarrow \infty} w_t$ . Instead, we look at the limit direction  $\bar{w}_\infty = \lim_{t \rightarrow \infty} \frac{w_t}{\|w_t\|}$  whenever the limit exists. We refer to existence of this limit as convergence in direction. Note that, the limit direction fully specifies the decision rule of the classifier that we care about.

In the remainder of the chapter, we focus on the following exponential loss  $\ell(u, y) = \exp(-uy)$ . However, our asymptotic results can be extended to loss functions with tight exponential tails, including logistic and sigmoid losses, along the lines of (author?) [SHS17] and (author?) [Tel13].

$$L(w) = \sum_{i=1}^n \exp(-y^{(i)} w^\top x^{(i)}). \quad (9.6)$$

### 9.3.1 Gradient Descent

(author?) [SHS17] showed that for almost all linearly separable datasets, gradient descent with *any initialization and any bounded step-size* converges in direction to maximum margin separator with unit  $\ell_2$  norm, i.e., the hard margin support vector machine classifier.

This characterization of the implicit bias is independent of both the step-size as well as the initialization. We already see a fundamentally difference from the implicit bias of gradient descent for losses with a unique finite root (Section ??) where the characterization depended on the initialization. The above result is rigorously proved as part of a more general result in Theorem 9.3.2. Below is a simpler statement and with a heuristic proof sketch intended to convey the intuition for such results.

**Theorem 9.3.1.** *For almost all dataset which is linearly separable, consider gradient descent updates with any initialization  $w_0$  and any step size that minimizes the exponential loss in eq. (9.6), i.e.,  $L(w_t) \rightarrow 0$ . The gradient descent iterates then converge in direction to the  $\ell_2$  max-margin vector, i.e.,  $\lim_{t \rightarrow \infty} \frac{w_t}{\|w_t\|_2} = \frac{\hat{w}}{\|\hat{w}\|}$ , where*

$$\hat{w} = \underset{w}{\operatorname{argmin}} \|w\|^2 \text{ s.t. } \forall i, w^\top x^{(i)} y^{(i)} \geq 1. \quad (9.7)$$

Without loss of generality assume that  $\forall i, y^{(i)} = 1$  as the sign for linear models can be absorbed into  $x^{(i)}$ .

*Proof Sketch* We first understand intuitively why an exponential tail of the loss entail asymptotic convergence to the max margin vector: examine the asymptotic regime of gradient descent in when the exponential loss is minimized, as we argued earlier, this required that  $\forall i : w^\top x^{(i)} \rightarrow \infty$ . Suppose  $w_t / \|w_t\|_2$  converges to some limit  $w_\infty$ , so we can write  $w_t = g(t)w_\infty + \rho(t)$  such that  $g(t) \rightarrow \infty, \forall i, w_\infty^\top x^{(i)} > 0$ , and  $\lim_{t \rightarrow \infty} \rho(t)/g(t) = 0$ . The gradients at  $w_t$  are given by:

$$\begin{aligned} -\nabla \mathcal{L}(w) &= \sum_{i=1}^n \exp(-w^\top x^{(i)}) x^{(i)} \\ &= \sum_{i=1}^n \exp(-g(t)w_\infty^\top x^{(i)}) \exp(-\rho(t)^\top x^{(i)}) x_n. \end{aligned} \quad (9.8)$$

As  $g(t) \rightarrow \infty$  and the exponents become more negative, only those samples with the largest (*i.e.*, least negative) exponents will contribute to the gradient. These are precisely the samples with the smallest margin  $\arg\min_i w_\infty^\top x^{(i)}$ , aka the “support vectors”. The accumulation of negative gradient, and hence  $w_t$ , would then asymptotically be dominated by a non-negative linear combination of support vectors. These are precisely the KKT conditions for the SVM problem (eq. 9.7). Making these intuitions rigorous constitutes the bulk of the proof in (author?) [SHS17], which uses a proof technique very different from that in the following section (Section 9.3.2).

### 9.3.2 Steepest Descent

. Recall that gradient descent is a special case of steepest descent (SD) w.r.t a generic norm  $\|\cdot\|$  with updates given by eq. (9.5). The optimality condition of  $\Delta w_t$  in eq. (9.5) requires

$$\langle \Delta w_t, -\nabla L(w_t) \rangle = \|\Delta w_t\|^2 = \|\nabla L(w_t)\|_{\star}^2, \quad (9.9)$$

where  $\|x\|_{\star} = \sup_{\|y\| \leq 1} x^\top y$  is the dual norm of  $\|\cdot\|$ . Examples of steepest descent include gradient descent, which is steepest descent w.r.t  $\ell_2$  norm and greedy coordinate descent (Gauss-Southwell selection rule), which is steepest descent w.r.t  $\ell_1$  norm. In general, the update  $\Delta w_t$  in eq. (9.5) is not uniquely defined and there could be multiple direction  $\Delta w_t$  that minimize eq. (9.5). In such cases, any minimizer of eq. (9.5) is a valid steepest descent update and satisfies eq. (9.9).

In the preliminary result in Theorem 9.3.1, we proved the limit direction of gradient flow on the exponential loss is the  $\ell_2$  max-margin solution. In the following theorem, we show the natural extension of this to all steepest descent algorithms.

**Theorem 9.3.2.** For any separable dataset  $\{x_i, y_i\}_{i=1}^n$  and any norm  $\|\cdot\|$ , consider the steepest descent updates from eq. (9.9) for minimizing  $L(w)$  in eq. (9.6) with the exponential loss  $\ell(u, y) = \exp(-uy)$ . For all initializations  $w_0$ , and all bounded step-sizes satisfying  $\eta_t \leq \min\{\eta_+, \frac{1}{B^2 L(w_t)}\}$ , where  $B := \max_n \|x_n\|_*$  and  $\eta_+ < \infty$  is any finite number, the iterates  $w_t$  satisfy the following,

$$\lim_{t \rightarrow \infty} \min_n \frac{y_i \langle w_t, y_i \rangle}{\|w_t\|} = \max_{w: \|w\| \leq 1} \min_n y_i \langle w, x_i \rangle =: \gamma.$$

In particular, if there is a unique maximum- $\|\cdot\|$  margin solution  $w^* = \arg \max_{\|w\| \leq 1} \min_i y_i \langle w, x_i \rangle$ , then the limit direction satisfies  $\lim_{t \rightarrow \infty} \frac{w_t}{\|w_t\|} = w^*$ .

A special case of Theorem 9.3.2 is for steepest descent w.r.t. the  $\ell_1$  norm, which as we already saw corresponds to greedy coordinate descent. More specifically, coordinate descent on the exponential loss with exact line search is equivalent to AdaBoost [SF12], where each coordinate represents the output of one “weak learner”. Indeed, initially mysterious generalization properties of boosting have been understood in terms of implicit  $\ell_1$  regularization [SF12], and later on AdaBoost with small enough step-size was shown to converge in direction precisely to the maximum  $\ell_1$  margin solution [ZY<sup>+</sup>05, SSS10, Tel13], just as guaranteed by Theorem 9.3.2. In fact, (author?) [Tel13] generalized the result to a richer variety of exponential tailed loss functions including logistic loss, and a broad class of non-constant step-size rules. Interestingly, coordinate descent with exact line search (AdaBoost) can result in infinite step-sizes, leading the iterates to converge in a different direction that is not a max- $\ell_1$ -margin direction [RDS04], hence the bounded step-sizes rule in Theorem 9.3.2.

Theorem 9.3.2 is a generalization of the result of (author?) to steepest descent with respect to other norms, and our proof follows the same strategy as (author?). We first prove a generalization of the duality result of (author?) [SSS10]: if there is a unit norm linear separator that achieves margin  $\gamma$ , then  $\|\nabla L(w)\|_* \geq \gamma L(w)$  for all  $w$ . By using this lower bound on the dual norm of the gradient, we are able to show that the loss decreases faster than the increase in the norm of the iterates, establishing convergence in a margin maximizing direction.

In the rest of this section, we discuss the proof of Theorem 9.3.2. The proof is divided into three steps:

1. Gradient domination condition: For all norms and any  $w$ ,  $\|\nabla L(w)\|_* \geq \gamma L(w)$
2. Optimization properties of steepest descent such as decrease of

loss function and convergence of the gradient in dual norm to zero.

3. Establishing sufficiently fast convergence of  $L(w_t)$  relative to the growth of  $\|w_t\|$  to prove the Theorem.

**Proposition 9.3.3.** *Gradient domination condition (Lemma 10 of [GLSS18a])*

Let  $\gamma = \max_{\|w\| \leq 1} \min_i y_i x_i^\top w$ . For all  $w$ ,

$$\|\nabla L(w)\|_* \geq \gamma L(w).$$

Next, we establish some optimization properties of the steepest descent algorithm including convergence of gradient norms and loss value.

**Proposition 9.3.4.** *(Lemma 11 and 12 of (author?) [GLSS18a]) Consider*

*the steepest descent iterates  $w_t$  on (9.6) with stepsize  $\eta \leq \frac{1}{B^2 L(w_0)}$ , where  $B = \max_i \|x_i\|_*$ . The following holds:*

1.  $L(w_{t+1}) \leq L(w_t)$ .
2.  $\sum_{t=0}^{\infty} \|\nabla L(w_t)\|^2 < \infty$  and hence  $\|\nabla L(w_t)\|_* \rightarrow 0$ .
3.  $L(w_t) \rightarrow 0$  and hence  $w_t^\top x_i \rightarrow \infty$ .
4.  $\sum_{t=0}^{\infty} \|\nabla L(w_t)\|_* = \infty$ .

Given these two Propositions, the proof proceeds in two steps. We first establish that the loss converges to zero sufficiently quickly to lower bound the unnormalized margin  $\min_i w_t^\top x_i$ . Next, we upper bound  $\|w_t\|$ . By dividing the lower bound in the first step by the upper bound in the second step, we can lower bound the normalized margin, which will complete the proof.

*Proof of Theorem 9.3.2. Step 1: Lower bound the unnormalized margin.* First, we establish the loss converges sufficiently quickly. Define  $\gamma_t = \|\nabla L(w_t)\|_*$ . From Taylor's theorem,

$$\begin{aligned} L(w_{t+1}) &\leq \\ L(w_t) + \eta_t \langle \nabla L(w_t), \Delta w_t \rangle + \sup_{\beta \in (0,1)} \frac{\eta_t^2}{2} \Delta w_t^\top \nabla^2 L(w_t + \beta \eta_t \Delta w_t) \Delta w_t \\ &\stackrel{(a)}{\leq} L(w_t) - \eta_t \|\nabla L(w_t)\|_*^2 + \frac{\eta_t^2 B^2}{2} \sup_{\beta \in (0,1)} L(w_t + \beta \eta_t \Delta w_t) \|\Delta w_t\|^2 \\ &\stackrel{(b)}{\leq} L(w_t) - \eta_t \|\nabla L(w_t)\|_*^2 + \frac{\eta_t^2 B^2}{2} L(w_t) \|\Delta w_t\|^2 \end{aligned} \tag{9.10}$$

where (a) uses  $v^\top \nabla^2 L(w)v \leq L(w)B^2\|v\|^2$  and (b) uses Proposition 9.3.4 part 1 and convexity to show  $\sup_{\beta \in (0,1)} L(w_t + \beta\eta_t \Delta w_t) \leq L(w_t)$ .

From eq. 9.10, using  $\gamma_t = \|\nabla L(w_t)\|_* = \|\Delta w_t\|$ , we have that

$$\begin{aligned} L(w_{t+1}) &\leq L(w_t) - \eta\gamma_t^2 + \frac{\eta^2 B^2 L(w_t) \gamma_t^2}{2} \\ &= L(w_t) \left[ 1 - \frac{\eta\gamma_t^2}{L(w_t)} + \frac{\eta^2 B^2 \gamma_t^2}{2} \right] \\ &\stackrel{(a)}{\leq} L(w_t) \exp\left(-\frac{\eta\gamma_t^2}{L(w_t)} + \frac{\eta^2 B^2 \gamma_t^2}{2}\right) \\ &\stackrel{(b)}{\leq} L(w_0) \exp\left(-\sum_{u \leq t} \frac{\eta\gamma_u^2}{L(w_u)} + \sum_{u \leq t} \frac{\eta^2 B^2 \gamma_u^2}{2}\right), \end{aligned} \quad (9.11)$$

where we get (a) by using  $(1+x) \leq \exp(x)$ , and (b) by recursing the argument.

Next, we lower bound the unnormalized margin. From eq. (9.11), we have,

$$\begin{aligned} \max_{n \in [N]} \exp(-\langle w_{t+1}, x_n \rangle) &\leq L(w_{(t+1)}) \\ &\leq L(w_0) \exp\left(-\sum_{u \leq t} \frac{\eta\gamma_u^2}{L(w_u)} + \sum_{u \leq t} \frac{\eta^2 B^2 \gamma_u^2}{2}\right) \end{aligned} \quad (9.12)$$

By applying  $-\log$ ,

$$\min_{n \in [N]} \langle w_{t+1}, x_n \rangle \geq \sum_{u \leq t} \frac{\eta\gamma_u^2}{L(w_u)} - \sum_{u \leq t} \frac{\eta^2 B^2 \gamma_u^2}{2} - \log L(w_0). \quad (9.13)$$

**Step 2: Upper bound  $\|w_{t+1}\|$ .** Using  $\|\Delta w_u\| = \|\nabla L(w_u)\|_* = \gamma_u$ , we have,

$$\|w_{t+1}\| \leq \|w_0\| + \sum_{u \leq t} \eta \|\Delta w_u\| \leq \|w_0\| + \sum_{u \leq t} \eta \gamma_u. \quad (9.14)$$

To complete the proof, we simply combine Equations (9.13) and (9.14) to lower bound the normalized margin.

$$\begin{aligned} \frac{\langle w_{t+1}, x_n \rangle}{\|w_{t+1}\|} &\geq \frac{\sum_{u \leq t} \frac{\eta\gamma_u^2}{L(w_u)}}{\sum_{u \leq t} \eta\gamma_u + \|w_0\|} - \left( \frac{\sum_{u \leq t} \frac{\eta^2 B^2 \gamma_u^2}{2} + \log L(w_0)}{\|w_{t+1}\|} \right). \\ &:= (I) \qquad \qquad \qquad + (II). \end{aligned} \quad (9.15)$$

For term (I), from Proposition 9.3.3, we have  $\gamma_u = \|\nabla L(w_u)\|_* \geq \gamma L(w_u)$ . Hence the numerator is lower bounded  $\sum_{u \leq t} \frac{\eta\gamma_u^2}{L(w_u)} \geq$



$\gamma \sum_{u \leq t} \eta \gamma_u$ . We have

$$\frac{\sum_{u \leq t} \frac{\eta \gamma_u^2}{L(w_u)}}{\sum_{u \leq t} \eta \gamma_u + \|w_0\|} \geq \gamma \frac{\sum_{u \leq t} \eta \gamma_u}{\sum_{u \leq t} \eta \gamma_u + \|w_0\|} \rightarrow \gamma, \quad (9.16)$$

using  $\sum_{u \leq t} \eta \gamma_u \rightarrow \infty$  and  $\|w_0\| < \infty$  from Proposition 9.3.4.

For term (II),  $\log L(w_0) < \infty$  and  $\sum_{u \leq t} \frac{\eta^2 B^2 \gamma_u^2}{2} < \infty$  using Proposition 9.3.3. Thus (II)  $\rightarrow 0$ .

Using the above in Equation (9.15), we obtain

$$\lim_{t \rightarrow \infty} \frac{w_{t+1}^\top x_i}{\|w_{t+1}\|} \geq \gamma := \max_{\|w\| \leq 1} \min_i \frac{w^\top x_i}{\|w\|}.$$

□

## 9.4 Homogeneous Models with Exponential Tailed Loss

«Suriya notes: Jason: I think we should give Kaifeng's proof here. Its more general and concurrent work.» In this section, we consider the asymptotic behavior of gradient descent when the prediction function is homogeneous in the parameters. Consider the loss

$$L(w) = \sum_{i=1}^n \exp(-y_i f_i(w)), \quad (9.17)$$

where  $f_i(cw) = c^\alpha f_i(w)$  is  $\alpha$ -homogeneous. Typically,  $f_i(w)$  is the output of the prediction function such as a deep network. Similar to the linear case in Section 6.5.1, there is a related maximum margin problem. Define the optimal margin as  $\gamma = \max_{\|w\|_2=1} \min_i y_i f_i(w)$ . The associated non-linear margin maximization is given by the following non-convex constrained optimization:

$$\min \|w\|^2 \text{ st } y_i f_i(w) \geq \gamma. \quad (\text{Max-Margin})$$

Analogous to Section 6.5.1, we expect that gradient descent on Equation (9.17) converges to the optimum of the Max-Margin problem (Max-Margin). However, the max-margin problem itself is a constrained non-convex problem, so we cannot expect to attain a global optimum. Instead, we show that gradient descent iterates converge to first-order stationary points of the max-margin problem.

**Definition 9.4.1** (First-order Stationary Point). *The first-order optimality conditions of Max-Margin are:*

1.  $\forall i, y_i f_i(w) \geq \gamma$
2. There exists Lagrange multipliers  $\lambda \in \mathbb{R}_+^N$  such that  $w = \sum_n \lambda_n \nabla f_n(w)$  and  $\lambda_n = 0$  for  $n \notin S_m(w) := \{i : y_i f_i(w) = \gamma\}$ , where  $S_m(w)$  is the set of support vectors.

We denote by  $\mathcal{W}^*$  the set of first-order stationary points.

Let  $w_t$  be the iterates of gradient flow (gradient descent with step-size tending to zero). Define  $\ell_{it} = \exp(-f_i(w_t))$  and  $\ell_t$  be the vector with entries  $\ell_{it}$ . The following two assumptions assume that the limiting direction  $\frac{w_t}{\|w_t\|}$  exist and the limiting direction of the losses  $\frac{\ell_t}{\|\ell_t\|_1}$  exist. Such assumptions are natural in the context of max-margin problems, since we want to argue that  $w_t$  converges to a max-margin direction, and also the losses  $\ell_t / \|\ell_t\|_1$  converges to an indicator vector of the support vectors. We will directly assume these limits exist, though this is proved in <sup>3</sup>.

3

**Assumption 9.4.2** (Smoothness). *We assume  $f_i(w)$  is a  $C^2$  function.*

**Assumption 9.4.3** (Asymptotic Formulas). *Assume that  $L(w_t) \rightarrow 0$ , that is we converge to a global minimizer. Further assume that  $\lim_{t \rightarrow \infty} \frac{w_t}{\|w_t\|_2}$  and*

*$\lim_{t \rightarrow \infty} \frac{\ell_t}{\|\ell_t\|_1}$  exist. Equivalently,*

$$\ell_{nt} = h_t a_n + h_t \epsilon_{nt} \quad (9.18)$$

$$w_t = g_t \bar{w} + g_t \delta_t, \quad (9.19)$$

*with  $\|a\|_1 = 1$ ,  $\|\bar{w}\|_2 = 1$ ,  $\lim_{t \rightarrow \infty} h(t) = 0$ ,  $\lim_{t \rightarrow \infty} \epsilon_{nt} = 0$ , and  $\lim_{t \rightarrow \infty} \delta_t = 0$ .*

**Assumption 9.4.4** (Linear Independence Constraint Qualification).

*Let  $w$  be a unit vector. LICQ holds at  $w$  if the vectors  $\{\nabla f_i(w)\}_{i \in \mathcal{S}_m(w)}$  are linearly independent.*

Constraint qualification allow the first-order optimality conditions of Definition 9.4.1 to be a necessary condition for optimality. Without constraint qualifications, even the global optimum may not satisfy the optimality conditions.

For example in linear SVM, LICQ is ensured if the support vectors  $x_i$  are linearly independent then LICQ holds. For data sampled from an absolutely continuous distribution, the linear SVM solution will always have linearly independent support vectors.

**Theorem 9.4.5.** *Define  $\bar{w} = \lim_{t \rightarrow \infty} \frac{w_t}{\|w_t\|}$ . Under Assumptions 9.4.2, 9.4.3, and 9.4.4,  $\bar{w} \in \mathcal{W}$  is a first-order stationary point of (Max-Margin).*

*Proof.* Define  $S = \{i : f_i(\bar{w}) = \gamma\}$ , where  $\gamma$  is the optimal margin attainable by a unit norm  $w$ .

**Lemma 9.4.6.** *Under the setting of Theorem 9.4.5,*

$$\nabla f_i(w_t) = \nabla f_i(g_t \bar{w}) + O(B g_t^{\alpha-1} \|\delta_t\|). \quad (9.20)$$

*For  $i \in S$ , the second term is asymptotically negligible as a function of  $t$ ,*

$$\nabla f_i(w_t) = \nabla f_i(g_t \bar{w}) + o(\nabla f_i(g_t \bar{w})).$$

**Lemma 9.4.7.** *Under the conditions of Theorem 9.4.5,  $a_i = 0$  for  $i \notin S$ .*

From the gradient flow dynamics,

$$\begin{aligned}\dot{w}(t) &= \sum_i \exp(-f_i(w_t)) \nabla f_i(w_t) \\ &= \sum_i (h_t a_i + h_t \epsilon_{it}) (\nabla f_i(g_t \bar{w}) + \Delta_{it}),\end{aligned}$$

where  $\Delta_i(t) = \int_{s=0}^{s=1} \nabla^2 f_i(g_t \bar{w} + s g_t \delta_t) g_t \delta_t ds$ . By expanding and using  $a_i = 0$  for  $n \notin S$  (Lemma 9.4.7),

$$\begin{aligned}\dot{w}_t &= \underbrace{\sum_{i \in S} h_t a_i \nabla f_i(g_t \bar{w})}_I \\ &\quad + \underbrace{h_t \sum_{i \in S} a_i \Delta_{it}}_{II} + \underbrace{h_t \sum_i \epsilon_{it} \nabla f_i(g_t \bar{w})}_{III} + \underbrace{\sum_i h_t \epsilon_{it} \Delta_{it}}_{IV}\end{aligned}$$

Via Assumption 9.4.4, term  $I = \Omega(g_t^{\alpha-1} h_t)$  and from Lemma 9.4.6,  $II = o(I)$ . Using these, the first term  $I$  is the largest and so after normalizing,

$$\frac{\dot{w}_t}{\|\dot{w}_t\|} = \sum_{i \in S} a_i \nabla f_i(g_t \bar{w}) + o(1). \quad (9.21)$$

Since  $\lim_t \frac{w_t}{\|w_t\|} = \lim_t \frac{\dot{w}_t}{\|\dot{w}_t\|}$  [GLSS18a], then

$$\lim_{t \rightarrow \infty} \frac{w_t}{\|w_t\|} = \sum_{i \in S} \nabla f_i(g_t \bar{w}). \quad (9.22)$$

Thus we have shown  $w$  satisfies the first-order optimality condition of Definition 9.4.1.  $\square$

## 9.5 Induced bias in function space

«Suriya notes: Jason: can you introduce the idea of induced biases and give special results for linear convnets, any relevant results from yours+tengyu's margin paper, and infinite width 2 layer ReLU network?»»



## *Ultra-wide Neural Networks and Neural Tangent Kernels*

This chapter concerns a model that seems ridiculous at first sight: one with infinitely many nodes at inner layers. To understand why it is actually interesting, let's recall the mysteries we're trying to understand.

Training a neural network is a non-convex optimization problem, and in the worst case, it is NP-hard [BR89]. On the other hand, empirically, simple gradient algorithms like stochastic gradient descent can often achieve zero training loss, i.e., the simple algorithm can find a neural network that fits all training data. Furthermore, one can still observe this phenomenon for nonsensical data, when the original labels are replaced with random labels [ZBH<sup>+</sup>16b].

*Key role of overparametrization.* The fact that networks can fit nonsensical data perfectly is not surprising because the nets are very over-parameterized. For example, Wide ResNet when trained on ImageNet has 100x more parameters than the number of training datapoints. Recall from Chapter 3 that under such conditions even linear regression (solved via gradient descent) can perfectly fit training data. But this does not explain real-life neural nets because we still need to prove that: (a) the low loss can be attained by a *gradient-based training* starting from a *random initialization* (b) that the trained net has good generalization when trained on proper data. Many traditional generalization bounds yield vacuous guarantees, as mentioned in Chapter 6.

*Teacher/Student Nets.* A difficulty that arises while addressing this research agenda is that clearly at some point the theory should take properties of the data into account, and real-life data (e.g., images) have no good description. One route taken by researchers is to assume that the labels for training data were computed via

some ground truth net sometimes referred to as *teacher net*. Thus the net being trained is thought of as a *student net* and the goal of good generalization is to be able to produce labels broadly in agreement with the teacher net.

Given the importance of overparametrization in real life, it is natural to allow the student net to be much deeper or wider than the teacher net.<sup>1</sup>

*Infinite nets and NTKs:* Now we explain the model of neural networks whose width in the inner layers is very large, essentially going to infinity. For example, imagine a standard net like AlexNet being allowed to expand its inner layers to have unlimited width. (In convolutional layers this means allowing the filters to have unlimited width.) This hugely inflated version of AlexNet architecture still takes the same input as before but its training and generalization could potentially change a lot. Researchers studied these architectures and quickly realized that at least one way of initialization/training leads to the net turning into a kernel classifier, called *Neural Tangent Kernel* (NTK)<sup>2</sup>. This chapter is an introduction to infinitely wide nets and NTKs. We'll see that behavior of NTKs can be computed efficiently by efficient algorithms for computing the kernel inner product for NTK. NTKs do exhibit good optimization (i.e. convergence to low training loss) and reasonable generalization behavior but are not as good as their finite counterparts. For example the NTK corresponding to AlexNet generalizes reasonably OK on image data but with far worse accuracy than AlexNet. We will discuss some practical uses of NTK for small-data tasks.

### 10.1 Evolution of the trained net

This section derives evolution of nets during training under least squares loss. It applies to any net, and the simple expression will play a key role in NTK theory.

We denote by  $f(w, x) \in \mathbb{R}$  the output of a neural network where  $w \in \mathbb{R}^N$  is all the parameters in the network and  $x \in \mathbb{R}^d$  is the input. Given a training dataset  $\{(x_i, y_i)\}_{i=1}^n \subset \mathbb{R}^d \times \mathbb{R}$ , consider training the neural network by minimizing the squared loss over training data:

$$\ell(w) = \frac{1}{2} \sum_{i=1}^n (f(w, x_i) - y_i)^2.$$

For simplicity, in this chapter, we study gradient flow, a.k.a., gradient decent with infinitesimally small learning rate. In this case, the dynamics can be described by an ordinary differential equation

<sup>1</sup> Indeed in experiments one finds that if synthetically labeled data is generated by passing inputs from a distribution through a teacher net, then teaching a new net from scratch to mimic the teacher is much easier in practice if the new net is allowed to be significantly bigger.

<sup>2</sup> Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and generalization in neural networks. In *Advances in neural information processing systems*, pages 8571–8580, 2018

(ODE):

$$\frac{dw(t)}{dt} = -\nabla \ell(w(t)).$$

Note this is the dynamics of the parameters. The following lemma describes the dynamics of the predictions on training data points.

**Lemma 10.1.1.** *Let  $u(t) = (f(w(t), x_i))_{i \in [n]} \in \mathbb{R}^n$  be the network outputs on all  $x_i$ 's at time  $t$ , and  $y = (y_i)_{i \in [n]}$  be the labels. Then  $u(t)$  follows the following evolution, where  $H(t)$  is an  $n \times n$  positive semidefinite matrix whose  $(i, j)$ -th entry is  $\left\langle \frac{\partial f(w(t), x_i)}{\partial w}, \frac{\partial f(w(t), x_j)}{\partial w} \right\rangle$ :*

$$\frac{du(t)}{dt} = -H(t) \cdot (u(t) - y). \quad (10.1)$$

*Proof of Lemma 10.1.1.* The parameters  $w$  evolve according to the differential equation

$$\frac{dw(t)}{dt} = -\nabla \ell(w(t)) = -\sum_{i=1}^n (f(w(t), x_i) - y_i) \frac{\partial f(w(t), x_i)}{\partial w}, \quad (10.2)$$

where  $t \geq 0$  is a continuous time index. Under Equation (10.2), the evolution of the network output  $f(w(t), x_i)$  can be written as

$$\frac{df(w(t), x_i)}{dt} = -\sum_{j=1}^n (f(w(t), x_j) - y_j) \left\langle \frac{\partial f(w(t), x_i)}{\partial w}, \frac{\partial f(w(t), x_j)}{\partial w} \right\rangle. \quad (10.3)$$

Since  $u(t) = (f(w(t), x_i))_{i \in [n]} \in \mathbb{R}^n$  is the network outputs on all  $x_i$ 's at time  $t$ , and  $y = (y_i)_{i \in [n]}$  is the desired outputs, Equation (10.3) can be written more compactly as

$$\frac{du(t)}{dt} = -H(t) \cdot (u(t) - y), \quad (10.4)$$

where  $H(t) \in \mathbb{R}^{n \times n}$  is a kernel matrix defined as  $[H(t)]_{ij} = \left\langle \frac{\partial f(w(t), x_i)}{\partial w}, \frac{\partial f(w(t), x_j)}{\partial w} \right\rangle$  ( $\forall i, j \in [n]$ ).  $\square$

### 10.1.1 Behavior in the infinite limit

Recall that we're interested in the limit of deep net training when the training set is fixed, the width goes to infinity, and for a suitable scale of initialization (which depends on the width). Under fairly general conditions it can be shown that the matrix  $H(t)$  remains roughly constant during training i.e., roughly equal to  $H(0)$ . Furthermore, the matrix  $H(0)$ , whose definition involves random weights used at initialization, converges in probability to the Gram Matrix  $H^*$  of the training dataset (see Chapter 3) with respect to certain kernel, called the *Neural Tangent Kernel*. Then Equation (10.1) becomes

$$\frac{du(t)}{dt} = -H^* \cdot (u(t) - y). \quad (10.5)$$

In other words, least squares kernel regression as described in Chapter 3, but with infinitesimally small learning rate. Recall that the final classifier is described as

$$f^*(x) = (k(x, x_1), \dots, k(x, x_n)) \cdot (H^*)^{-1}y. \quad (10.6)$$

## 10.2 NTK: Simple 2-layer example

In this section, we develop the theory in context of a simple two-layer neural network of the following form:

$$f(a, W, x) = \frac{1}{\sqrt{m}} \sum_{r=1}^m a_r \sigma(w_r^\top x) \quad (10.7)$$

where  $\sigma(\cdot)$  is the activation function. Here we assume  $|\dot{\sigma}(z)|$  and  $|\ddot{\sigma}(z)|$  are bounded by 1 for all  $z \in \mathbb{R}$  and For example, soft-plus activation function,  $\sigma(z) = \log(1 + \exp(z))$ , satisfies this assumption.<sup>3</sup> We also assume input  $x$  is a unit vector, i.e.,  $\|x\|_2 = 1$ . The scaling  $1/\sqrt{m}$  will play an important role in proving  $H(t)$  stays close to the Gram Matrix  $H^*$  for a fixed kernel. Throughout the section, to measure the closeness of two matrices  $A$  and  $B$ , we use the operator norm  $\|\cdot\|_2$ . We will use the fact that if corresponding entries are close then so are their spectral properties  $A, B$ . This follows from the fact that the sum of squared differences across coordinates,  $\|A - B\|_F^2$ , is also an upper bound on  $\|A - B\|_2^2$ .

We use random initialization  $w_r(0) \sim N(0, I)$  and  $a_r \sim \text{Unif}[\{-1, 1\}]$ . For simplicity, we will only optimize the first layer, i.e.,  $W = [w_1, \dots, w_m]$ . Note this is still a non-convex optimization problem.

We can first calculate  $H(0)$  and show as  $m \rightarrow \infty$ ,  $H(0)$  converges to a fixed matrix  $H^*$ . Note  $\frac{\partial f(a, W, x_i)}{\partial w_r} = \frac{1}{\sqrt{m}} a_r x_i \sigma'(w_r^\top x_i)$ . Therefore, each entry of  $H(0)$  admits the formula

$$\begin{aligned} [H(0)]_{ij} &= \sum_{r=1}^m \left\langle \frac{\partial f(a, W(0), x_i)}{\partial w_r(0)}, \frac{\partial f(a, W(0), x_j)}{\partial w_r(0)} \right\rangle \\ &= \sum_{r=1}^m \left\langle \frac{1}{\sqrt{m}} a_r x_i \dot{\sigma}(w_r(0)^\top x_i), \frac{1}{\sqrt{m}} a_r x_j \dot{\sigma}(w_r(0)^\top x_j) \right\rangle \\ &= x_i^\top x_j \cdot \frac{\sum_{r=1}^m \dot{\sigma}'(w_r(0)^\top x_i) \dot{\sigma}'(w_r(0)^\top x_j)}{m} \end{aligned}$$

Here the last step we used  $a_r^2 = 1$  for all  $r = 1, \dots, m$  because we initialize  $a_r \sim \text{Unif}[\{-1, 1\}]$ . Recall every  $w_r(0)$  is i.i.d. sampled from a standard Gaussian distribution. Therefore, one can view  $[H(0)]_{ij}$  as the average of  $m$  i.i.d. random variables. If  $m$  is large, then by the law of large number, we know this average is close to the expectation of the random variable. Here the expectation is the NTK evaluated on  $x_i$

<sup>3</sup>Note rectified linear unit (ReLU) activation function does not satisfy this assumption. However, one can use a specialized analysis of ReLU to show  $H(t) \approx H^*$  [DZPS18].



and  $x_j$ :

$$H_{ij}^* \triangleq x_i^\top x_j \cdot \mathbb{E}_{w \sim \mathcal{N}(0, I)} \left[ \sigma'(w^\top x_i) \sigma'(w^\top x_j) \right]$$

**Problem 10.2.1.** *If the activation  $\sigma$  is ReLU then (noting that it is differentiable everywhere except at one point) then show that*

$$H_{ij}^* = \mathbb{E}_{w \sim \mathcal{N}(0, I)} \left[ \dot{\sigma} w^\top x \dot{\sigma} w^\top x' \right] = \frac{\pi - \arccos \left( \frac{x^\top x'}{\|x\|_2 \|x'\|_2} \right)}{2\pi}. \quad (10.8)$$

Using Hoeffding inequality and the union bound, one can easily obtain the following bound characterizing  $m$  and the closeness of  $H(0)$  and  $H^*$ .

**Lemma 10.2.2** (Perturbation on the Initialization, [DZPS19, SY19]). *Fix some  $\epsilon > 0$ . If  $m = \Omega(\epsilon^{-2} n^2 \log(n/\delta))$ , then with probability at least  $1 - \delta$  over  $w_1(0), \dots, w_m(0)$ , we have*

$$\|H(0) - H^*\|_2 \leq \epsilon.$$

*Proof of Lemma 10.2.2.* We first fixed an entry  $(i, j)$ . Note

$$\left| x_i^\top x_j \sigma'(w_t(0)^\top x_i) \sigma'(w_r(0)^\top x_j) \right| \leq 1.$$

Applying Hoeffding inequality, we have with probability  $1 - \frac{\delta}{n^2}$ ,

$$|[H(0)]_{i,j} - H_{i,j}^*| \leq \left( \frac{2}{m} \log(2n^2/\delta) \right)^{1/2} \leq 4 \left( \frac{\log(n/\delta)}{m} \right)^{1/2} \leq \frac{\epsilon}{n}.$$

Next, applying the union bound over all pairs  $(i, j) \in [n] \times [n]$ , we have for all  $(i, j)$ ,  $|[H(0)]_{i,j} - H_{i,j}^*| \leq \frac{\epsilon}{n^2}$ . To establish the operator norm bound, we simply use the following chain of inequalities

$$\begin{aligned} \|H(0) - H^*\|_2 &\leq \|H(0) - H^*\|_F \\ &= \left( \sum_{ij} |[H(0)]_{i,j} - H_{i,j}^*|^2 \right)^{1/2} \\ &\leq (n^2 \cdot \frac{\epsilon^2}{n^2})^{1/2} = \epsilon. \end{aligned}$$

□

Now we proceed to show during training,  $H(t)$  is close to  $H(0)$ . Formally, we prove the following lemma.

**Lemma 10.2.3.** *Assume  $y_i = O(1)$  for all  $i = 1, \dots, n$ . Given  $t > 0$ , suppose that for all  $0 \leq \tau \leq t$ ,  $u_i(\tau) = O(1)$  for all  $i = 1, \dots, n$ . If  $m = \Omega\left(\frac{n^6 t^2}{\epsilon^2}\right)$ , we have*

$$\|H(t) - H(0)\|_2 \leq \epsilon.$$

*Proof of Lemma 10.2.3.* The first key idea is to show that *every weight vector only moves little if  $m$  is large*. To show this, let us calculate the movement of a single weight vector  $w_r$ .

$$\begin{aligned}
\|w_r(t) - w_r(0)\|_2 &= \left\| \int_0^t \frac{dw_r(\tau)}{d\tau} d\tau \right\|_2 \\
&= \left\| \int_0^t \frac{1}{\sqrt{m}} \sum_{i=1}^n (u_i(\tau) - y_i) a_r x_i \dot{\sigma}(w_r(\tau)^\top x_i) d\tau \right\|_2 \\
&\leq \frac{1}{\sqrt{m}} \int_0^t \left\| \sum_{i=1}^n (u_i(\tau) - y_i) a_r x_i \dot{\sigma}(w_r(\tau)^\top x_i) \right\|_2 d\tau \\
&\leq \frac{1}{\sqrt{m}} \sum_{i=1}^n \int_0^t \|u_i(\tau) - y_i a_r x_i \dot{\sigma}(w_r(\tau)^\top x_i)\|_2 d\tau \\
&\leq \frac{1}{\sqrt{m}} \sum_{i=1}^n \int_0^t O(1) d\tau \\
&= O\left(\frac{tn}{\sqrt{m}}\right).
\end{aligned}$$

This calculation shows that at any given time  $t$ ,  $w_r(t)$  is close to  $w_r(0)$ , as long as  $m$  is large. Next, we show this implies the kernel matrix  $H(t)$  is close  $H(0)$ . We calculate the difference on a single entry.

$$\begin{aligned}
&[H(t)]_{ij} - [H(0)]_{ij} \\
&= \left| \frac{1}{m} \sum_{r=1}^m \left( \dot{\sigma}(w_r(t)^\top x_i) \dot{\sigma}(w_r(t)^\top x_j) - \dot{\sigma}(w_r(0)^\top x_i) \dot{\sigma}(w_r(0)^\top x_j) \right) \right| \\
&\leq \frac{1}{m} \sum_{r=1}^m \left| \dot{\sigma}(w_r(t)^\top x_i) \left( \dot{\sigma}(w_r(t)^\top x_j) - \dot{\sigma}(w_r(0)^\top x_j) \right) \right| \\
&\quad + \frac{1}{m} \sum_{r=1}^m \left| \dot{\sigma}(w_r(0)^\top x_j) \left( \dot{\sigma}(w_r(t)^\top x_i) - \dot{\sigma}(w_r(0)^\top x_i) \right) \right| \\
&\leq \frac{1}{m} \sum_{r=1}^m \left| \max_r \dot{\sigma}(w_r(t)^\top x_i) \|x_i\|_2 \|w_r(t) - w_r(0)\|_2 \right| \\
&\quad + \frac{1}{m} \sum_{r=1}^m \left| \max_r \dot{\sigma}(w_r(t)^\top x_i) \|x_i\|_2 \|w_r(t) - w_r(0)\|_2 \right| \\
&= \frac{1}{m} \sum_{r=1}^m O\left(\frac{tn}{\sqrt{m}}\right) \\
&= O\left(\frac{tn}{\sqrt{m}}\right).
\end{aligned}$$

Therefore, using the same argument as in Lemma 10.2.2, we have

$$\|H(t) - H(0)\|_2 \leq \sum_{i,j} |[H(t)]_{ij} - [H(0)]_{ij}| = O\left(\frac{tn^3}{\sqrt{m}}\right).$$

Plugging our assumption on  $m$ , we finish the proof.  $\square$

Several remarks are in sequel.

**Remark 1:** The assumption that  $y_i = O(1)$  is a mild assumption because in practice most labels are bounded by an absolute constant.

**Remark 2:** The assumption on  $u_i(\tau) = O(1)$  for all  $\tau \leq t$  and  $m$ 's dependency on  $t$  can be relaxed. This requires a more refined analysis. See [DZPS19].

**Remark 3:** One can generalize the proof for multi-layer neural network. See [ADH<sup>+</sup>19b] for more details.

**Remark 4:** While we only prove the continuous time limit, it is not hard to show with small learning rate (discrete time) gradient descent,  $H(t)$  is close to  $H^*$ . See [DZPS19].

### 10.3 Explaining Optimization and Generalization of Ultra-wide Neural Networks via NTK

Now we have established the following approximation

$$\frac{du(t)}{dt} \approx -H^* \cdot (u(t) - y) \quad (10.9)$$

where  $H^*$  is the NTK matrix. Now we use this approximation to analyze the optimization and generalization behavior of ultra-wide neural networks.

*Understanding Optimization* The dynamics of  $u(t)$  that follows

$$\frac{du(t)}{dt} = -H^* \cdot (u(t) - y)$$

is actually linear dynamical system. For this dynamics, there is a standard analysis. We denote the eigenvalue decomposition of  $H^*$  as

$$H^* = \sum_{i=1}^n \lambda_i v_i v_i^\top$$

where  $\lambda_1 \geq \dots \geq \lambda_n \geq 0$  are eigenvalues and  $v_1, \dots, v_n$  are eigenvectors. With this decomposition, we consider the dynamics of  $u(t)$  on each eigenvector *separately*. Formally, fixing an eigenvector  $v_i$  and multiplying both side by  $v_i$ , we obtain

$$\begin{aligned} \frac{dv_i^\top u(t)}{dt} &= -v_i^\top H^* \cdot (u(t) - y) \\ &= -\lambda_i \left( v_i^\top (u(t) - y) \right). \end{aligned}$$

Observe that the dynamics of  $v_i^\top u(t)$  only depends on itself and  $\lambda_i$ , so this is actually a one dimensional ODE. Moreover, this ODE

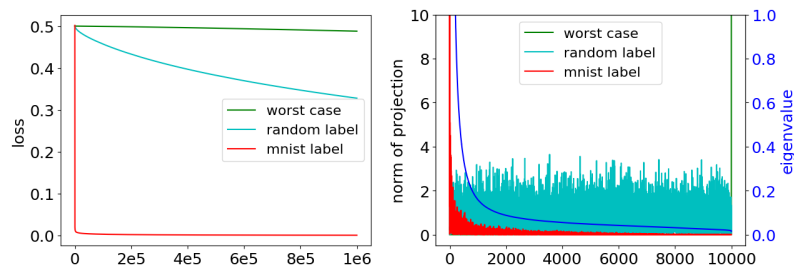


Figure 10.1: Convergence rate vs. projections onto eigenvectors of the kernel matrix.

admits an analytical solution

$$v_i^\top (u(t) - y) = \exp(-\lambda_i t) \left( v_i^\top (u(0) - y) \right). \quad (10.10)$$

Now we use Equation (10.10) to explain why we can find a zero training error solution. We need to assume  $\lambda_i > 0$  for all  $i = 1, \dots, n$ , i.e., all eigenvalues of this kernel matrix are strictly positive. One can prove this under fairly general conditions. See [DZPS19, DLL<sup>+</sup>18].

Observe that  $(u(t) - y)$  is the difference between predictions and training labels at time  $t$  and the algorithm finds a 0 training error solutions means as  $t \rightarrow \infty$ , we have  $u(t) - y \rightarrow 0$ . Equation (10.10) implies that each component of this difference, i.e.,  $v_i^\top (u(t) - y)$  is converging to 0 exponentially fast because of the  $\exp(-\lambda_i t)$  term. Furthermore, notice that  $\{v_1, \dots, v_n\}$  forms an orthonormal basis of  $\mathbb{R}^n$ , so  $(u(t) - y) = \sum_{i=1}^n v_i^\top (u(t) - y) v_i$ . Since we know each  $v_i^\top (u_i(t) - y) \rightarrow 0$ , we can conclude that  $(u(t) - y) \rightarrow 0$  as well.

Equation (10.10) actually gives us more information about the convergence. Note each component  $v_i^\top (u(t) - y)$  converges to 0 at a different rate. The component that corresponds to larger  $\lambda_i$  converges to 0 at a faster rate than the one with a smaller  $\lambda_i$ . For a set of labels, in order to have faster convergence, we would like the projections of  $y$  onto the top eigenvectors to be larger.<sup>4</sup> Therefore, we obtain the following intuitive rule to compare the convergence rates in a qualitative manner (for fixed  $\|y\|_2$ ):

- For a set of labels  $y$ , if they align with top eigenvectors, i.e.,  $(v_i^\top y)$  is large for large  $\lambda_i$ , then gradient descent converges quickly.
- For a set of labels, if the projections on eigenvectors  $\{(v_i^\top y)\}_{i=1}^n$  are uniform, or labels align with eigenvectors with respect to small eigenvalues, then gradient descent converges with a slow rate.

We can verify this phenomenon experimentally. In Figure 10.1, we compare convergence rates of gradient descent between using original labels, random labels and the worst case labels (normalized eigenvector of  $H^*$  corresponding to  $\lambda_n$ ). We use the neural network architecture defined in Equation (10.7) with ReLU activation function

<sup>4</sup> Here we ignore the effect of  $u(0)$  for simplicity. See [ADH<sup>+</sup>19a] on how to mitigate the effect on  $u(0)$ .

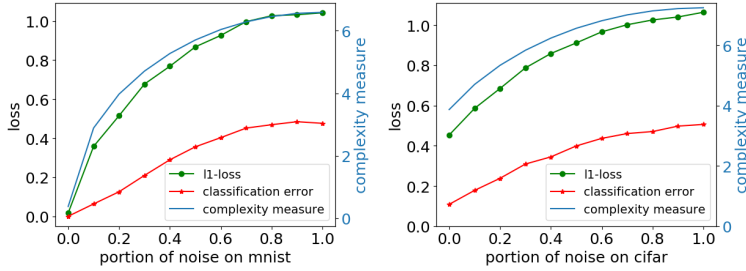


Figure 10.2: Generalization error vs. complexity measure.

and only train the first layer. In the right figure, we plot the eigenvalues of  $H^*$  as well as projections of true, random, and worst case labels on different eigenvectors of  $H^*$ . The experiments use gradient descent on data from two classes of MNIST. The plots demonstrate that original labels have much better alignment with top eigenvectors, thus enjoying faster convergence.

### 10.3.1 Understanding Generalization in 2-layer setting

The approximation in Equation (10.9) implies the final prediction function of ultra-wide neural network is approximately the kernel prediction function defined in Equation (10.6). Therefore, we can just use the generalization theory for kernels to analyze the generalization behavior of ultra-wide neural networks. For the kernel prediction function defined in Equation (10.6), we can use Rademacher complexity bound to derive the following generalization bound for 1-Lipschitz loss function (which is an upper bound of classification error):

$$\frac{\sqrt{2y^\top (H^*)^{-1} y \cdot \text{tr}(H^*)}}{n}. \quad (10.11)$$

This is a *data-dependent* complexity measure that upper bounds the generalization error.

We can check this complexity measure empirically. In Figure 10.2, we compare the generalization error ( $\ell_1$  loss and classification error) with this complexity measure. We vary the portion of random labels in the dataset to see how the generalization error and the complexity measure change. We use the neural network architecture defined in Equation (10.7) with ReLU activation function and only train the first layer. The left figure uses data from two classes of MNIST and the right figure uses two classes from CIFAR. This complexity measure almost matches the trend of generalization error as the portion of random labels increases.

*Learning from simple teacher nets.* Now we explain why NTK can learn some functions that can be expressed as two-layer nets. Since we know the optimization error goes to 0 for any label, so it is sufficient to study what teacher nets enables the generalization error to be small. Concretely, we give some examples of two-layer teach nets that make generalization bound (10.11) goes to 0 as  $n \rightarrow \infty$ .

**Linear Function:** We begin with a simple example. Assume the label  $y = \beta^\top x$  for some vector  $\beta$ . Then one can show that (10.11) is upper bounded by  $O\left(\frac{\|\beta\|_2}{\sqrt{n}}\right)$ . Therefore, NTK can learn linear functions with bounded coefficients.

**Two-layer nets with Polynomial Activation:** Consider  $y = \sum_{j=1}^k \alpha_j \left(\beta_j^\top x\right)^p$ , i.e., a two-layer net with the activation function being  $z^p$  with  $p$  being an even number. Similar to the linear function, one can show that (10.11) is upper bounded by  $O\left(\frac{p \sum_{j=1}^k |\alpha_j| \|\beta_j\|_2^p}{\sqrt{n}}\right)$ . Therefore, we can argue NTK can learn two-layer polynomial nets with bounded coefficients.

**Cosine activation** Beyond polynomials, one can also show NTK can learn somewhat bizarre function. For example, if  $y = \sum_{j=1}^k \alpha_j \left(\cos\left(\beta_j^\top x\right) - 1\right)$ , then we can bound (10.11) is by  $O\left(\frac{p \sum_{j=1}^k |\alpha_j| \|\beta_j\|_2 \sinh(\|\beta_j\|_2^2)}{\sqrt{n}}\right)$ .

All the these examples can be proved by a general technique based on Taylor expansion of NTK. See [ADH<sup>+</sup>19a].

#### 10.4 NTK formula for Multilayer Fully-connected Neural Network

In this section we show case the NTK formulas of fully-connected neural network. We first define a fully-connected neural net formally. Let  $x \in \mathbb{R}^d$  be the input, and denote  $g^{(0)}(x) = x$  and  $d_0 = d$  for notational convenience. We define an  $L$ -hidden-layer fully-connected neural network recursively, for  $h = 1, 2, \dots, L$ :

$$f^{(h)}(x) = W^{(h)} g^{(h-1)}(x) \in \mathbb{R}^{d_h}, g^{(h)}(x) = \sqrt{\frac{c_\sigma}{d_h}} \sigma\left(f^{(h)}(x)\right) \in \mathbb{R}^{d_h} \quad (10.12)$$

where  $W^{(h)} \in \mathbb{R}^{d_h \times d_{h-1}}$  is the weight matrix in the  $h$ -th layer ( $h \in [L]$ ),  $\sigma : \mathbb{R} \rightarrow \mathbb{R}$  is a coordinate-wise activation function, and  $c_\sigma = \left(\mathbb{E}_{z \sim \mathcal{N}(0,1)} [\sigma z^2]\right)^{-1}$ . The last layer of the neural network is

$$\begin{aligned} f(w, x) &= f^{(L+1)}(x) = W^{(L+1)} \cdot g^{(L)}(x) \\ &= W^{(L+1)} \cdot \sqrt{\frac{c_\sigma}{d_L}} \sigma W^{(L)} \cdot \sqrt{\frac{c_\sigma}{d_{L-1}}} \sigma W^{(L-1)} \dots \sqrt{\frac{c_\sigma}{d_1}} \sigma W^{(1)} x, \end{aligned}$$

where  $W^{(L+1)} \in \mathbb{R}^{1 \times d_L}$  is the weights in the final layer, and  $w = \left(W^{(1)}, \dots, W^{(L+1)}\right)$  represents all the parameters in the network.

We initialize all the weights to be i.i.d.  $\mathcal{N}(0, 1)$  random variables, and consider the limit of large hidden widths:  $d_1, d_2, \dots, d_L \rightarrow \infty$ . The scaling factor  $\sqrt{c_\sigma/d_h}$  in Equation (10.12) ensures that the norm of  $g^{(h)}(x)$  for each  $h \in [L]$  is approximately preserved at initialization (see [DLL<sup>+</sup>18]). In particular, for ReLU activation, we have

$$\mathbb{E} \left[ \left\| g^{(h)}(x) \right\|_2^2 \right] = \|x\|_2^2 \quad (\forall h \in [L]).$$

Recall from Lemma 10.1.1 that we need to compute the value that  $\left\langle \frac{\partial f(w, x)}{\partial w}, \frac{\partial f(w, x')}{\partial w} \right\rangle$  converges to at random initialization in the infinite width limit. We can write the partial derivative with respect to a particular weight matrix  $W^{(h)}$  in a compact form:

$$\frac{\partial f(w, x)}{\partial W^{(h)}} = b^{(h)}(x) \cdot \left( g^{(h-1)}(x) \right)^\top, \quad h = 1, 2, \dots, L+1,$$

where

$$b^{(h)}(x) = \begin{cases} 1 \in \mathbb{R}, & h = L+1, \\ \sqrt{\frac{c_\sigma}{d_h}} D^{(h)}(x) \left( W^{(h+1)} \right)^\top b^{(h+1)}(x) \in \mathbb{R}^{d_h}, & h = 1, \dots, L, \end{cases} \quad (10.13)$$

$$D^{(h)}(x) = \text{diag} \left( \dot{\sigma} \left( f^{(h)}(x) \right) \right) \in \mathbb{R}^{d_h \times d_h}, \quad h = 1, \dots, L. \quad (10.14)$$

Then, for any two inputs  $x$  and  $x'$ , and any  $h \in [L+1]$ , we can compute

$$\begin{aligned} & \left\langle \frac{\partial f(w, x)}{\partial W^{(h)}}, \frac{\partial f(w, x')}{\partial W^{(h)}} \right\rangle \\ &= \left\langle b^{(h)}(x) \cdot \left( g^{(h-1)}(x) \right)^\top, b^{(h)}(x') \cdot \left( g^{(h-1)}(x') \right)^\top \right\rangle \\ &= \left\langle g^{(h-1)}(x), g^{(h-1)}(x') \right\rangle \cdot \left\langle b^{(h)}(x), b^{(h)}(x') \right\rangle. \end{aligned}$$

Note the first term  $\left\langle g^{(h-1)}(x), g^{(h-1)}(x') \right\rangle$  is the covariance between  $x$  and  $x'$  at the  $h$ -th layer. When the width goes to infinity,  $\left\langle g^{(h-1)}(x), g^{(h-1)}(x') \right\rangle$  will converge to a fix number, which we denote as  $\Sigma^{(h-1)}(x, x')$ . This covariance admits a recursive formula, for  $h \in [L]$ ,

$$\begin{aligned} \Sigma^{(0)}(x, x') &= x^\top x', \\ \Lambda^{(h)}(x, x') &= \begin{pmatrix} \Sigma^{(h-1)}(x, x) & \Sigma^{(h-1)}(x, x') \\ \Sigma^{(h-1)}(x', x) & \Sigma^{(h-1)}(x', x') \end{pmatrix} \in \mathbb{R}^{2 \times 2}, \quad (10.15) \\ \Sigma^{(h)}(x, x') &= c_\sigma \mathbb{E}_{(u, v) \sim \mathcal{N}(0, \Lambda^{(h)})} [\sigma(u) \sigma(v)]. \end{aligned}$$

Now we proceed to derive this formula. The intuition is that  $\left[ f^{(h+1)}(x) \right]_i = \sum_{j=1}^{d_h} \left[ W^{(h+1)} \right]_{i,j} \left[ g^{(h)}(x) \right]_j$  is a centered Gaussian

process conditioned on  $f^{(h)}$  ( $\forall i \in [d_{h+1}]$ ), with covariance

$$\begin{aligned} & \mathbb{E} \left[ \left[ f^{(h+1)}(x) \right]_i \cdot \left[ f^{(h+1)}(x') \right]_i \mid f^{(h)} \right] \\ &= \left\langle g^{(h)}(x), g^{(h)}(x') \right\rangle \\ &= \frac{c_\sigma}{d_h} \sum_{j=1}^{d_h} \sigma \left( \left[ f^{(h)}(x) \right]_j \right) \sigma \left( \left[ f^{(h)}(x') \right]_j \right), \end{aligned} \quad (10.16)$$

which converges to  $\Sigma^{(h)}(x, x')$  as  $d_h \rightarrow \infty$  given that each  $\left[ f^{(h)} \right]_j$  is a centered Gaussian process with covariance  $\Sigma^{(h-1)}$ . This yields the inductive definition in Equation (10.15).

Next we deal with the second term  $\left\langle b^{(h)}(x), b^{(h)}(x') \right\rangle$ . From Equation (10.13) we get

$$\begin{aligned} & \left\langle b^{(h)}(x), b^{(h)}(x') \right\rangle \\ &= \left\langle \sqrt{\frac{c_\sigma}{d_h}} D^{(h)}(x) \left( W^{(h+1)} \right)^\top b^{(h+1)}(x), \sqrt{\frac{c_\sigma}{d_h}} D^{(h)}(x') \left( W^{(h+1)} \right)^\top b^{(h+1)}(x') \right\rangle. \end{aligned} \quad (10.17)$$

Although  $W^{(h+1)}$  and  $b_{h+1}(x)$  are dependent, the Gaussian initialization of  $W^{(h+1)}$  allows us to replace  $W^{(h+1)}$  with a fresh new sample  $\tilde{W}^{(h+1)}$  without changing its limit: (See [ADH<sup>+</sup>19c] for the precise proof.)

$$\begin{aligned} & \left\langle \sqrt{\frac{c_\sigma}{d_h}} D^{(h)}(x) \left( W^{(h+1)} \right)^\top b^{(h+1)}(x), \sqrt{\frac{c_\sigma}{d_h}} D^{(h)}(x') \left( W^{(h+1)} \right)^\top b^{(h+1)}(x') \right\rangle \\ &\approx \left\langle \sqrt{\frac{c_\sigma}{d_h}} D^{(h)}(x) \left( \tilde{W}^{(h+1)} \right)^\top b^{(h+1)}(x), \sqrt{\frac{c_\sigma}{d_h}} D^{(h)}(x') \left( \tilde{W}^{(h+1)} \right)^\top b^{(h+1)}(x') \right\rangle \\ &\rightarrow \frac{c_\sigma}{d_h} \text{tr} D^{(h)}(x) D^{(h)}(x') \left\langle b^{(h+1)}(x), b^{(h+1)}(x') \right\rangle \\ &\rightarrow \dot{\Sigma}^{(h)}(x, x') \left\langle b^{(h+1)}(x), b^{(h+1)}(x') \right\rangle. \end{aligned}$$

Applying this approximation inductively in Equation (10.17), we get

$$\left\langle b^{(h)}(x), b^{(h)}(x') \right\rangle \rightarrow \prod_{h'=h}^L \dot{\Sigma}^{(h')}(x, x').$$

Finally, since  $\left\langle \frac{\partial f(w, x)}{\partial w}, \frac{\partial f(w, x')}{\partial w} \right\rangle = \sum_{h=1}^{L+1} \left\langle \frac{\partial f(w, x)}{\partial W^{(h)}}, \frac{\partial f(w, x')}{\partial W^{(h)}} \right\rangle$ , we obtain the final NTK expression for the fully-connected neural network:

$$\Theta^{(L)}(x, x') = \sum_{h=1}^{L+1} \left( \Sigma^{(h-1)}(x, x') \cdot \prod_{h'=h}^{L+1} \dot{\Sigma}^{(h')}(x, x') \right).$$

### 10.5 NTK in Practice

Up to now we have shown an ultra-wide neural network with certain initialization scheme and trained by gradient flow correspond to a



kernel with a particular kernel function. A natural question is: *why don't we use this kernel classifier directly?*

A recent line of work showed that NTKs can be empirically useful, especially on small to medium scale datasets. Arora et al. [ADL<sup>+</sup>19] tested the NTK classifier on 90 small to medium scale datasets from UCI database.<sup>5</sup> They found NTK can beat neural networks, other kernels like Gaussian and the best previous classifier, random forest under various metrics, including average rank, average accuracy, etc. This suggests the NTK classifier should belong in any list of off-the-shelf machine learning methods.

<sup>5</sup> <https://archive.ics.uci.edu/ml/datasets.php>

For every neural network architecture, one can derive a corresponding kernel function. Du et al. [DHS<sup>+</sup>19] derived graph NTK (GNTK) for graph classification tasks. On various social network and bioinformatics datasets, GNTK can outperform graph neural networks.

Similarly, Arora et al. [ADH<sup>+</sup>19c] derived convolutional NTK (CNTK) formula that corresponds to convolutional neural networks. For image classification task, in small-scale data and low-shot settings, CNTKs can be quite strong [ADL<sup>+</sup>19]. However, for large scale data, Arora et al. [ADH<sup>+</sup>19c] found there is still a performance gap between CNTK and CNN. It is an open problem to explain this phenomenon theoretically. This may need to go beyond the NTK framework.

## 10.6 Exercises

1. NTK formula for ReLU activation function: prove

$$\mathbb{E}_{w \sim \mathcal{N}(0, I)} \left[ (\dot{\sigma}(w^\top x) \dot{\sigma}(w^\top x')) \right] = \frac{\pi - \arccos \left( \frac{x^\top x'}{\|x\|_2 \|x'\|_2} \right)}{2\pi}.$$

2. Prove Equation (10.11). (Hint: The generalization bound depends upon the norm of the difference between the final  $W$  matrix and the initial  $W$ , which you can upper bound using a sum/integral similar to the analysis of kernel regression in Chapter 3.)
3. Why NTK can learn a linear function: in this question, you are asked to prove that NTK can learn a linear function (the technique here can be used to show NTK can learn other functions listed in Section 10.3.1). In this question, we assume we have  $n$  data points,  $\{(x_i, y_i)\}_{i=1}^n$  where every input  $x_i$  has norm 1 and  $y_i = \beta^\top x_i$  for some  $\beta$ . Each entry of neural tangent kernel matrix  $H^*$  (induced by the two-layer ReLU neural network) is  $H_{ij}^* = \frac{\pi - \arccos \left( \frac{x_i^\top x_j}{\|x_i\|_2 \|x_j\|_2} \right)}{2\pi}$ .

- (a) Taylor Expansion: use the Taylor expansion of  $\arccos(z) = \frac{\pi}{2} - z - \sum_{\ell=1}^{\infty} \frac{(2\ell-3)!!z}{(2\ell-2)!!}$  to show that  $H^*$  admits the following form

$$H^* = \frac{K}{4} + \sum_{\ell=1}^{\infty} \frac{1}{2\pi} \frac{(2\ell-3)!!}{(2\ell-2)!!} \cdot \frac{K^{\circ 2\ell}}{2\ell-1}$$

where  $K_{ij} = x_i^\top x_j$  and  $K_{ij}^\ell = (x_i^\top x_j)^\ell$ .

- (b) Assume  $H^*$  and  $K$  are invertible. Show  $(H^*)^{-1} \preceq 4K^{-1}$ .
- (c) Show  $\text{tr}(H^*) \leq n$ .
- (d) Using the assumption  $y_i = x_i^\top \beta$  to show

$$\frac{\sqrt{2y^\top (H^*)^{-1} y \cdot \text{tr}(H^*)}}{n} \leq \frac{2\sqrt{2} \|\beta\|_2}{\sqrt{n}}.$$

## *Inductive Biases due to Algorithmic Regularization*

Many successful modern machine learning systems based on deep neural networks are over-parametrized, i.e., the number of parameters is typically much larger than the sample size. In other words, there exist (infinitely) many (approximate) minimizers of the empirical risk, many of which would not generalize well on the unseen data. For learning to succeed then, it is crucial to bias the learning algorithm towards “simpler” hypotheses by trading off empirical loss with a certain complexity term that ensures that empirical and population risks are close. Several explicit regularization strategies have been used in practice to help these systems generalize, including  $\ell_1$  and  $\ell_2$  regularization of the parameters [NH92].

Besides explicit regularization techniques, practitioners have used a spectrum of algorithmic approaches to improve the generalization ability of over-parametrized models. This includes early stopping of back-propagation [CLG01], batch normalization [IS15], dropout [SHK<sup>+</sup>14], and more<sup>1</sup>. While these heuristics have enjoyed tremendous success in training deep networks, a theoretical understanding of how these heuristics provide regularization in deep learning remains somewhat limited.

In this chapter, we investigate regularization due to Dropout, an algorithmic heuristic recently proposed by [SHK<sup>+</sup>14]. The basic idea when training a neural network using dropout, is that during a forward pass, we randomly drop neurons in the neural network, independently and identically according to a Bernoulli distribution. Specifically, at each round of the back-propagation algorithm, for each neuron, independently, with probability  $p$  we “drop” the neuron, so it does not participate in making a prediction for the given data point, and with probability  $1 - p$  we retain that neuron<sup>2</sup>.

Deep learning is a field where key innovations have been driven by practitioners, with several techniques motivated by drawing insights from other fields. For instance, Dropout was introduced as a way of breaking up “co-adaptation” among neurons, drawing

<sup>1</sup> We refer the reader to [KGC17] for an excellent exposition of over 50 of such proposals.

<sup>2</sup> The parameter  $p$  is treated as a hyper-parameter which we typically tune for based on a validation set.

insights from the success of the sexual reproduction model in the evolution of advanced organisms. Another motivation that was cited by [SHK<sup>+</sup>14] was in terms of “balancing networks”. Despite several theoretical works aimed at explaining Dropout<sup>3</sup>, it remains unclear what kind of regularization does Dropout provide or what kinds of networks does Dropout prefer and how that helps with generalization. In this chapter, we work towards that goal by instantiating explicit forms of regularizers due to Dropout and how they provide capacity control in various machine learning including linear regression (Section 11.4), matrix sensing (Section 11.1.1), matrix completion (Section 11.1.2), and deep learning (Section 11.2).

3

### 11.1 Matrix Sensing

We begin with understanding dropout for matrix sensing, a problem which arguably is an important instance of a matrix learning problem with lots of applications, and is well understood from a theoretical perspective. Here is the problem setup.

Let  $M_* \in \mathbb{R}^{d_2 \times d_0}$  be a matrix with rank  $r_* := \text{Rank}(M_*)$ . Let  $A^{(1)}, \dots, A^{(n)}$  be a set of measurement matrices of the same size as  $M_*$ . The goal of matrix sensing is to recover the matrix  $M_*$  from  $n$  observations of the form  $y_i = \langle M_*, A^{(i)} \rangle$  such that  $n \ll d_2 d_0$ . The learning algorithm we consider is empirical risk minimization, and we choose to represent the parameter matrix  $M \in \mathbb{R}^{d_2 \times d_0}$  in terms of product of its factors  $U \in \mathbb{R}^{d_2 \times d_1}, V \in \mathbb{R}^{d_0 \times d_1}$ :

$$\min_{U \in \mathbb{R}^{d_2 \times d_1}, V \in \mathbb{R}^{d_0 \times d_1}} \widehat{L}(U, V) := \frac{1}{n} \sum_{i=1}^n (y_i - \langle UV^T, A^{(i)} \rangle)^2. \quad (11.1)$$

When  $d_1 \gg r_*$ , there exist many “bad” empirical minimizers, i.e., those with a large true risk. However, recently, [LMZ18] showed that under restricted isometry property, despite the existence of such poor ERM solutions, gradient descent with proper initialization is *implicitly* biased towards finding solutions with minimum nuclear norm – this is an important result which was first conjectured and empirically verified by [GWB<sup>+</sup>17].

We propose solving the ERM problem (11.1) with algorithmic regularization due to dropout, where at training time, the corresponding columns of  $U$  and  $V$  are dropped independently and identically according to a Bernoulli random variable. As opposed to the *implicit* effect of gradient descent, this dropout heuristic *explicitly* regularizes the empirical objective. It is then natural to ask, in the case of matrix sensing, if dropout also biases the ERM towards certain low norm solutions. To answer this question, we begin with the observation that dropout can be viewed as an instance of SGD on the following

objective:

$$\widehat{L}_{\text{drop}}(\mathbf{U}, \mathbf{V}) = \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{\mathbf{B}} (y_i - \langle \mathbf{UBV}^{\top}, \mathbf{A}^{(i)} \rangle)^2, \quad (11.2)$$

where  $\mathbf{B} \in \mathbb{R}^{d_1 \times d_1}$  is a diagonal matrix whose diagonal elements are Bernoulli random variables distributed as  $B_{jj} \sim \frac{1}{1-p} \text{Ber}(1-p)$ , for  $j \in [d_1]$ . In this case, we can show that for any  $p \in [0, 1)$ :

$$\widehat{L}_{\text{drop}}(\mathbf{U}, \mathbf{V}) = \widehat{L}(\mathbf{U}, \mathbf{V}) + \frac{p}{1-p} \widehat{R}(\mathbf{U}, \mathbf{V}), \quad (11.3)$$

where

$$\widehat{R}(\mathbf{U}, \mathbf{V}) = \sum_{j=1}^{d_1} \frac{1}{n} \sum_{i=1}^n (\mathbf{u}_j^{\top} \mathbf{A}^{(i)} \mathbf{v}_j)^2 \quad (11.4)$$

is a data-dependent term that captures the *explicit* regularizer due to dropout.

*Proof.* Consider one of the summands in the Dropout objective in Equation 11.2. Then, we can write

$$\begin{aligned} \mathbb{E}_{\mathbf{B}} [(y_i - \langle \mathbf{UBV}^{\top}, \mathbf{A}^{(i)} \rangle)^2] &= \left( \mathbb{E}_{\mathbf{B}} [y_i - \langle \mathbf{UBV}^{\top}, \mathbf{A}^{(i)} \rangle] \right)^2 \\ &\quad + \text{Var}(y_i - \langle \mathbf{UBV}^{\top}, \mathbf{A}^{(i)} \rangle) \end{aligned} \quad (11.5)$$

For Bernoulli random variable  $B_{jj}$ , we have that  $\mathbb{E}[B_{jj}] = 1$  and  $\text{Var}(B_{jj}) = \frac{p}{1-p}$ . Thus, the first term on right hand side is equal to  $(y_i - \langle \mathbf{UV}^{\top}, \mathbf{A}^{(i)} \rangle)^2$ . For the second term we have

$$\begin{aligned} \text{Var}(y_i - \langle \mathbf{UBV}^{\top}, \mathbf{A}^{(i)} \rangle) &= \text{Var}(\langle \mathbf{UBV}^{\top}, \mathbf{A}^{(i)} \rangle) \\ &= \text{Var}(\langle \mathbf{B}, \mathbf{U}^{\top} \mathbf{A}^{(i)} \mathbf{V} \rangle) \\ &= \text{Var}\left(\sum_{j=1}^{d_1} B_{jj} \mathbf{u}_j^{\top} \mathbf{A}^{(i)} \mathbf{v}_j\right) \\ &= \sum_{j=1}^{d_1} (\mathbf{u}_j^{\top} \mathbf{A}^{(i)} \mathbf{v}_j)^2 \text{Var}(B_{jj}) \\ &= \frac{p}{1-p} \sum_{j=1}^{d_1} (\mathbf{u}_j^{\top} \mathbf{A}^{(i)} \mathbf{v}_j)^2 \end{aligned}$$

Using the facts above in Equation (11.2), we get

$$\begin{aligned} \widehat{L}_{\text{drop}}(\mathbf{U}, \mathbf{V}) &= \frac{1}{n} \sum_{i=1}^n (y_i - \langle \mathbf{UV}^{\top}, \mathbf{A}^{(i)} \rangle)^2 + \frac{1}{n} \sum_{i=1}^n \frac{p}{1-p} \sum_{j=1}^{d_1} (\mathbf{u}_j^{\top} \mathbf{A}^{(i)} \mathbf{v}_j)^2 \\ &= \widehat{L}(\mathbf{U}, \mathbf{V}) + \frac{p}{1-p} \widehat{R}(\mathbf{U}, \mathbf{V}). \end{aligned}$$

which completes the proof.  $\square$

Provided that the sample size  $n$  is large enough, the *explicit* regularizer on a given sample behaves much like its expected value with respect to the underlying data distribution<sup>4</sup>. Further, given that we seek a minimum of  $\widehat{L}_{\text{drop}}$ , it suffices to consider the factors with the minimal value of the regularizer among all that yield the same empirical loss. This motivates studying the the following distribution-dependent *induced* regularizer:

$$\Theta(\mathbf{M}) := \min_{\mathbf{U}\mathbf{V}^\top = \mathbf{M}} R(\mathbf{U}, \mathbf{V}), \text{ where } R(\mathbf{U}, \mathbf{V}) := \mathbb{E}_\mathbf{A}[\widehat{R}(\mathbf{U}, \mathbf{V})].$$

Next, we consider two two important examples of random sensing matrices.

### 11.1.1 Gaussian Sensing Matrices

We assume that the entries of the sensing matrices are independently and identically distributed as standard Gaussian, i.e.,  $A_{k\ell}^{(i)} \sim \mathcal{N}(0, 1)$ . For Gaussian sensing matrices, we show that the induced regularizer due to Dropout provides nuclear-norm regularization. Formally, we show that

$$\Theta(\mathbf{M}) = \frac{1}{d_1} \|\mathbf{M}\|_*^2. \quad (11.6)$$

*Proof.* We recall the general form for the dropout regularizer for the matrix sensing problem in Equation 11.4, and take expectation with respect to the distribution on the sensing matrices. Then, for any pair of factors  $(\mathbf{U}, \mathbf{V})$ , it holds that the expected regularizer is given as follows.

$$\begin{aligned} R(\mathbf{U}, \mathbf{V}) &= \sum_{j=1}^{d_1} \mathbb{E}(\mathbf{u}_j^\top \mathbf{A} \mathbf{v}_j)^2 \\ &= \sum_{j=1}^{d_1} \mathbb{E} \left( \sum_{k=1}^{d_2} \sum_{\ell=1}^{d_0} \mathbf{U}_{kj} \mathbf{A}_{k\ell} \mathbf{V}_{\ell j} \right)^2 \\ &= \sum_{j=1}^{d_1} \sum_{k,k'=1}^{d_2} \sum_{\ell,\ell'=1}^{d_0} \mathbf{U}_{kj} \mathbf{U}_{k'j} \mathbf{V}_{\ell j} \mathbf{V}_{\ell'j} \mathbb{E}[\mathbf{A}_{k\ell} \mathbf{A}_{k'\ell'}] \\ &= \sum_{j=1}^{d_1} \sum_{k=1}^{d_2} \sum_{\ell=1}^{d_0} \mathbf{U}_{kj}^2 \mathbf{V}_{\ell j}^2 \mathbb{E}[\mathbf{A}_{k\ell}^2] \\ &= \sum_{j=1}^{d_1} \sum_{k=1}^{d_2} \sum_{\ell=1}^{d_0} \mathbf{U}_{kj}^2 \mathbf{V}_{\ell j}^2 \\ &= \sum_{j=1}^{d_1} \|\mathbf{u}_j\|^2 \|\mathbf{v}_j\|^2 \end{aligned}$$

<sup>4</sup> Under mild assumptions, we can formally show that the dropout regularizer is well concentrated around its mean

Now, using the Cauchy-Schwartz inequality, we can bound the expected regularizer as

$$\begin{aligned}
 R(\mathbf{U}, \mathbf{V}) &\geq \frac{1}{d_1} \left( \sum_{i=1}^{d_1} \|\mathbf{u}_i\| \|\mathbf{v}_i\| \right)^2 \\
 &= \frac{1}{d_1} \left( \sum_{i=1}^{d_1} \|\mathbf{u}_i \mathbf{v}_i^\top\|_* \right)^2 \\
 &\geq \frac{1}{d_1} \left( \left\| \sum_{i=1}^{d_1} \mathbf{u}_i \mathbf{v}_i^\top \right\|_* \right)^2 = \frac{1}{d_1} \|\mathbf{U} \mathbf{V}^\top\|_*^2
 \end{aligned}$$

where the equality follows because for any pair of vectors  $\mathbf{a}, \mathbf{b}$ , it holds that  $\|\mathbf{a} \mathbf{b}^\top\|_* = \|\mathbf{a} \mathbf{b}^\top\|_F = \|\mathbf{a}\| \|\mathbf{b}\|$ , and the last inequality is due to triangle inequality.

Next, we need the following key result from [MAV18].

**Theorem 11.1.1.** For any pair of matrices  $\mathbf{U} \in \mathbb{R}^{d_2 \times d_1}, \mathbf{V} \in \mathbb{R}^{d_0 \times d_1}$ , there exists a rotation matrix  $\mathbf{Q} \in \text{SO}(d_1)$  such that matrices  $\tilde{\mathbf{U}} := \mathbf{U} \mathbf{Q}, \tilde{\mathbf{V}} := \mathbf{V} \mathbf{Q}$  satisfy  $\|\tilde{\mathbf{u}}_i\| \|\tilde{\mathbf{v}}_i\| = \frac{1}{d_1} \|\mathbf{U} \mathbf{V}^\top\|_*$ , for all  $i \in [d_1]$ .

Using Theorem 11.1.1 on  $(\mathbf{U}, \mathbf{V})$ , the expected dropout regularizer at  $\mathbf{U} \mathbf{Q}, \mathbf{V} \mathbf{Q}$  is given as

$$\begin{aligned}
 R(\mathbf{U} \mathbf{Q}, \mathbf{V} \mathbf{Q}) &= \sum_{i=1}^{d_1} \|\mathbf{U} \mathbf{q}_i\|^2 \|\mathbf{V} \mathbf{q}_i\|^2 \\
 &= \sum_{i=1}^{d_1} \frac{1}{d_1^2} \|\mathbf{U} \mathbf{V}^\top\|_*^2 \\
 &= \frac{1}{d_1} \|\mathbf{U} \mathbf{V}^\top\|_*^2 \\
 &\leq \Theta(\mathbf{U} \mathbf{V}^\top)
 \end{aligned}$$

which completes the proof.  $\square$

For completeness we provide a proof of Theorem 11.1.1.

*Proof.* Define  $\mathbf{M} := \mathbf{U} \mathbf{V}^\top$ . Let  $\mathbf{M} = \mathbf{W} \mathbf{\Sigma} \mathbf{Y}^\top$  be compact SVD of  $\mathbf{M}$ . Define  $\hat{\mathbf{U}} := \mathbf{W} \mathbf{\Sigma}^{1/2}$  and  $\hat{\mathbf{V}} := \mathbf{Y} \mathbf{\Sigma}^{1/2}$ . Let  $\mathbf{G}_\mathbf{U} = \hat{\mathbf{U}}^\top \hat{\mathbf{U}}$  and  $\mathbf{G}_\mathbf{V} = \hat{\mathbf{V}}^\top \hat{\mathbf{V}}$  be respective Gram matrices. Observe that  $\mathbf{G}_\mathbf{U} = \mathbf{G}_\mathbf{V} = \mathbf{\Sigma}$ . We will show that there exists a rotation  $\mathbf{Q}$  such that for  $\tilde{\mathbf{U}} = \hat{\mathbf{U}} \mathbf{Q}, \tilde{\mathbf{V}} = \hat{\mathbf{V}} \mathbf{Q}$ , it holds that

$$\|\tilde{\mathbf{u}}_j\|^2 = \frac{1}{d_1} \|\tilde{\mathbf{U}}\|_F^2 = \frac{1}{d_1} \text{Tr}(\tilde{\mathbf{U}}^\top \tilde{\mathbf{U}}) = \frac{1}{d_1} \text{Tr}(\mathbf{\Sigma}) = \frac{1}{d_1} \|\mathbf{M}\|_*$$

and

$$\|\tilde{\mathbf{v}}_j\|^2 = \frac{1}{d_1} \|\tilde{\mathbf{V}}\|_F^2 = \frac{1}{d_1} \text{Tr}(\tilde{\mathbf{V}}^\top \tilde{\mathbf{V}}) = \frac{1}{d_1} \text{Tr}(\mathbf{\Sigma}) = \frac{1}{d_1} \|\mathbf{M}\|_*$$

Consequently, it holds that  $\|\tilde{\mathbf{u}}_i\| \|\tilde{\mathbf{v}}_i\| = \frac{1}{d_1} \|\mathbf{M}\|_*$ .

All that remains is to give a construction of matrix  $\mathbf{Q}$ . We note that a rotation matrix  $\mathbf{Q}$  satisfies the desired properties above if and only if all diagonal elements of  $\mathbf{Q}^\top \mathbf{G}_U \mathbf{Q}$  are equal<sup>5</sup>, and equal to  $\frac{\text{Tr} \mathbf{G}_U}{d_1}$ . The key idea is that for the trace zero matrix  $\mathbf{G}_1 := \mathbf{G}_U - \frac{\text{Tr} \mathbf{G}_U}{d_1} \mathbf{I}_{d_1}$ , if  $\mathbf{G}_1 = \sum_{i=1}^r \lambda_i \mathbf{e}_i \mathbf{e}_i^\top$  is an eigendecomposition of  $\mathbf{G}_1$ , then for the average of the eigenvectors, i.e. for  $\mathbf{w}_{11} = \frac{1}{\sqrt{r}} \sum_{i=1}^r \mathbf{e}_i$ , it holds that  $\mathbf{w}_{11}^\top \mathbf{G}_1 \mathbf{w}_{11} = 0$ . We use this property recursively to exhibit an orthogonal transformation  $\mathbf{Q}$ , such that  $\mathbf{Q}^\top \mathbf{G}_1 \mathbf{Q}$  is zero on its diagonal.

<sup>5</sup> since  $(\mathbf{Q}^\top \mathbf{G}_U \mathbf{Q})_{jj} = \|\tilde{\mathbf{u}}_j\|^2$

To verify the claim, first notice that  $\mathbf{w}_{11}$  is unit norm

$$\|\mathbf{w}_{11}\|^2 = \left\| \frac{1}{\sqrt{r}} \sum_{i=1}^r \mathbf{e}_i \right\|^2 = \frac{1}{r} \sum_{i=1}^r \|\mathbf{e}_i\|^2 = 1.$$

Further, it is easy to see that

$$\mathbf{w}_{11}^\top \mathbf{G} \mathbf{w}_{11} = \frac{1}{r} \sum_{i,j=1}^r \mathbf{e}_i \mathbf{G} \mathbf{e}_j = \frac{1}{r} \sum_{i,j=1}^r \lambda_j \mathbf{e}_i^\top \mathbf{e}_j = \frac{1}{r} \sum_{i=1}^r \lambda_i = 0.$$

Complete  $\mathbf{W}_1 := [\mathbf{w}_{11}, \mathbf{w}_{12}, \dots, \mathbf{w}_{1d}]$  be such that  $\mathbf{W}_1^\top \mathbf{W}_1 = \mathbf{W}_1 \mathbf{W}_1^\top = \mathbf{I}_d$ . Observe that  $\mathbf{W}_1^\top \mathbf{G}_1 \mathbf{W}_1$  has zero on its first diagonal elements

$$\mathbf{W}_1^\top \mathbf{G}_1 \mathbf{W}_1 = \begin{bmatrix} 0 & \mathbf{b}_1^\top \\ \mathbf{b}_1 & \mathbf{G}_2 \end{bmatrix}$$

The principal submatrix  $\mathbf{G}_2$  also has a zero trace. With a similar argument, let  $\mathbf{w}_{22} \in \mathbb{R}^{d-1}$  be such that  $\|\mathbf{w}_{22}\| = 1$  and  $\mathbf{w}_{22}^\top \mathbf{G}_2 \mathbf{w}_{22} = 0$  and define  $\mathbf{W}_2 = \begin{bmatrix} 1 & 0 & 0 & \dots & 0 \\ 0 & \mathbf{w}_{22} & \mathbf{w}_{23} & \dots & \mathbf{w}_{2d} \end{bmatrix} \in \mathbb{R}^{d \times d}$  such that  $\mathbf{W}_2^\top \mathbf{W}_2 = \mathbf{W}_2 \mathbf{W}_2^\top = \mathbf{I}_d$ , and observe that

$$(\mathbf{W}_1 \mathbf{W}_2)^\top \mathbf{G}_1 (\mathbf{W}_1 \mathbf{W}_2) = \begin{bmatrix} 0 & \cdot & \dots \\ \cdot & 0 & \dots \\ \vdots & \vdots & \mathbf{G}_3 \end{bmatrix}.$$

This procedure can be applied recursively so that for the matrix  $\mathbf{Q} = \mathbf{W}_1 \mathbf{W}_2 \dots \mathbf{W}_d$  we have

$$\mathbf{Q}^\top \mathbf{G}_1 \mathbf{Q} = \begin{bmatrix} 0 & \cdot & \dots & \cdot \\ \cdot & 0 & \dots & \cdot \\ \vdots & \vdots & \ddots & \vdots \\ \cdot & \cdot & \cdot & 0 \end{bmatrix},$$

so that  $\text{Tr}(\tilde{\mathbf{U}} \tilde{\mathbf{U}}^\top) = \text{Tr}(\mathbf{Q}^\top \mathbf{G}_U \mathbf{Q}) = \text{Tr}(\Sigma) = \text{Tr}(\mathbf{Q}^\top \mathbf{G}_V \mathbf{Q}) = \text{Tr}(\tilde{\mathbf{V}}^\top \tilde{\mathbf{V}})$ .  $\square$



### 11.1.2 Matrix Completion

Next, we consider the problem of matrix completion which can be formulated as a special case of matrix sensing with sensing matrices that random indicator matrices. Formally, we assume that for all  $j \in [n]$ , let  $A^{(j)}$  be an indicator matrix whose  $(i, k)$ -th element is selected randomly with probability  $p(i)q(k)$ , where  $p(i)$  and  $q(k)$  denote the probability of choosing the  $i$ -th row and the  $j$ -th column, respectively.

We will show next that in this setup Dropout induces the *weighted trace-norm* studied by [SS10] and [FSSS11]. Formally, we show that

$$\Theta(M) = \frac{1}{d_1} \|\text{diag}(\sqrt{p})UV^\top \text{diag}(\sqrt{q})\|_*^2. \quad (11.7)$$

*Proof.* For any pair of factors  $(U, V)$  it holds that

$$\begin{aligned} R(U, V) &= \sum_{j=1}^{d_1} \mathbb{E}(\mathbf{u}_j^\top A \mathbf{v}_j)^2 \\ &= \sum_{j=1}^{d_1} \sum_{k=1}^{d_2} \sum_{\ell=1}^{d_0} p(k)q(\ell) (\mathbf{u}_j^\top \mathbf{e}_k \mathbf{e}_\ell^\top \mathbf{v}_j)^2 \\ &= \sum_{j=1}^{d_1} \sum_{k=1}^{d_2} \sum_{\ell=1}^{d_0} p(k)q(\ell) U(k, j)^2 V(\ell, j)^2 \\ &= \sum_{j=1}^{d_1} \|\sqrt{\text{diag}(p)} \mathbf{u}_j\|^2 \|\sqrt{\text{diag}(q)} \mathbf{v}_j\|^2 \\ &\geq \frac{1}{d_1} \left( \sum_{j=1}^{d_1} \|\sqrt{\text{diag}(p)} \mathbf{u}_j\| \|\sqrt{\text{diag}(q)} \mathbf{v}_j\| \right)^2 \\ &= \frac{1}{d_1} \left( \sum_{j=1}^{d_1} \|\sqrt{\text{diag}(p)} \mathbf{u}_j \mathbf{v}_j^\top \sqrt{\text{diag}(q)}\|_* \right)^2 \\ &\geq \frac{1}{d_1} \left( \|\sqrt{\text{diag}(p)} \sum_{j=1}^{d_1} \mathbf{u}_j \mathbf{v}_j^\top \sqrt{\text{diag}(q)}\|_* \right)^2 \\ &= \frac{1}{d_1} \|\sqrt{\text{diag}(p)} UV^\top \sqrt{\text{diag}(q)}\|_*^2 \end{aligned}$$

where the first inequality is due to Cauchy-Schwartz and the second inequality follows from the triangle inequality. The equality right after the first inequality follows from the fact that for any two vectors  $\mathbf{a}, \mathbf{b}$ ,  $\|\mathbf{a}\mathbf{b}^\top\|_* = \|\mathbf{a}\mathbf{b}^\top\|_F = \|\mathbf{a}\| \|\mathbf{b}\|$ . Since the inequalities hold for any  $U, V$ , it implies that

$$\Theta(UV^\top) \geq \frac{1}{d_1} \|\sqrt{\text{diag}(p)} UV^\top \sqrt{\text{diag}(q)}\|_*^2.$$

Applying Theorem 11.1.1 on  $(\sqrt{\text{diag}(p)}U, \sqrt{\text{diag}(q)}V)$ , there

exists a rotation matrix  $Q$  such that

$$\|\sqrt{\text{diag}(p)}Uq_j\| \|\sqrt{\text{diag}(q)}Vq_j\| = \frac{1}{d_1} \|\sqrt{\text{diag}(p)}UV^\top \sqrt{\text{diag}(q)}\|_*.$$

We evaluate the expected dropout regularizer at  $UQ, VQ$ :

$$\begin{aligned} R(UQ, VQ) &= \sum_{j=1}^{d_1} \|\sqrt{\text{diag}(p)}Uq_j\|^2 \|\sqrt{\text{diag}(q)}Vq_j\|^2 \\ &= \sum_{j=1}^{d_1} \frac{1}{d_1^2} \|\sqrt{\text{diag}(p)}UV^\top \sqrt{\text{diag}(q)}\|_*^2 \\ &= \frac{1}{d_1} \|\sqrt{\text{diag}(p)}UV^\top \sqrt{\text{diag}(q)}\|_*^2 \\ &\leq \Theta(UV^\top) \end{aligned}$$

which completes the proof.  $\square$

The results above are interesting because they connect Dropout, an algorithmic heuristic in deep learning, to strong complexity measures that are empirically effective as well as theoretically well understood. To illustrate, here we give a generalization bound for matrix completion with dropout in terms of the value of the *explicit* regularizer at the minimum of the empirical problem.

**Theorem 11.1.2.** Without loss of generality, assume that  $d_2 \geq d_0$  and  $\|M_*\| \leq 1$ . Furthermore, assume that  $\min_{i,j} p(i)q(j) \geq \frac{\log(d_2)}{n\sqrt{d_2d_0}}$ . Let  $(U, V)$  be a minimizer of the dropout ERM objective in equation (11.2), and assume that  $\max_i \|U(i, \cdot)\|^2 \leq \gamma$ ,  $\max_i \|V(i, \cdot)\|^2 \leq \gamma$ . Let  $\alpha$  be such that  $\hat{R}(U, V) \leq \alpha/d_1$ . Then, for any  $\delta \in (0, 1)$ , the following generalization bounds holds with probability at least  $1 - 2\delta$  over a sample of size  $n$ :

$$L(U, V) \leq \hat{L}(U, V) + C(1 + \gamma) \sqrt{\frac{\alpha d_2 \log(d_2)}{n}} + C'(1 + \gamma^2) \sqrt{\frac{\log(2/\delta)}{2n}}$$

as long as  $n = \Omega((d_1\gamma^2/\alpha)^2 \log(2/\delta))$ , where  $C, C'$  are some absolute constants.

The proof of Theorem 11.1.2 follows from standard generalization bounds for  $\ell_2$  loss [MRT18] based on the Rademacher complexity [BMo2] of the class of functions with weighted trace-norm bounded by  $\sqrt{\alpha}$ , i.e.  $\mathcal{M}_\alpha := \{M : \|\text{diag}(\sqrt{p})M\text{diag}(\sqrt{q})\|_*^2 \leq \alpha\}$ . A bound on the Rademacher complexity of this class was established by [FSS11]. The technicalities here include showing that the explicit regularizer is well concentrated around its expected value, as well as deriving a bound on the supremum of the predictions. A few remarks are in order.

We require that the sampling distributions be non-degenerate, as specified by the condition  $\min_{i,j} p(i)q(j) \geq \frac{\log(d_2)}{n\sqrt{d_2d_0}}$ . This is a natural requirement for bounding the Rademacher complexity of  $\mathcal{M}_\alpha$ , as discussed in [FSSS11].

We note that for large enough sample size,  $\widehat{R}(U, V) \approx R(U, V) \approx \Theta(UV^\top) = \frac{1}{d_1} \|\text{diag}(\sqrt{p})UV^\top \text{diag}(\sqrt{q})\|_*^2$ , where the second approximation is due the fact that the pair  $(U, V)$  is a minimizer. That is, compared to the weighted trace-norm, the value of the explicit regularizer at the minimizer roughly scales as  $1/d_1$ . Hence the assumption  $\widehat{R}(U, V) \leq \alpha/d_1$  in the statement of the corollary.

In practice, for models that are trained with dropout, the training error  $\widehat{L}(U, V)$  is negligible. Moreover, given that the sample size is large enough, the third term can be made arbitrarily small. Having said that, the second term, which is  $\widetilde{O}(\gamma\sqrt{\alpha d_2/n})$ , dominates the right hand side of generalization error bound in Theorem 11.1.2.

The assumption  $\max_i \|U(i, \cdot)\|^2 \leq \gamma$ ,  $\max_i \|V(i, \cdot)\|^2 \leq \gamma$  is motivated by the practice of deep learning; such *max-norm* constraints are typically used with dropout, where the norm of the vector of incoming weights at each hidden unit is constrained to be bound by a constant [SHK<sup>+</sup>14]. In this case, if a dropout update violates this constraint, the weights of the hidden unit are projected back to the constraint norm ball. In proofs, we need this assumption to give a concentration bound for the empirical [explicit regularizer](#), as well as bound the supremum deviation between the predictions and the true values. We remark that the value of  $\gamma$  also determines the complexity of the function class. On one hand, the generalization gap explicitly depends on and increases with  $\gamma$ . However, when  $\gamma$  is large, the constraints on  $U, V$  are milder, so that  $\widehat{L}(U, V)$  can be made smaller.

Finally, the required sample size heavily depends on the value of the explicit regularizer at the optima ( $\alpha/d_1$ ), and hence, on the dropout rate  $p$ . In particular, increasing the dropout rate increases the regularization parameter  $\lambda := \frac{p}{1-p}$ , thereby intensifies the penalty due to the explicit regularizer. Intuitively, a larger dropout rate  $p$  results in a smaller  $\alpha$ , thereby a tighter generalization gap can be guaranteed. We show through experiments that that is indeed the case in practice.

## 11.2 Deep neural networks

Next, we focus on neural networks with multiple hidden layers. Let  $\mathcal{X} \subseteq \mathbb{R}^{d_0}$  and  $\mathcal{Y} \subseteq \mathbb{R}^{d_k}$  denote the input and output spaces, respectively. Let  $\mathcal{D}$  denote the joint probability distribution on  $\mathcal{X} \times \mathcal{Y}$ . Given  $n$  examples  $\{(x_i, y_i)\}_{i=1}^n \sim \mathcal{D}^n$  drawn i.i.d. from the joint distribution and a loss function  $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$ , the goal of learning

is to find a hypothesis  $f_w : \mathcal{X} \rightarrow \mathcal{Y}$ , parameterized by  $w$ , that has a small *population risk*  $L(w) := \mathbb{E}_{\mathcal{D}}[\ell(f_w(x), y)]$ .

We focus on the squared  $\ell_2$  loss, i.e.,  $\ell(y, y') = \|y - y'\|^2$ , and study the generalization properties of the dropout algorithm for minimizing the *empirical risk*  $\widehat{L}(w) := \frac{1}{n} \sum_{i=1}^n [\|y_i - f_w(x_i)\|^2]$ . We consider the hypothesis class associated with feed-forward neural networks with  $k$  layers, i.e., functions of the form  $f_w(x) = W_k \sigma(W_{k-1} \sigma(\dots W_2 \sigma(W_1 x) \dots))$ , where  $W_i \in \mathbb{R}^{d_i \times d_{i-1}}$ , for  $i \in [k]$ , is the weight matrix at  $i$ -th layer. The parameter  $w$  is the collection of weight matrices  $\{W_k, W_{k-1}, \dots, W_1\}$  and  $\sigma : \mathbb{R} \rightarrow \mathbb{R}$  is an activation function applied entrywise to an input vector.

In modern machine learning systems, rather than talk about a certain network topology, we should think in terms of layer topology where each layer could have different characteristics – for example, fully connected, locally connected, or convolutional. In convolutional neural networks, it is a common practice to apply dropout only to the fully connected layers and not to the convolutional layers. Furthermore, in deep regression, it has been observed that applying dropout to only one of the hidden layers is most effective [LMAPH19]. In our study, dropout is applied on top of the learned representations or *features*, i.e. the output of the top hidden layer. In this case, dropout updates can be viewed as stochastic gradient descent iterates on the *dropout objective*:

$$\widehat{L}_{\text{drop}}(w) := \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{\text{B}} \|y_i - W_k \text{B} \sigma(W_{k-1} \sigma(\dots W_2 \sigma(W_1 x_i) \dots))\|^2 \quad (11.8)$$

where  $\text{B}$  is a diagonal random matrix with diagonal elements distributed identically and independently as  $\text{B}_{ii} \sim \frac{1}{1-p} \text{Bern}(1-p)$ ,  $i \in [d_{k-1}]$ , for some *dropout rate*  $p$ . We seek to understand the *explicit regularizer* due to dropout:

$$\widehat{R}(w) := \widehat{L}_{\text{drop}}(w) - \widehat{L}(w) \quad (\text{explicit regularizer})$$

We denote the output of the  $i$ -th hidden node in the  $j$ -th hidden layer on an input vector  $x$  by  $a_{i,j}(x) \in \mathbb{R}$ ; for example,  $a_{1,2}(x) = \sigma(W_2(1, \cdot)^\top \sigma(W_1 x))$ . Similarly, the vector  $a_j(x) \in \mathbb{R}^{d_j}$  denotes the activation of the  $j$ -th layer on input  $x$ . Using this notation, we can conveniently rewrite the Dropout objective (see Equation 11.8) as  $\widehat{L}_{\text{drop}}(w) := \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{\text{B}} \|y_i - W_k \text{B} a_{k-1}(x_i)\|^2$ . It is then easy to show that the *explicit regularizer* due to dropout is given as follows.

**Proposition 11.2.1** (Dropout regularizer in deep regression).

$$\widehat{L}_{\text{drop}}(w) = \widehat{L}(w) + \widehat{R}(w), \quad \text{where } \widehat{R}(w) = \lambda \sum_{j=1}^{d_{k-1}} \|W_k(:, j)\|^2 \widehat{a}_j^2.$$

where  $\hat{a}_j = \sqrt{\frac{1}{n} \sum_{i=1}^n a_{j,k-1}(x_i)^2}$  and  $\lambda = \frac{p}{1-p}$  is the regularization parameter.

*Proof.* Recall that  $\mathbb{E}[\mathbf{B}_{ii}] = 1$  and  $\text{Var}(\mathbf{B}_{ii}) = \frac{p}{1-p}$ . Conditioned on  $x, y$  in the current mini-batch, we have that

$$\mathbb{E}_{\mathbf{B}}[\|y - \mathbf{W}_k \mathbf{B} \mathbf{a}_{k-1}(x)\|^2] = \sum_{i=1}^{d_k} \mathbb{E}_{\mathbf{B}}(y_i - \mathbf{W}_k(i, :)^{\top} \mathbf{B} \mathbf{a}_{k-1}(x))^2. \quad (11.9)$$

The following holds for each of the summands above:

$$\begin{aligned} \mathbb{E}_{\mathbf{B}}(y_i - \mathbf{W}_k(i, :)^{\top} \mathbf{B} \mathbf{a}_{k-1}(x))^2 &= \left( \mathbb{E}_{\mathbf{B}}[y_i - \mathbf{W}_k(i, :)^{\top} \mathbf{B} \mathbf{a}_{k-1}(x)] \right)^2 \\ &\quad + \text{Var}(y_i - \mathbf{W}_k(i, :)^{\top} \mathbf{B} \mathbf{a}_{k-1}(x)). \end{aligned}$$

Since  $\mathbb{E}[\mathbf{B}] = \mathbf{I}$ , the first term on right hand side is equal to  $(y_i - \mathbf{W}_k(i, :)^{\top} \mathbf{a}_{k-1}(x))^2$ . For the second term we have

$$\begin{aligned} \text{Var}(y_i - \mathbf{W}_k(i, :)^{\top} \mathbf{B} \mathbf{a}_{k-1}(x)) &= \text{Var}(\mathbf{W}_k(i, :)^{\top} \mathbf{B} \mathbf{a}_{k-1}(x)) \\ &= \text{Var}\left(\sum_{j=1}^{d_{k-1}} \mathbf{W}_k(i, j) \mathbf{B}_{jj} a_{j,k-1}(x)\right) \\ &= \sum_{j=1}^{d_{k-1}} (\mathbf{W}_k(i, j) a_{j,k-1}(x))^2 \text{Var}(\mathbf{B}_{jj}) \\ &= \frac{p}{1-p} \sum_{j=1}^{d_{k-1}} \mathbf{W}_k(i, j)^2 a_{j,k-1}(x)^2 \end{aligned}$$

Plugging the above into Equation (11.9)

$$\begin{aligned} \mathbb{E}_{\mathbf{B}}[\|y - \mathbf{W}_k \mathbf{B} \mathbf{a}_{k-1}(x)\|^2] &= \|y - \mathbf{W}_k \mathbf{a}_{k-1}(x)\|^2 \\ &\quad + \frac{p}{1-p} \sum_{j=1}^{d_{k-1}} \|\mathbf{W}_k(:, j)\|^2 a_{j,k-1}(x)^2 \end{aligned}$$

Now taking the empirical average with respect to  $x, y$ , we get

$$\hat{L}_{\text{drop}}(\mathbf{w}) = \hat{L}(\mathbf{w}) + \frac{p}{1-p} \sum_{j=1}^{d_{k-1}} \|\mathbf{W}_k(:, j)\|^2 \hat{a}_j^2 = \hat{L}(\mathbf{w}) + \hat{R}(\mathbf{w})$$

which completes the proof.  $\square$

The **explicit regularizer**  $\hat{R}(\mathbf{w})$  is the summation over hidden nodes, of the product of the squared norm of the outgoing weights with the empirical second moment of the output of the corresponding neuron. For a two layer neural network with ReLU, when the input distribution is symmetric and isotropic, the expected regularizer is equal to the squared  $\ell_2$  path-norm of the network [NTS15b]. Such a connection has been previously established for deep linear networks [MAV18, MA19]; here we extend that result to single hidden layer ReLU networks.

**Proposition 11.2.2.** Consider a two layer neural network  $f_w(\cdot)$  with ReLU activation functions in the hidden layer. Furthermore, assume that the marginal input distribution  $\mathbb{P}_{\mathcal{X}}(x)$  is symmetric and isotropic, i.e.,  $\mathbb{P}_{\mathcal{X}}(x) = \mathbb{P}_{\mathcal{X}}(-x)$  and  $\mathbb{E}[xx^\top] = I$ . Then the expected [explicit regularizer](#) due to dropout is given as

$$R(w) := \mathbb{E}[\widehat{R}(w)] = \frac{\lambda}{2} \sum_{i_0, i_1, i_2=1}^{d_0, d_1, d_2} W_2(i_2, i_1)^2 W_1(i_1, i_0)^2, \quad (11.10)$$

*Proof of Proposition 11.2.2.* Using Proposition 11.2.1, we have that:

$$R(w) = \mathbb{E}[\widehat{R}(w)] = \lambda \sum_{j=1}^{d_1} \|W_2(:, j)\|^2 \mathbb{E}[\sigma(W_1(j, :)^\top x)^2]$$

It remains to calculate the quantity  $\mathbb{E}_x[\sigma(W_1(j, :)^\top x)^2]$ . By symmetry assumption, we have that  $\mathbb{P}_{\mathcal{X}}(x) = \mathbb{P}_{\mathcal{X}}(-x)$ . As a result, for any  $v \in \mathbb{R}^{d_0}$ , we have that  $\mathbb{P}(v^\top x) = \mathbb{P}(-v^\top x)$  as well. That is, the random variable  $z_j := W_1(j, :)^\top x$  is also symmetric about the origin. It is easy to see that  $\mathbb{E}_z[\sigma(z)^2] = \frac{1}{2} \mathbb{E}_z[z^2]$ .

$$\begin{aligned} \mathbb{E}_z[\sigma(z)^2] &= \int_{-\infty}^{\infty} \sigma(z)^2 d\mu(z) \\ &= \int_0^{\infty} \sigma(z)^2 d\mu(z) = \int_0^{\infty} z^2 d\mu(z) \\ &= \frac{1}{2} \int_{-\infty}^{\infty} z^2 d\mu(z) = \frac{1}{2} \mathbb{E}_z[z^2]. \end{aligned}$$

Plugging back the above identity in the expression of  $R(w)$ , we get that

$$R(w) = \frac{\lambda}{2} \sum_{j=1}^{d_1} \|W_2(:, j)\|^2 \mathbb{E}[(W_1(j, :)^\top x)^2] = \frac{\lambda}{2} \sum_{j=1}^{d_1} \|W_2(:, j)\|^2 \|W_1(j, :)\|^2$$

where the second equality follows from the assumption that the distribution is isotropic.  $\square$

### 11.3 Landscape of the Optimization Problem

While the focus in Section 11.2 was on understanding the implicit bias of dropout in terms of the global optima of the resulting regularized learning problem, here we focus on computational aspects of dropout as an optimization procedure. Since dropout is a first-order method and the landscape of the Dropout objective (e.g., Problem 11.11) is highly non-convex, we can perhaps only hope to find a *local* minimum, that too provided if the problem has no degenerate saddle points [LSJR16, GHJY15b]. Therefore, in this section, we pose the following questions: *What is the implicit bias of dropout in*

terms of local minima? Do local minima share anything with global minima structurally or in terms of the objective? Can dropout find a local optimum?

For the sake of simplicity of analysis, we focus on the case of single hidden layer *linear* autoencoders with tied weights, i.e.  $U = V$ . We assume that the input distribution is isotropic, i.e.  $C_x = I$ . In this case, the population risk reduces to

$$\begin{aligned} \mathbb{E}[\|y - UU^\top x\|^2] &= \text{Tr}(C_y) - 2\langle C_{yx}, UU^\top \rangle + \|UU^\top\|_F^2 \\ &= \|M - UU^\top\|_F^2 + \text{Tr}(C_y) - \|C_{yx}\|_F^2 \end{aligned}$$

where  $M = \frac{C_{yx} + C_{xy}}{2}$ . Ignoring the terms that are independent of the weights matrix  $U$ , the goal is to minimize  $L(U) = \|M - UU^\top\|_F^2$ . Using Dropout amounts to solving the following problem:

$$\min_{U \in \mathbb{R}^{d_0 \times d_1}} L_\theta(U) := \|M - UU^\top\|_F^2 + \lambda \underbrace{\sum_{i=1}^{d_1} \|u_i\|^4}_{R(U)} \quad (11.11)$$

We can characterize the global optima of the problem above as follows.

**Theorem 11.3.1.** For any  $j \in [r]$ , let  $\kappa_j := \frac{1}{j} \sum_{i=1}^j \lambda_i(C_{yx})$ . Furthermore, define  $\rho := \max\{j \in [r] : \lambda_j(C_{yx}) > \frac{\lambda_j \kappa_j}{r + \lambda_j}\}$ . Then, if  $U_*$  is a global optimum of Problem 11.11, it satisfies that  $U_* U_*^\top = \mathcal{S}_{\frac{\lambda \rho \kappa_\rho}{r + \lambda \rho}}(C_{yx})$ .

Next, it is easy to see that the gradient of the objective of Problem 11.11 is given by

$$\nabla L_\theta(U) = 4(UU^\top - M)U + 4\lambda U \text{diag}(U^\top U).$$

We also make the following important observation about the critical points of Problem 11.11. Lemma 11.3.2 allows us to bound different norms of the critical points, as will be seen later in the proofs.

**Lemma 11.3.2.** If  $U$  is a critical point of Problem 11.11, then it holds that  $UU^\top \preceq M$ .

*Proof of Lemma 11.3.2.* Since  $\nabla L_\theta(U) = 0$ , we have that

$$(M - UU^\top)U = \lambda U \text{diag}(U^\top U)$$

multiply both sides from right by  $U^\top$  and rearrange to get

$$MUU^\top = UU^\top UU^\top + \lambda U \text{diag}(U^\top U)U^\top \quad (11.12)$$

Note that the right hand side is symmetric, which implies that the left hand side must be symmetric as well, i.e.

$$MUU^\top = (MUU^\top)^\top = UU^\top M,$$

so that  $M$  and  $UU^\top$  commute. Note that in Equation (11.12),  $U\text{diag}(U^\top U)U^\top \succeq 0$ . Thus,  $MUU^\top \succeq UU^\top UU^\top$ . Let  $UU^\top = W\Gamma W^\top$  be a compact eigen-decomposition of  $UU^\top$ . We get

$$MUU^\top = MWW^\top \succeq UU^\top UU^\top = W\Gamma^2 W^\top.$$

Multiplying from right and left by  $W\Gamma^{-1}$  and  $W^\top$  respectively, we get that  $W^\top MW \succeq \Gamma$  which completes the proof.  $\square$

We show in Section 11.3.1 that (a) local minima of Problem 11.11 inherit the same implicit bias as the global optima, i.e. all local minima are equalized. Then, in Section 11.3.2, we show that for sufficiently small regularization parameter, (b) there are no spurious local minima, i.e. all local minima are global, and (c) all saddle points are non-degenerate (see Definition 11.3.4).

### 11.3.1 Implicit bias in local optima

Recall that the population risk  $L(U)$  is rotation invariant, i.e.  $L(UQ) = L(U)$  for any rotation matrix  $Q$ . Now, if the weight matrix  $U$  were not equalized, then there exist indices  $i, j \in [r]$  such that  $\|u_i\| > \|u_j\|$ . We show that it is easy to design a rotation matrix (equal to identity everywhere except for columns  $i$  and  $j$ ) that moves mass from  $u_i$  to  $u_j$  such that the difference in the norms of the corresponding columns of  $UQ$  decreases strictly while leaving the norms of other columns invariant. In other words, this rotation strictly reduces the regularizer and hence the objective. Formally, this implies the following result.

**Lemma 11.3.3.** All local minima of Problem 11.11 are equalized, i.e. if  $U$  is a local optimum, then  $\|u_i\| = \|u_j\| \forall i, j \in [r]$ .

Lemma 11.3.3 unveils a fundamental property of dropout. As soon as we perform dropout in the hidden layer – *no matter how small the dropout rate* – all local minima become equalized. We illustrate this using a toy problem in Figure 11.1.

*Proof of Lemma 11.3.3.* We show that if  $U$  is not equalized, then any  $\epsilon$ -neighborhood of  $U$  contains a point with dropout objective strictly smaller than  $L_\theta(U)$ . More formally, for any  $\epsilon > 0$ , we exhibit a rotation  $Q_\epsilon$  such that  $\|U - UQ_\epsilon\|_F \leq \epsilon$  and  $L_\theta(UQ_\epsilon) < L_\theta(U)$ . Let  $U$  be a critical point of Problem 11.11 that is not equalized, i.e. there exists two columns of  $U$  with different norms. Without loss of generality, let  $\|u_1\| > \|u_2\|$ . We design a rotation matrix  $Q$  such that it is almost an isometry, but it moves mass from  $u_1$  to  $u_2$ . Consequently, the new factor becomes “less un-equalized” and achieves a smaller regularizer,



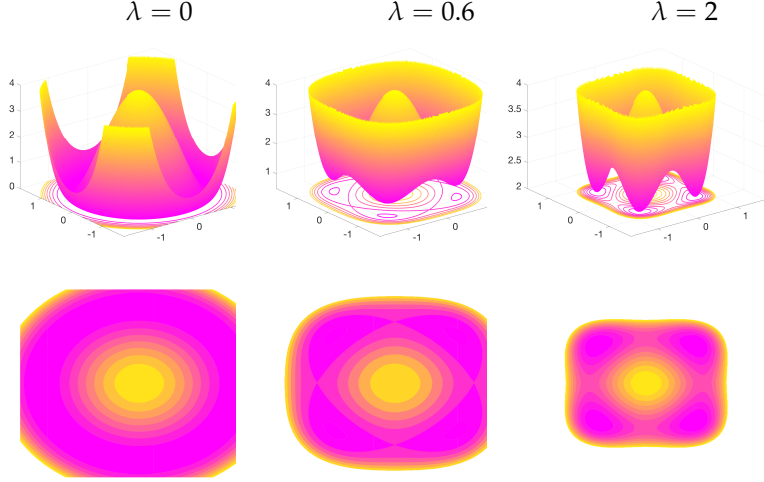


Figure 11.1: Optimization landscape (top) and contour plot (bottom) for a single hidden-layer linear autoencoder network with one dimensional input and output and a hidden layer of width  $r = 2$  with dropout, for different values of the regularization parameter  $\lambda$ . Left: for  $\lambda = 0$  the problem reduces to squared loss minimization, which is rotation invariant as suggested by the level sets. Middle: for  $\lambda > 0$  the global optima shrink toward the origin. All local minima are global, and are equalized, i.e. the weights are parallel to the vector  $(\pm 1, \pm 1)$ . Right: as  $\lambda$  increases, global optima shrink further.

while preserving the value of the loss. To that end, define

$$Q_\delta := \begin{bmatrix} \sqrt{1-\delta^2} & -\delta & 0 \\ \delta & \sqrt{1-\delta^2} & 0 \\ 0 & 0 & I_{r-2} \end{bmatrix}$$

and let  $\widehat{U} := UQ_\delta$ . It is easy to verify that  $Q_\delta$  is indeed a rotation. First, we show that for any  $\epsilon$ , as long as  $\delta^2 \leq \frac{\epsilon^2}{2\text{Tr}(M)}$ , we have  $\widehat{U} \in \mathcal{B}_\epsilon(U)$ :

$$\begin{aligned} \|U - \widehat{U}\|_F^2 &= \sum_{i=1}^r \|u_i - \widehat{u}_i\|^2 \\ &= \|u_1 - \sqrt{1-\delta^2}u_1 - \delta u_2\|^2 + \|u_2 - \sqrt{1-\delta^2}u_2 + \delta u_1\|^2 \\ &= 2(1 - \sqrt{1-\delta^2})(\|u_1\|^2 + \|u_2\|^2) \\ &\leq 2\delta^2 \text{Tr}(M) \leq \epsilon^2 \end{aligned}$$

where the second to last inequality follows from Lemma 11.3.2, because  $\|u_1\|^2 + \|u_2\|^2 \leq \|U\|_F^2 = \text{Tr}(UU^\top) \leq \text{Tr}(M)$ , and also the fact that  $1 - \sqrt{1-\delta^2} = \frac{1-\delta^2}{1+\sqrt{1-\delta^2}} \leq \delta^2$ .

Next, we show that for small enough  $\delta$ , the value of  $L_\theta$  at  $\widehat{U}$  is strictly smaller than that of  $U$ . Observe that

$$\begin{aligned} \|\widehat{u}_1\|^2 &= (1-\delta^2)\|u_1\|^2 + \delta^2\|u_2\|^2 + 2\delta\sqrt{1-\delta^2}u_1^\top u_2 \\ \|\widehat{u}_2\|^2 &= (1-\delta^2)\|u_2\|^2 + \delta^2\|u_1\|^2 - 2\delta\sqrt{1-\delta^2}u_1^\top u_2 \end{aligned}$$

and the remaining columns will not change, i.e. for  $i = 3, \dots, r$ ,  $\widehat{u}_i = u_i$ . Together with the fact that  $Q_\delta$  preserves the norms, i.e.  $\|U\|_F = \|UQ_\delta\|_F$ , we get

$$\|\widehat{u}_1\|^2 + \|\widehat{u}_2\|^2 = \|u_1\|^2 + \|u_2\|^2. \quad (11.13)$$

Let  $\delta = -c \cdot \text{sgn}(\mathbf{u}_1^\top \mathbf{u}_2)$  for a small enough  $c > 0$  such that  $\|\mathbf{u}_2\| < \|\hat{\mathbf{u}}_2\| \leq \|\hat{\mathbf{u}}_1\| < \|\mathbf{u}_1\|$ . Using Equation (11.13), This implies that  $\|\hat{\mathbf{u}}_1\|^4 + \|\hat{\mathbf{u}}_2\|^4 < \|\mathbf{u}_1\|^4 + \|\mathbf{u}_2\|^4$ , which in turn gives us  $R(\hat{\mathbf{U}}) < R(\mathbf{U})$  and hence  $L_\theta(\hat{\mathbf{U}}) < L_\theta(\mathbf{U})$ . Therefore, a non-equalized critical point cannot be local minimum, hence the first claim of the lemma.  $\square$

### 11.3.2 Landscape properties

Next, we characterize the solutions to which dropout converges. We do so by understanding the optimization landscape of Problem 11.11. Central to our analysis, is the following notion of *strict saddle property*.

**Definition 11.3.4** (Strict saddle point/property). Let  $f : \mathcal{U} \rightarrow \mathbb{R}$  be a twice differentiable function and let  $\mathbf{U} \in \mathcal{U}$  be a critical point of  $f$ . Then,  $\mathbf{U}$  is a *strict saddle point* of  $f$  if the Hessian of  $f$  at  $\mathbf{U}$  has at least one negative eigenvalue, i.e.  $\lambda_{\min}(\nabla^2 f(\mathbf{U})) < 0$ . Furthermore,  $f$  satisfies *strict saddle property* if all saddle points of  $f$  are strict saddle.

Strict saddle property ensures that for any critical point  $\mathbf{U}$  that is not a local optimum, the Hessian has a significant negative eigenvalue which allows first order methods such as gradient descent (GD) and stochastic gradient descent (SGD) to escape saddle points and converge to a local minimum [LSJR16, GHJY15b]. Following this idea, there has been a flurry of works on studying the landscape of different machine learning problems, including low rank matrix recovery [BNS16b], generalized phase retrieval problem [SQW16b], matrix completion [GLM16b], deep linear networks [Kaw16], matrix sensing and robust PCA [GJZ17b] and tensor decomposition [GHJY15b], making a case for global optimality of first order methods.

For the special case of no regularization (i.e.  $\lambda = 0$ ; equivalently, no dropout), Problem 11.11 reduces to standard squared loss minimization which has been shown to have no spurious local minima and satisfy strict saddle property (see, e.g. [BH89, JGN<sup>+</sup>17]). However, the regularizer induced by dropout can potentially introduce new spurious local minima as well as degenerate saddle points. Our next result establishes that that is not the case, at least when the dropout rate is sufficiently small.

**Theorem 11.3.5.** Let  $r := \text{Rank}(\mathbf{M})$ . Assume that  $d_1 \leq d_0$  and that the regularization parameter satisfies  $\lambda < \frac{r\lambda_r(\mathbf{M})}{(\sum_{i=1}^r \lambda_i(\mathbf{M})) - r\lambda_r(\mathbf{M})}$ . Then it holds for Problem 11.11 that

1. all local minima are global,
2. all saddle points are strict saddle points.

A few remarks are in order. First, the assumption  $d_1 \leq d_0$  is by no means restrictive, since the network map  $\mathbf{U}\mathbf{U}^\top \in \mathbb{R}^{d_0 \times d_0}$  has rank at

most  $d_0$ , and letting  $d_1 > d_0$  does not increase the expressivity of the function class represented by the network. Second, Theorem 11.3.5 guarantees that any critical point  $U$  that is not a global optimum is a strict saddle point, i.e.  $\nabla^2 L(U, U)$  has a negative eigenvalue. This property allows first order methods, such as dropout, to escape such saddle points. Third, note that the guarantees in Theorem 11.3.5 hold when the regularization parameter  $\lambda$  is sufficiently small. Assumptions of this kind are common in the literature (see, for example [GJZ17b]). While this is a *sufficient* condition for the result in Theorem 11.3.5, it is not clear if it is *necessary*.

*Proof of Theorem 11.3.5.* Here we outline the main steps in the proof of Theorem 11.3.5.

1. In Lemma 11.3.3, we show that the set of non-equalized critical points does not include any local optima. Furthermore, Lemma 11.3.6 shows that all such points are strict saddles.
2. In Lemma 11.3.7, we give a closed-form characterization of all the equalized critical points in terms of the eigendecomposition of  $M$ . We then show that if  $\lambda$  is chosen appropriately, all such critical points that are not global optima, are strict saddle points.
3. It follows from Item 1 and Item 2 that if  $\lambda$  is chosen appropriately, then all critical points that are not global optimum, are strict saddle points.

□

**Lemma 11.3.6.** *All critical points of Problem 11.11 that are not equalized, are strict saddle points.*

*Proof of Lemma 11.3.6.* By Lemma 11.3.3, the set of non-equalized critical points does not include any local optima. We show that all such points are strict saddles. Let  $U$  be a critical point that is not equalized. To show that  $U$  is a strict saddle point, it suffices to show that the Hessian has a negative eigenvalue. In here, we exhibit a curve along which the second directional derivative is negative. Assume, without loss of generality that  $\|u_1\| > \|u_2\|$  and consider the curve

$$\Delta(t) := [(\sqrt{1-t^2}-1)u_1 + tu_2, (\sqrt{1-t^2}-1)u_2 - tu_1, o_{d,r-2}]$$

It is easy to check that for any  $t \in \mathbb{R}$ ,  $L(U + \Delta(t)) = L(U)$  since  $U + \Delta(t)$  is essentially a rotation on  $U$  and  $L$  is invariant under

rotations. Observe that

$$\begin{aligned}
g(t) &:= L_\theta(\mathbf{U} + \Delta(t)) \\
&= L_\theta(\mathbf{U}) + \|\sqrt{1-t^2}\mathbf{u}_1 + t\mathbf{u}_2\|^4 - \|\mathbf{u}_1\|^4 + \|\sqrt{1-t^2}\mathbf{u}_2 - t\mathbf{u}_1\|^4 - \|\mathbf{u}_2\|^4 \\
&= L_\theta(\mathbf{U}) - 2t^2(\|\mathbf{u}_1\|^4 + \|\mathbf{u}_2\|^4) + 8t^2(\mathbf{u}_1\mathbf{u}_2)^2 + 4t^2\|\mathbf{u}_1\|^2\|\mathbf{u}_2\|^2 \\
&\quad + 4t\sqrt{1-t^2}\mathbf{u}_1^\top\mathbf{u}_2(\|\mathbf{u}_1\|^2 - \|\mathbf{u}_2\|^2) + O(t^3).
\end{aligned}$$

The derivative of  $g$  then is given as

$$\begin{aligned}
g'(t) &= -4t(\|\mathbf{u}_1\|^4 + \|\mathbf{u}_2\|^4) + 16t(\mathbf{u}_1\mathbf{u}_2)^2 + 8t\|\mathbf{u}_1\|^2\|\mathbf{u}_2\|^2 \\
&\quad + 4\left(\sqrt{1-t^2} - \frac{t^2}{\sqrt{1-t^2}}\right)(\mathbf{u}_1^\top\mathbf{u}_2)(\|\mathbf{u}_1\|^2 - \|\mathbf{u}_2\|^2) + O(t^2).
\end{aligned}$$

Since  $\mathbf{U}$  is a critical point and  $L_\theta$  is continuously differentiable, it should hold that

$$g'(0) = 4(\mathbf{u}_1^\top\mathbf{u}_2)(\|\mathbf{u}_1\|^2 - \|\mathbf{u}_2\|^2) = 0.$$

Since by assumption  $\|\mathbf{u}_1\|^2 - \|\mathbf{u}_2\|^2 > 0$ , it should be the case that  $\mathbf{u}_1^\top\mathbf{u}_2 = 0$ . We now consider the second order directional derivative:

$$\begin{aligned}
g''(0) &= -4(\|\mathbf{u}_1\|^4 + \|\mathbf{u}_2\|^4) + 16(\mathbf{u}_1\mathbf{u}_2)^2 + 8\|\mathbf{u}_1\|^2\|\mathbf{u}_2\|^2 \\
&= -4(\|\mathbf{u}_1\|^2 - \|\mathbf{u}_2\|^2)^2 < 0
\end{aligned}$$

which completes the proof.  $\square$

We now focus on the critical points that are equalized, i.e. points  $\mathbf{U}$  such that  $\nabla L_\theta(\mathbf{U}) = \mathbf{o}$  and  $\text{diag}(\mathbf{U}^\top\mathbf{U}) = \frac{\|\mathbf{U}\|_F^2}{d_1}\mathbf{I}$ .

**Lemma 11.3.7.** *Let  $r := \text{Rank}(M)$ . Assume that  $d_1 \leq d_0$  and  $\lambda < \frac{r\lambda_r}{\sum_{i=1}^r(\lambda_i - \lambda_r)}$ . Then all equalized local minima are global. All other equalized critical points are strict saddle points.*

*Proof of Lemma 11.3.7.* Let  $\mathbf{U}$  be a critical point that is equalized.

Furthermore, let  $r'$  be the rank of  $\mathbf{U}$ , and  $\mathbf{U} = \mathbf{W}\Sigma\mathbf{V}^\top$  be its rank- $r'$  SVD, i.e.  $\mathbf{W} \in \mathbb{R}^{d_0 \times r'}$ ,  $\mathbf{V} \in \mathbb{R}^{d_1 \times r'}$  are such that  $\mathbf{U}^\top\mathbf{U} = \mathbf{V}^\top\mathbf{V} = \mathbf{I}_{r'}$  and  $\Sigma \in \mathbb{R}^{r' \times r'}$ , is a positive definite diagonal matrix whose diagonal entries are sorted in descending order. We have:

$$\begin{aligned}
\nabla L_\theta(\mathbf{U}) &= 4(\mathbf{U}\mathbf{U}^\top - \mathbf{M})\mathbf{U} + 4\lambda\mathbf{U}\text{diag}(\mathbf{U}^\top\mathbf{U}) = \mathbf{o} \\
\implies \mathbf{U}\mathbf{U}^\top\mathbf{U} + \frac{\lambda\|\mathbf{U}\|_F^2}{d_1}\mathbf{U} &= \mathbf{M}\mathbf{U} \\
\implies \mathbf{W}\Sigma^3\mathbf{V}^\top + \frac{\lambda\|\Sigma\|_F^2}{d_1}\mathbf{W}\Sigma\mathbf{V}^\top &= \mathbf{M}\mathbf{W}\Sigma\mathbf{V}^\top \\
\implies \Sigma^2 + \frac{\lambda\|\Sigma\|_F^2}{d_1}\mathbf{I} &= \mathbf{W}^\top\mathbf{M}\mathbf{W}
\end{aligned}$$

Since the left hand side of the above equality is diagonal, it implies that  $W \in \mathbb{R}^{d_0 \times r'}$  corresponds to some  $r'$  eigenvectors of  $M$ . Let  $\mathcal{E} \subseteq [d_0]$ ,  $|\mathcal{E}| = r'$  denote the set of eigenvectors of  $M$  that are present in  $W$ . The above equality is equivalent of the following system of linear equations:

$$\left(I + \frac{\lambda}{d_1} \mathbf{1}\mathbf{1}^\top\right) \text{diag}(\Sigma^2) = \vec{\lambda},$$

where  $\vec{\lambda} = \text{diag}(W^\top MW)$ . The solution to the linear system of equations above is given by

$$\text{diag}(\Sigma^2) = \left(I - \frac{\lambda}{d_1 + \lambda r'}\right) \vec{\lambda} = \vec{\lambda} - \frac{\lambda \sum_{i=1}^{r'} \lambda_i}{d_1 + \lambda r'} \mathbf{1}_{r'}. \quad (11.14)$$

Thus, the set  $\mathcal{E}$  belongs to one of the following categories:

0.  $\mathcal{E} = [r']$ ,  $r' > \rho$
1.  $\mathcal{E} = [r']$ ,  $r' = \rho$
2.  $\mathcal{E} = [r']$ ,  $r' < \rho$
3.  $\mathcal{E} \neq [r']$

We provide a case by case analysis for the above partition here.

**Case 0.** [ $\mathcal{E} = [r']$ ,  $r' > \rho$ ]. We show that  $\mathcal{E}$  cannot belong to this class, i.e. when  $\mathcal{E} = [r']$ , it should hold that  $r' \leq \rho$ . To see this, consider the  $r'$ -th linear equation in Equation (11.14):

$$\sigma_{r'}^2 = \lambda_{r'} - \frac{\lambda \sum_{i=1}^{r'} \lambda_i}{d_1 + \lambda r'}.$$

Since  $\text{Rank } U = r'$ , it follows that  $\sigma_{r'} > 0$ , which in turn implies that

$$\lambda_{r'} > \frac{\lambda \sum_{i=1}^{r'} \lambda_i}{d_1 + \lambda r'} = \frac{\lambda r' \kappa_{r'}}{d_1 + \lambda r'}.$$

It follows from maximality of  $\rho$  in Theorem 11.3.1 that  $r' \leq \rho$ .

**Case 1.** [ $\mathcal{E} = [r']$ ,  $r' = \rho$ ]. When  $W$  corresponds to the top- $\rho$  eigenvectors of  $M$ , we retrieve a global optimum described by Theorem 11.3.1. Therefore, all such critical points are global minima.

**Case 2.** [ $\mathcal{E} = [r']$ ,  $r' < \rho$ ]. Let  $W_{d_0} := [W, W_\perp]$  be a complete eigenbasis for  $M$  corresponding to eigenvalues of  $M$  in descending order, where  $W_\perp \in \mathbb{R}^{d_0 \times d_0 - r'}$  constitutes a basis for the orthogonal subspace of  $W$ . For rank deficient  $M$ ,  $W_\perp$  contains the null-space of  $M$ , and hence eigenvectors corresponding to zero eigenvalues of  $M$ . Similarly, let  $V_\perp \in \mathbb{R}^{d_1 \times d_1 - r'}$  span the orthogonal subspace of  $V$ , such that  $V_{d_1} := [V, V_\perp]$  forms an orthonormal basis for  $\mathbb{R}^{d_1}$ . Note that both  $W_\perp$  and  $V_\perp$  are well-defined since  $r' \leq \min\{d_0, d_1\}$ . Define

$U(t) = W_{d_0} \Sigma' V_{d_1}^\top$  where  $\Sigma' \in \mathbb{R}^{d_0 \times d_1}$  is diagonal with non-zero diagonal elements given as  $\sigma'_i = \sqrt{\sigma_i^2 + t^2}$  for  $i \leq d_1$ . Observe that

$$U(t)^\top U(t) = V \Sigma^2 V^\top + t^2 V_{d_1}^\top V_{d_1} = U^\top U + t^2 \mathbf{I}_{d_1}.$$

Thus, the parametric curve  $U(t)$  is equalized for all  $t$ . The population risk at  $U(t)$  equals:

$$\begin{aligned} L(U(t)) &= \sum_{i=1}^{d_1} (\lambda_i - \sigma_i^2 - t^2)^2 + \sum_{i=d_1+1}^{d_0} \lambda_i^2 \\ &= L(U) + d_1 t^4 - 2t^2 \sum_{i=1}^{d_1} (\lambda_i - \sigma_i^2). \end{aligned}$$

Furthermore, since  $U(t)$  is equalized, we obtain the following form for the regularizer:

$$\begin{aligned} R(U(t)) &= \frac{\lambda}{d_1} \|U(t)\|_F^4 = \frac{\lambda}{d_1} \left( \|U\|_F^2 + d_1 t^2 \right)^2 \\ &= R(U) + \lambda d_1 t^4 + 2\lambda t^2 \|U\|_F^2. \end{aligned}$$

Define  $g(t) := L(U(t)) + R(U(t))$ . We have that

$$g(t) = L(U) + R(U) + d_1 t^4 - 2t^2 \sum_{i=1}^{d_1} (\lambda_i - \sigma_i^2) + \lambda d_1 t^4 + 2\lambda t^2 \|U\|_F^2.$$

It is easy to verify that  $g'(0) = 0$ . Moreover, the second derivative of  $g$  at  $t = 0$  is given as:

$$g''(0) = -4 \sum_{i=1}^{d_1} (\lambda_i - \sigma_i^2) + 4\lambda \|U\|_F^2 = -4 \sum_{i=1}^{d_1} \lambda_i + 4(1 + \lambda) \|U\|_F^2 \quad (11.15)$$

We use  $\|U\|_F^2 = \sum_{i=1}^{r'} \sigma_i^2$  and Equation (11.14) to arrive at

$$\|U\|_F^2 = \text{tr} \Sigma^2 = \sum_{i=1}^{r'} \left( \lambda_i - \frac{\lambda \sum_{j=1}^{r'} \lambda_j}{d_1 + \lambda r'} \right) = \left( \sum_{i=1}^{r'} \lambda_i \right) \left( 1 - \frac{\lambda r'}{d_1 + \lambda r'} \right) = \frac{d_1 \sum_{i=1}^{r'} \lambda_i}{d_1 + \lambda r'}$$

Plugging back the above equality in Equation (11.15), we get

$$g''(0) = -4 \sum_{i=1}^{d_1} \lambda_i + 4 \frac{d_1 + d_1 \lambda}{d_1 + \lambda r'} \sum_{i=1}^{r'} \lambda_i = -4 \sum_{i=r'+1}^{d_1} \lambda_i + 4 \frac{(d_1 - r') \lambda}{d_1 + \lambda r'} \sum_{i=1}^{r'} \lambda_i$$

To get a sufficient condition for  $U$  to be a strict saddle point, it suf-

fices that  $g''(t)$  be negative at  $t = 0$ , i.e.

$$\begin{aligned}
 g''(0) < 0 &\implies \frac{(d_1 - r')\lambda}{d_1 + \lambda r'} \sum_{i=1}^{r'} \lambda_i < \sum_{i=r'+1}^{d_1} \lambda_i \\
 &\implies \lambda < \frac{(d_1 + \lambda r') \sum_{i=r'+1}^{d_1} \lambda_i}{(d_1 - r') \sum_{i=1}^{r'} \lambda_i} \\
 &\implies \lambda \left(1 - \frac{r' \sum_{i=r'+1}^{d_1} \lambda_i}{(d_1 - r') \sum_{i=1}^{r'} \lambda_i}\right) < \frac{d_1 \sum_{i=r'+1}^{d_1} \lambda_i}{(d_1 - r') \sum_{i=1}^{r'} \lambda_i} \\
 &\implies \lambda < \frac{d_1 \sum_{i=r'+1}^{d_1} \lambda_i}{(d_1 - r') \sum_{i=1}^{r'} \lambda_i - r' \sum_{i=r'+1}^{d_1} \lambda_i} \\
 &\implies \lambda < \frac{d_1 h(r')}{\sum_{i=1}^{r'} (\lambda_i - h(r'))}
 \end{aligned}$$

where  $h(r') := \frac{\sum_{i=r'+1}^{d_1} \lambda_i}{d_1 - r'}$  is the average of the tail eigenvalues  $\lambda_{r'+1}, \dots, \lambda_{d_1}$ . It is easy to see that the right hand side is monotonically decreasing with  $r'$ , since  $h(r')$  monotonically decreases with  $r'$ . Hence, it suffices to make sure that  $\lambda$  is smaller than the right hand side for the choice of  $r' = r - 1$ , where  $r := \text{Rank}(\mathbf{M})$ . That is,  $\lambda < \frac{r\lambda_r}{\sum_{i=1}^r (\lambda_i - \lambda_r)}$ .

**Case 3.** [ $\mathcal{E} \neq [r']$ ] We show that all such critical points are strict saddle points. Let  $w'$  be one of the top  $r'$  eigenvectors that are missing in  $\mathcal{W}$ . Let  $j \in \mathcal{E}$  be such that  $w_j$  is not among the top  $r'$  eigenvectors of  $\mathbf{M}$ . For any  $t \in [0, 1]$ , let  $W(t)$  be identical to  $W$  in all the columns but the  $j^{\text{th}}$  one, where  $w_j(t) = \sqrt{1 - t^2}w_j + tw'$ . Note that  $W(t)$  is still an orthogonal matrix for all values of  $t$ . Define the parametrized curve  $U(t) := W(t)\Sigma V^\top$  for  $t \in [0, 1]$  and observe that:

$$\begin{aligned}
 \|U - U(t)\|_F^2 &= \sigma_j^2 \|w_j - w_j(t)\|^2 \\
 &= 2\sigma_j^2(1 - \sqrt{1 - t^2}) \leq t^2 \text{Tr } \mathbf{M}
 \end{aligned}$$

That is, for any  $\epsilon > 0$ , there exist a  $t > 0$  such that  $U(t)$  belongs to the  $\epsilon$ -ball around  $U$ . We show that  $L_\theta(U(t))$  is strictly smaller than  $L_\theta(U)$ , which means  $U$  cannot be a local minimum. Note that this construction of  $U(t)$  guarantees that  $R(U') = R(U)$ . In particular, it is easy to see that  $U(t)^\top U(t) = U^\top U$ , so that  $U(t)$  remains equalized for all values of  $t$ . Moreover, we have that

$$\begin{aligned}
 L_\theta(U(t)) - L_\theta(U) &= \|\mathbf{M} - U(t)U(t)^\top\|_F^2 - \|\mathbf{M} - UU^\top\|_F^2 \\
 &= -2 \text{Tr}(\Sigma^2 W(t)^\top \mathbf{M} W(t)) + 2 \text{Tr}(\Sigma^2 W^\top \mathbf{M} W) \\
 &= -2\sigma_j^2 t^2 (w_j(t)^\top \mathbf{M} w_j(t) - w_j^\top \mathbf{M} w_j) < 0,
 \end{aligned}$$

where the last inequality follows because by construction  $w_j(t)^\top \mathbf{M} w_j(t) > w_j^\top \mathbf{M} w_j$ . Define  $g(t) := L_\theta(U(t)) = L(U(t)) + R(U(t))$ . To see that

such saddle points are non-degenerate, it suffices to show  $g''(0) < 0$ . It is easy to check that the second directional derivative at the origin is given by

$$g''(0) = -4\sigma_j^2(w_j(t)^\top M w_j(t) - w_j^\top M w_j) < 0,$$

which completes the proof.  $\square$

#### 11.4 Role of Parametrization

For least squares linear regression (i.e., for  $k = 1$  and  $\mathbf{u} = W_1^\top \in \mathbb{R}^{d_0}$  in Problem 11.8), we can show that using dropout amounts to solving the following regularized problem:

$$\min_{\mathbf{u} \in \mathbb{R}^{d_0}} \frac{1}{n} \sum_{i=1}^n (y_i - \mathbf{u}^\top \mathbf{x}_i)^2 + \lambda \mathbf{u}^\top \hat{\mathbf{C}} \mathbf{u}.$$

All the minimizers of the above problem are solutions to the following system of linear equations  $(1 + \lambda)X^\top X \mathbf{u} = X^\top \mathbf{y}$ , where  $X = [\mathbf{x}_1, \dots, \mathbf{x}_n]^\top \in \mathbb{R}^{n \times d_0}$ ,  $\mathbf{y} = [y_1, \dots, y_n]^\top \in \mathbb{R}^{n \times 1}$  are the design matrix and the response vector, respectively. Unlike Tikhonov regularization which yields solutions to the system of linear equations  $(X^\top X + \lambda I) \mathbf{u} = X^\top \mathbf{y}$  (a useful prior, discards the directions that account for small variance in data even when they exhibit good discriminability), the dropout regularizer manifests merely as a scaling of the parameters. This suggests that parametrization plays an important role in determining the nature of the resulting regularizer. However, a similar result was shown for deep linear networks [MA19] that the data dependent regularization due to dropout results in merely scaling of the parameters. At the same time, in the case of matrix sensing we see a richer class of regularizers. One potential explanation is that in the case of linear networks, we require a convolutional structure in the network to yield rich inductive biases. For instance, matrix sensing can be written as a two layer network in the following convolutional form:

$$\langle \mathbf{U} \mathbf{V}^\top, \mathbf{A} \rangle = \langle \mathbf{U}^\top, \mathbf{V}^\top \mathbf{A}^\top \rangle = \langle \mathbf{U}^\top, (\mathbf{I} \otimes \mathbf{V}^\top) \mathbf{A}^\top \rangle.$$

##### 11.4.1 Related Work



## *Unsupervised learning: Distribution Learning*

Much of the book so far concerned supervised learning —i.e., where training dataset consists of datapoints and a label indicating which class they belong to, and the model has to learn to produce the right label given an input. This chapter is an introduction to unsupervised learning, where one has randomly sampled datapoints but no labels or classes. We survey possible goals for this form of learning, and then focus on *distribution learning*, which addresses many of these goals.

### *12.1 Possible goals of unsupervised learning*

*Learn hidden/latent structure of data.* An example would be *Principal Component Analysis (PCA)*, concerned with finding the most important directions in the data. Other examples of structure learning can include sparse coding (aka dictionary learning) or nonnegative matrix factorization (NMF).

*Learn the distribution of the data.* A classic example is Pearson's 1893 contribution to theory of evolution by studying data about the crab population on Malta island. Biologists had sampled a thousand crabs in the wild, and measured 23 attributes (e.g., length, weight, etc.) for each. The presumption was that these datapoints should exhibit Gaussian distribution, but Pearson could not find a good fit to a Gaussian. He was able to show however that the distribution was actually *mixture* of two Gaussians. Thus the population consisted of two distinct species, which had diverged not too long ago in evolutionary terms.

In general, in density estimation the hypothesis is that the unlabeled dataset consists of iid samples from a fixed distribution, and model  $\theta$  learns representation of some distribution  $p_\theta(\cdot)$  that assigns a probability  $p_\theta(x)$  to datapoint  $x$ . This is the general problem of *density estimation*.

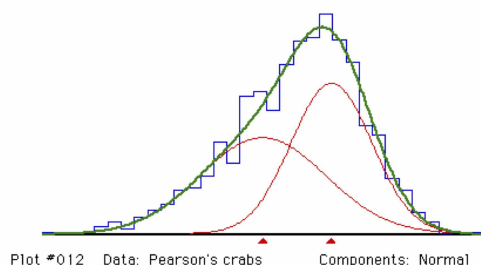
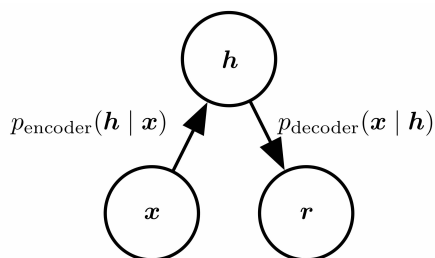


Figure 12.1: Visualization of Pearson's Crab Data as mixture of two Gaussians. (Credit: MIX homepage at McMaster University.)

One form of density estimation is to learn a *generative model*, where the learnt distribution has the form  $p_{\theta}(h, x)$  where  $x$  is the observable (i.e., datapoint) and  $h$  consists of a vector of hidden variables, often called *latent* variables. Then the density distribution of  $x$  is  $\int p_{\theta}(h, x)dh$ . In the crab example, the distribution is a mixture of Gaussians  $\mathcal{N}(\mu_1, \Sigma_1), \mathcal{N}(\mu_2, \Sigma_2)$  where the first contributes  $\rho_1$  fraction of samples and the other contributes  $1 - \rho_2$  fraction. Then  $\theta$  vector consists of parameters of the two Gaussians as well as  $\rho_1$ . The visible part  $x$  consists of attribute vector for a crab. Hidden vector  $h$  consists of a bit, indicating which of the two Gaussians this  $x$  was generated from, as well as the value of the gaussian random variable that generated  $x$ .

*Learning good representation/featurization of data* For example, the pixel representation of images may not be very useful in other tasks and one may desire a more “high level” representation that allows downstream tasks to be solved in a data-efficient way. One would hope to learn such featurization using unlabeled data.

In some settings, featurization is learnt via generative models: one assumes a data distribution  $p_{\theta}(h, x)$  as above and the featurization of the visible samplepoint  $x$  is assumed to be the hidden variable  $h$  that was used to generate it. More precisely, the hidden variable is a sample from the conditional distribution  $p(h|x)$ . This view of representation learning is used in the *autoencoders* described later.



For example, topic models are a simple probabilistic model of

Figure 12.2: Autoencoder defined using a density distribution  $p(h, x)$ , where  $h$  is the latent feature vector corresponding to visible vector  $x$ . The process of computing  $h$  given  $x$  is called “encoding” and the reverse is called “decoding.” In general applying the encoder on  $x$  followed by the decoder would not give  $x$  again, since the composed transformation is a sample from a distribution.

text generation, where  $x$  is some piece of text, and  $h$  is the proportion of specific topics (“sports,” “politics” etc.). Then one could imagine that  $h$  is some short and more high-level descriptor of  $x$ .

Many techniques for density estimation —such as variational methods, described later —also give a notion of a representation: the method for learning the distribution often also come with a candidate distribution for  $p(h|x)$ . This why students sometimes conflate representation learning with density estimation. But many of today’s approaches to representation learning do not boil down to

## 12.2 Training Objective for Learning Distributions: Log Likelihood

We wish to infer the best  $\theta$  given the set  $S$  of i.i.d. samples (“evidence”) from the distribution. One standard way to quantify “best” is pick  $\theta$  is according to the *maximum likelihood principle*, which says that the best model is one that assigns the highest probability to the training dataset.<sup>1</sup>

<sup>1</sup> The maximum likelihood principle is a philosophical stance, not a consequence of some mathematical analysis.

$$\max_{\theta} \prod_{x^{(i)} \in S} p_{\theta}(x^{(i)}) \quad (12.1)$$

Because log is monotone, this is also equivalent to minimizing the *log likelihood*, which is a sum over training samples and thus similar in form to the training objectives seen so far in the book:

$$\max_{\theta} \sum_{x^{(i)} \in S} \log p_{\theta}(x^{(i)}) \quad (\text{log likelihood}) \quad (12.2)$$

Often one uses average log likelihood per datapoint, which means dividing (12.2) by  $\|S\|$ .

As in supervised learning, one has to keep track of training log-likelihood in addition to generalization, and choose among models that maximize it. In general such an optimization is computationally intractable for even fairly simple settings, and variants of gradient descent are used in practice.

Of course, the more important question is how well does the trained model learn the data distribution. Clearly, we need a notion of “goodness” for unsupervised learning that is analogous to *generalization* in supervised learning.

### 12.2.1 Notion of goodness for distribution learning

The most obvious notion of generalization follows from the log likelihood objective. The notion of generalization most analogous

to the one in supervised learning is to evaluate the log likelihood objective on *held-out* data: reserve some of the data for testing and compare the average log likelihood of the model on training data with that on test data.

**Example 12.2.1.** *The log likelihood objective makes sense for fitting any parametric model to the training data. For example, it is always possible to fit a simple Gaussian distribution  $\mathcal{N}(\mu, \sigma^2 I)$  to the training data in  $\mathbb{R}^d$ . The log-likelihood objective is*

$$\sum_i \frac{|x_i - \mu|^2}{\sigma^2},$$

which is minimized by setting  $\mu$  to  $\frac{1}{m} \sum_i x_i$  and  $\sigma^2$  to  $\sum_i \frac{1}{n} |x_i - \mu|^2$ .

Suppose we carry this out for the distribution of real-life images. What do we learn? The mean  $\mu$  will be the vector of average pixel values, and  $\sigma^2$  will correspond to the average variance per pixel. Thus a random sample from the learned distribution will look like some noisy version of the average pixel.

This example also shows that matching average log-likelihood for training and held-out data is insufficient for actually learning the distribution. The gaussian model only has  $d + 1$  parameters and simple  $\epsilon$ -cover arguments as in Chapter 6 show under fairly general conditions (such as coordinates of  $x_i$ 's being bounded) that if the number of training samples is moderately high then the log-likelihood on the average test sample is similar to that on the average training sample. However, the learned distribution may be nothing like the true distribution.

This is reminiscent of the situation in supervised learning whereby a nonsensical model —e.g., one that outputs random labels—has excellent generalization as well because it has similar loss on training as well as test data.

But how can we know that log likelihood objective is in principle capable of learning the distribution? The following theorem shows so.

**Theorem 12.2.2.** *Given enough training data, the  $\theta$  maximizing (12.2) minimizes the KL divergence  $KL(P||Q)$  where  $P$  is the true distribution and  $Q$  is the learnt distribution.*

*Proof.* This follows from

$$\begin{aligned} KL(P||Q) &= \mathbb{E}_{x \sim P} \left[ \log \frac{P(x)}{Q(x)} \right] \\ &= \mathbb{E}_{x \sim P} [\log P(x)] - \mathbb{E}_{x \sim P} [\log Q(x)]. \end{aligned}$$

Notice that  $\mathbb{E}_{x \sim P} [\log P(x)]$  is a constant that depends only upon the data distribution, and that computing log-likelihood using iid samples from  $P$  is like estimating the second term. We conclude that

given enough samples, minimizing  $KL(P||Q)$  amounts to maximising log likelihood up to an additive constant.  $\square$

Note that except for low-dimensional settings, the previous Theorem does not give any meaningful bounds on the number of training datapoints needed for proper learning.

### 12.3 Variational method

As sketched above, we are assuming a ground truth generative model  $p(x, h)$  and we are assuming we have samples of  $x$  obtained by generating pairs of  $(x, h)$  according to the ground truth and discarding the  $h$  part. The *variational method* tries to learn  $p(x)$  from such samples, where “variational” in the title refers to calculus of variations. It leverages *duality*, a widespread principle in math. The idea is to maintain a distribution  $q(h|x)$  as an attempt to model  $p(h|x)$  and improve a certain lower bound on  $p(x)$ . The key fact is the following. <sup>2</sup>

<sup>2</sup> See the blog post on [offconvex.org](http://offconvex.org) by Arora and Risteski on how algorithms try to use some form of gradient descent or local improvement to improve  $q(h|x)$ .

**Lemma 12.3.1** (ELBO Bound). *For any distribution  $q(h|x)$*

$$\log p(x) \geq \mathbb{E}_{q(h|x)}[\log(p(x, h))] + H[q(h|x)], \quad (12.3)$$

where  $H$  is the Shannon Entropy. (Note: equality is attained when  $q(h|x) = p(h|x)$ .)

*Proof.* Since

$$KL[q(h|x) || p(h|x)] = \mathbb{E}_{q(h|x)} \left[ \log \frac{q(h|x)}{p(h|x)} \right] \quad (12.4)$$

and  $p(x)p(h|x) = p(x, h)$  (Bayes’ Rule) we have:

$$KL[q(h|x)|p(h|x)] = \mathbb{E}_{q(h|x)} \left[ \log \frac{q(h|x)}{p(x, h)} \cdot p(x) \right] \quad (12.5)$$

$$= \underbrace{\mathbb{E}_{q(h|x)}[\log(q(h|x))]}_{-H(q(h|x))} - \mathbb{E}_{q(h|x)}[\log(p(x, h))] + \mathbb{E}_{q(h|x)}[\log p(x)] \quad (12.6)$$

But since KL divergence is always nonnegative, so we get:

$$\mathbb{E}_{q(h|x)}[\log(p(x))] - \mathbb{E}_{q(h|x)}[\log(p(x, h))] - H(q(h|x)) \geq 0 \quad (12.7)$$

which leads to the desired inequality since  $\log(p(x))$  is constant over  $q(h|x)$  and thus  $\mathbb{E}_{q(h|x)}[\log(p(x))] = \log p(x)$ .  $\square$

## 12.4 Autoencoders and Variational Autoencoder (VAEs)

Autoencoders find a compressed latent representation  $h$  of the datapoint  $x$  such that  $x$  can be approximately recovered from  $h$ . They can be defined in multiple ways by changing the formalization of what “approximate recovery” means.

In this section we formalize them using latent variable generative models. A popular instantiation of this in deep learning is *Variational Autoencoder (VAE)* <sup>3</sup>. As its name suggests two core classical ideas rest behind the design of VAEs: autoencoders – the original data  $x \in \mathbb{R}^n$  is mapped into a high-level descriptor  $z \in \mathbb{R}^d$  on a low dimensional (hopefully) meaningful manifold; variational inference – the objective to maximize is a lower bound on log-likelihood instead of the log-likelihood itself.

Recall that in density estimation we are given a data sample  $x_1, \dots, x_m$  and a parametric model  $p_\theta(x)$ , and our goal is to maximize the log-likelihood of the data:  $\max_\theta \sum_{i=1}^m \log p_\theta(x_i)$ . As a variational method, VAEs use the evidence lower bound (ELBO) as a training objective instead. For any distributions  $p$  on  $(x, z)$  and  $q$  on  $z|x$ , ELBO is derived from the fact that  $KL(q(z|x) || p(z|x)) \geq 0$

$$\log p(x) \geq \mathbb{E}_{q(z|x)}[\log p(x, z)] - \mathbb{E}_{q(z|x)}[\log q(z|x)] = ELBO \quad (12.8)$$

where equality holds if and only if  $q(z|x) \equiv p(z|x)$ . In the VAE setting, the distribution  $q(z|x)$  acts as the encoder, mapping a given data point  $x$  to a distribution of high-level descriptors, while  $p(x, z) = p(z)p(x|z)$  acts as the decoder, reconstructing a distribution on data  $x$  given a random seed  $z \sim p(z)$ . Deep learning comes in play for VAEs when constructing the aforementioned encoder  $q$  and decoder  $p$ . In particular,

$$q(z|x) = \mathcal{N}(z; \mu_x, \sigma_x^2 I_d), \quad \mu_x, \sigma_x = E_\phi(x) \quad (12.9)$$

$$p(x|z) = \mathcal{N}(x; \mu_z, \sigma_z^2 I_n), \quad \mu_z, \sigma_z = D_\theta(z), \quad p(z) = \mathcal{N}(z; 0, I_d) \quad (12.10)$$

where  $E_\phi$  and  $D_\theta$  are the encoder and decoder neural networks parameterized by  $\phi$  and  $\theta$  respectively,  $\mu_x, \mu_z$  are vectors of corresponding dimensions, and  $\sigma_x, \sigma_z$  are (nonnegative) scalars. The particular choice of Gaussians is not a necessity in itself for the model and can be replaced with any other relevant distribution. However, Gaussians provide, as is often the case, computational ease and intuitive backing. The intuitive argument behind the use of Gaussian distributions is that under mild regularity conditions every distribution can be approximated (in distribution) by a mixture of Gaussians. This follows from the fact that by approximating the CDF of a distribution by step functions one obtains an approximation in distribution by

a mixture of constants, i.e. mixture of Gaussians with  $\approx 0$  variance. The computational ease, on the other hand, is more clearly seen in the training process of VAEs.

#### 12.4.1 Training VAEs

As previously mentioned, the training of variational autoencoders involves maximizing the RHS of (12.8), the ELBO, over the parameters  $\phi, \theta$  under the model described by (12.9), (12.10). Given that the parametric model is based on two neural networks  $E_\phi, D_\theta$ , the objective optimization is done via gradient-based methods. Since the objective involves expectation over  $q(z|x)$ , computing an exact estimate of it, and consequently its gradient, is intractable so we resort to (unbiased) gradient estimators and eventually use a stochastic gradient-based optimization method (e.g. SGD).

In this section, use the notation  $\mu_\phi(x), \sigma_\phi(x) = E_\phi(x)$  and  $\mu_\theta(z), \sigma_\theta(z) = D_\theta(z)$  to emphasize the dependence on the parameters  $\phi, \theta$ . Given training data  $x_1, \dots, x_m \in \mathbb{R}^n$ , consider an arbitrary data point  $x_i, i \in [m]$  and pass it through the encoder neural network  $E_\phi$  to obtain  $\mu_\phi(x_i), \sigma_\phi(x_i)$ . Next, sample  $s$  points  $z_{i1}, \dots, z_{is}$ , where  $s$  is the batch size, from the distribution  $q(z|x = x_i) = \mathcal{N}(z; \mu_\phi(x_i), \sigma_\phi(x_i)^2 I_d)$  via the reparameterization trick<sup>4</sup> by sampling  $\epsilon_1, \dots, \epsilon_s \sim \mathcal{N}(0, I_d)$  from the standard Gaussian and using the transformation  $z_{ij} = \mu_\phi(x_i) + \sigma_\phi(x_i) \cdot \epsilon_j$ . The reason behind the reparameterization trick is that the gradient w.r.t. parameter  $\phi$  of an unbiased estimate of expectation over a general distribution  $q_\phi$  is not necessarily an unbiased estimate of the gradient of expectation. This is the case, however, when the distribution  $q_\phi$  can separate the parameter  $\phi$  from the randomness in the distribution, i.e. it's a deterministic transformation that depends on  $\phi$  of a parameter-less distribution. With the  $s$  i.i.d. samples from  $q(z|x = x_i)$  we obtain an unbiased estimate of the objective ELBO

$$\sum_{j=1}^s \log p(x_i, z_{ij}) - \sum_{j=1}^s \log q(z_{ij}|x_i) = \sum_{j=1}^s [\log p(x_i|z_{ij}) + \log p(z_{ij}) - \log q(z_{ij}|x_i)] \quad (12.11)$$

Here the batch size  $s$  indicates the fundamental tradeoff between computational efficiency and accuracy in estimation. Since each of the terms in the sum in (12.11) is a Gaussian distribution, we can write the ELBO estimate explicitly in terms of the parameter-dependent  $\mu_\phi(x_i), \sigma_\phi(x_i), \mu_\theta(z_{ij}), \sigma_\theta(z_{ij})$  (while skipping some constants). A single term for  $j \in [s]$  is given by

$$-\frac{1}{2} \left[ \frac{\|x_i - \mu_\theta(z_{ij})\|^2}{\sigma_\theta(z_{ij})^2} + n \log \sigma_\theta(z_{ij})^2 + \|z_{ij}\|^2 - \frac{\|z_{ij} - \mu_\phi(x_i)\|^2}{\sigma_\phi(x_i)^2} - d \log \sigma_\phi(x_i)^2 \right] \quad (12.12)$$

Notice that (12.12) is differentiable with respect to all the components  $\mu_\phi(x_i), \sigma_\phi(x_i), \mu_\theta(z_{ij}), \sigma_\theta(z_{ij})$  while each of these components, being an output of a neural network with parameters  $\phi$  or  $\theta$ , is differentiable with respect to the parameters  $\phi$  or  $\theta$ . Thus, the tractable gradient of the batch sum (12.11) w.r.t.  $\phi$  (or  $\theta$ ) is, *due to the reparameterization trick*, an unbiased estimate of  $\nabla_\phi ELBO$  (or  $\nabla_\theta ELBO$ ) which can be used in any stochastic gradient-based optimization algorithm to maximize the objective ELBO and train the VAE.

## 12.5 Normalizing Flows

The limitation of VAE is that instead of direct log likelihood, it optimizes a lower bound to it. Ideally we would want to get around this limitation while staying with a deep model with sophisticated representation capability. (The simple Gaussian fit as described at the start of the chapter also optimizes log likelihood directly but it cannot represent complicated distributions.) *Normalizing flows* can do this,

The idea in *Normalizing Flows* is to make the deep net *invertible*. Specifically, it computes a function  $f_\theta: \mathbb{R}^d \rightarrow \mathbb{R}^d$  that is parametrized by trainable parameter vector  $\theta$  and maps image  $x$  to its representation  $h = f_\theta(x)$  (note: both have the same dimension). Importantly,  $f$  is an invertible map (i.e., one-to-one and onto) and differentiable (or almost everywhere differentiable). The advantage of such a transformation is that it gives a clear connection between the probability densities of  $x$  and  $h$ . In generative models  $h$  is assumed to have some prescribed probability density  $\mu(h)$ , usually uniform gaussian. Via the invertible map, this translates to a density  $\rho(\cdot)$  on  $x$  given by

$$\rho(x) = \mu(f(x)) \left| \det\left(\frac{\partial f}{\partial x}\right) \right| \quad (12.13)$$

where  $J$  is the Jacobian of  $f$  namely, whose  $(i, j)$  entry is  $\partial f(x)_i / \partial x_j$  and  $\det(\cdot)$  denotes determinant of the matrix. This exact expression for likelihood of the training datapoints allows usual gradient-based training.

Which raises the question: how does one constrain nets to be invertible? Note that it suffices to constrain individual layers to be invertible, because the overall Jacobian is the composition of layer Jacobians.<sup>5</sup> To make layers invertible one often uses a variant of the following trick from the models NICE<sup>6</sup> and Real NVP<sup>7</sup>. If  $z^l$  is the input to layer  $l$  and  $z^{l+1}$  its output, then identify a special set of coordinates  $A$  in  $z^l$  and  $z^{l+1}$  and impose the restriction (where  $z_A$  denotes

<sup>5</sup> Since  $\det(AB) = \det(A)\det(B)$  the determinant of the net Jacobian is the product of the determinants of the layers.

<sup>6</sup> Laurent Dinh, David Krueger, and Yoshua Bengio. NICE: Nonlinear Independent Component Analysis. *Proc. ICLR*, 2015

<sup>7</sup> Laurent Dinh, Jascha Sohl-Dickstein, and Samy Bengio. Density Estimation using Real NVP. *Proc. ICLR*, 2017



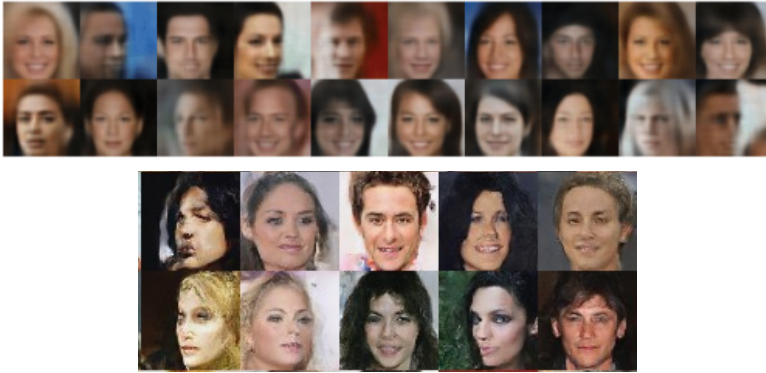


Figure 12.3: Faces in the top row were produced by a VAE based method and those in the second row by RealNVP using normalizing flows. VAE is known for producing blurry images. RealNVP’s output is much better, but still has visible artifacts.

portion of  $z$  in the coordinates given by  $A$ , and  $B$  is shorthand for  $\bar{A}$ ).

$$z_A^{l+1} = z_A^l \quad (12.14)$$

$$z_B^{l+1} = z_B^l \odot h_\theta(z_A^l) + s_\theta(z_A^l) \quad (12.15)$$

where  $\odot$  denotes component-wise product and  $h_\theta()$  is a function whose each output is nonnegative, with a convenient choice being to make it  $\exp(r_\theta(z_A^l))$  for some other function  $r_\theta()$ .

This layer is invertible because given  $z^{l+1}$  one can recover  $z^l$  as follows:

$$z_A^l = z_A^{l+1} \quad (12.16)$$

$$z_B^l = (z_B^{l+1} - s_\theta(z_A^l)) \odot h_\theta(z_A^l) \quad (12.17)$$

Note that the choice of  $A, B$  can change from layer to layer, so all coordinates may get updated as they go through multiple layers. Furthermore, denoting by  $z|_A$  the portion of the layer vector on coordinates  $A$ , the Jacobian for the layer mapping is lower triangular. Hence the determinant is the product of the diagonal entries.

$$\frac{\partial}{\partial z^l} z^{l+1} = \begin{pmatrix} I_{|A| \times |A|} & 0 \\ \frac{\partial}{\partial z^l|_A} z^{l+1}|_B & \text{diag}(h_\theta(z_A^l)) \end{pmatrix}$$

Normalizing flows can be extended to convolutional nets by restricting the convolutions are  $1 \times 1$ . Then convolutional filters just involve scalings of channel values, and the corresponding Jacobian is a diagonal nonzero matrix. Also the split of coordinates into  $A$  and  $B$  split can be done within channels as well. This is the main idea in GLOW model<sup>8</sup>, which can generate better images than its predecessors.

More recent auto-regressive models such as PixelCNN are capable of producing very realistic-looking images from random seeds. However, they do not fit into the distribution learning paradigm described above so we do not discuss them here.

<sup>8</sup> Diederik P. Kingma and Prafulla Dhariwal. GLOW: Generative Flow with Invertible  $1 \times 1$  convolutions. *Proc. Neurips*, 2019

**Problem 12.5.1.** Let  $(z_1, z_2, z_3, z_4)$  be distributed as a standard Gaussian  $\mathcal{N}(0, I)$  in  $\mathbb{R}^4$ . Let  $f : \mathbb{R}^4 \rightarrow \mathbb{R}^4$  be an invertible function which maps  $(z_1, z_2, z_3, z_4)$  to  $(z_1, z_2, e^{a_0} z_3 + a_1 z_1^2 + a_2 z_2^2, e^{b_0} z_4 + b_1 z_1^2 + b_2 z_2^2)$  for some coefficients  $a_0, a_1, a_2, b_0, b_1, b_2 \in \mathbb{R}$ . Compute the probability density function of  $f(z_1, z_2, z_3, z_4)$ .

## Generative Adversarial Nets

Chapter 12 described some classical approaches to generative models, which are often trained using a log-likelihood approach. We also saw that they often do not suffice for high-fidelity learning of complicated distributions such as the distribution of real-life images. *Generative Adversarial Nets (GANs)* is an approach that generates more realistic samples. It relies upon power of deep nets at discriminative tasks. For convenience throughout this chapter we assume the model is trying to generate images.

Before developing the theory of GANs we survey various notions of how to test similarity of two distributions, because that discussion feeds directly into the setup used in GANs. This is relevant because Example 12.2.1 illustrated how the the notion of *generalization* from supervised learning can lead us astray when it comes to reasoning about correctness of distribution learning.

### 13.1 Distance between Distributions

How can we measure how different two distributions  $P$  and  $Q$  are? If we have access to a formula for computing the density of each distribution, then one can compute an  $f$ -divergence for any suitable  $f: \mathfrak{R} \rightarrow \mathfrak{R}$  that is convex.

$$D_f(P||Q) = \int f\left(\frac{P(x)}{Q(x)}\right)Q(x)dx \quad (13.1)$$

**Problem 13.1.1.** (i) Show that  $f$ -divergence is nonnegative. (ii) Show that the  $f$ -divergence for  $f(t) = t \log t$  coincides with  $KL(P||Q)$ . (iii) Show that the  $f$ -divergence for  $f = \frac{1}{2}|t - 1|$  coincides with total variation (or  $\ell_1$ ) distance:  $|P - Q|_1 = \int |P(x) - Q(x)|dx$ .

However, in practice one often doesn't have a formula for the probability density and must estimate distance from samples. A natural idea is to compare the expectation of a class of *test* functions on the two distributions.

**Example 13.1.2.** The expectation  $\mathbb{E}_P[x]$  is the mean of the distribution  $P$ , and expectation of monomials of form  $x_{i_1}x_{i_2}\cdots x_{i_k}$  constitutes the  $k$ th moment. Moments can be estimated from samples (under fairly general conditions) and the difference of means for distributions  $P, Q$  can be seen as some measure of their difference.

Unfortunately, accurate estimation of all higher moments for multivariate distributions gets expensive with respect to both computation time and sample complexity. This motivates a notion of distance that arises in *transportation metrics*.<sup>1</sup> If  $\mathcal{F}$  is a class of functions, then define the distance between  $P$  and  $Q$  using the highest different in expectation achievable over functions in  $\mathcal{F}$ .

$$d(P, Q) = \sup_{f \in \mathcal{F}} |\mathbb{E}_{x \sim P} f(x) - \mathbb{E}_{x \sim Q} f(x)| \quad (13.2)$$

For example  $\mathcal{F}$  could be polynomials of degree at most  $k$ .

**Problem 13.1.3.** Show that the distance defined above satisfies triangle inequality.

<sup>1</sup> Please look up Wasserstein metrics and Earth-Mover distance on the internet.

## 13.2 Introducing GANs

Generative Adversarial Nets (Goodfellow et al.<sup>2</sup>) seeks to learn generative models via the definition in (13.2). Specifically, one tries to train a generative model  $G$  (which, as in Chapter 12) produces an image  $G(u)$  using random vector  $u$ . This yields a distribution on images, and we check the quality of this distribution by using a class  $\mathcal{F}$  of deep nets (with a fixed size and architecture) and trying to compute the distance in (13.2) by trying to find a net that is a maximiser of the expression. Now we spell out the main components of GANs.

*Idea 1:* Since deep nets are good at recognizing images —e.g., distinguishing pictures of people from pictures of cats—why not let a deep net be the judge of the outputs of  $G$ ?

More concretely, let  $P_{real}$  be the distribution over real images, and  $P_{synth}$  the synthetic ones (i.e., the distribution of  $G(h)$  when  $h$  is a random seed). We could try to train a discriminator deep net  $D$  that maps images to numbers in  $[0, 1]$  and tries to discriminate between these distributions in the following sense. Its expected output  $E_x[D(x)]$  as high as possible when  $x$  is drawn from  $P_{real}$  and as low as possible when  $x$  is drawn from  $P_{synth}$ . This is just supervised learning (i.e., regression) with two labels, and can be done using backpropagation. If the two distributions are identical then of course no such deep net can exist, and so the training will end in failure. If, on the other hand, we are able to train a good

<sup>2</sup>

discriminator deep net—one whose average output is noticeably different between real and synthetic samples—then this is proof positive that the two distributions are different.<sup>3</sup>

*Idea 2: If a good discriminator net has been trained, use it to provide “gradient feedback” that improves the generative model.*

The natural goal for the generator is to make  $E_h[D(G(h))]$  as high as possible, because that means it is doing better at fooling the discriminator  $D$ . So if we fix  $D$  the natural way to improve  $G$  is to pick a few random seeds  $h$ , and slightly adjust the trainable parameters of  $G$  to increase this objective. Note that this gradient computation involves backpropagation through the composed net  $D(G(\cdot))$ .

*Idea 3: Turn the training of the generative model and the accompanying discriminator net into a game of many moves (i.e., rounds of parameter updates).*

Each move for the discriminator consists of taking a few samples from  $P_{real}$  and  $P_{synth}$  and improving its ability to discriminate between them. Each move for the generator consists of producing a few samples from  $P_{synth}$  and updating its parameters so that  $E_u[D(G(h))]$  goes up a bit.

Notice, the discriminator always uses the generator as a black box—i.e., never examines its internal parameters—whereas the generator needs the discriminator’s parameters to compute its gradient direction. Specifically, the gradient for  $G$  is computed by backpropagating through  $D$ . Also, the generator does not ever use real images from  $P_{real}$  for its computation. (Though of course it does rely on the real images indirectly since the discriminator is trained using them.)

One can fill in the above framework in multiple ways. The most obvious is that the generator could try to maximize  $E_u[f(D(G(h)))]$  where  $f$  is some increasing function. (We call this the “measuring function.”) Concretely, if  $D, G$  are deep nets with specified architecture and whose number of parameters is fixed in advance by the algorithm designer, then the training objective is:

$$\min_G \max_D E_{x \sim P_{real}} [f(D(x))] + E_h [f(1 - D(G(h)))]. \quad (13.3)$$

**Effect of  $f$ :** The measuring function has the effect of giving different importance to different samples. Goodfellow et al. originally used  $f(x) = \log(x)$ , which, since the derivative of  $\log x$  is  $1/x$ , implicitly gives much more importance to synthetic data  $G(u)$  where the discriminator outputs very low values  $D(G(h))$ . In other words, using  $f(x) = \log x$  makes the training more sensitive to instances which the

<sup>3</sup> There is an in-between case, whereby the distributions are different but the discriminator net doesn’t detect a difference. This case is going to be important in the story very soon.

discriminator finds terrible than to instances which the discriminator finds so-so. By contrast,  $f(x) = x$  gives the same importance to all samples and results in *Wasserstein GAN*.

**Problem 13.2.1.** *Show that if the deep net has unbounded capacity then*

### 13.2.1 Game-theoretic interpretation and implications for trainings

A serious practical difficulty in implementing training as above is that it can be oscillatory, meaning the above objective can go up and down. This is unlike usual deep net training, where training (at least in cases where it works) steadily improves the objective. The reason is that implicitly the discriminator and generator are playing a 2-person zero sum game<sup>4</sup> where their “moves” are the two circuits  $D, G$  and the payoff from generator (minimizer) to discriminator (maximizer) is the loss. Thus generator is picking moves to minimize the following payoff

$$\max_D E_{x \sim P_{real}} [f(D(x))] + E_h [f(1 - D(G(h)))]$$

whereas discriminator is maximising

$$\min_G E_{x \sim P_{real}} [f(D(x))] + E_h [f(1 - D(G(h)))].$$

Such games do not always have an *equilibrium*, i.e., a point where both players are reacting optimally to the other and thus lack an incentive to change. (It is akin to a saddle point in optimization.)

Although equilibrium may not exist when viewed as a two-person game, of course during training both players are under the control of the training algorithm. Thus suitable modifications to gradient-based training could conceivably allow convergence to some solution even though it is not an equilibrium. (For example, even if  $D$  is not optimal response to the current  $G$ , gradients may not allow a way to improve on the current  $D$ .) An extensive list of papers present such ideas; some examples are: [NEED SOME REFERENCES HERE](#)

<sup>4</sup> Please read up about zero sum games online, including the famous Min-Max Theorem about equilibria.

## 13.3 Do GANs learn the distribution?

Section 12.2.1 discussed complications arising when we try to learn distributions from finite samples. In the GAN setting the training objective (13.3) was described using the full distribution but of course in practice the discriminator  $D$  is discriminating between finite samples of the two distributions.

**Example 13.3.1.** *Since the objective (13.3) allows maximization over all neural nets  $D$  of the allowed architecture, even two different samples from*

the same distribution can look very different to a neural net. For example if we take two finite samples from the  $d$ -dimensional Gaussian  $\mathcal{N}(0, I)$  then even if the samples have size say  $d^3$ , they are distinguishable by a deep net that is somewhat larger than  $d^3$ . The reason is that the samples are two discrete sets in which all pairs of points are almost orthogonal (with high probability). While this is an artificial example, it is the case that in real-life GAN training, the objective does not usually drop to zero—the net is indeed able to somewhat distinguish the samples from the two distributions.

For several years, GANs practitioners felt that any problems in training arise from using too small a training set. We now describe theoretical analysis showing that the quality of the learnt distribution is inherently limited by the *representational capacity* of the discriminator.

But as is usual when discussing generalization, let us assume the net has some finite size, say  $N$ , and the sample size of the distributions is large enough for nets of size  $N$  to generalize. The following two problems are drawn from <sup>5</sup>.

**Problem 13.3.2.** *Suppose the training loss is at most  $\epsilon_1$ , the loss is  $C$ -Lipschitz and the parameter vector has  $\ell_2$ -norm at most  $L$ . Show that a discriminator trained on at least ?? samples from the two distributions has test loss at most  $\epsilon_1 + \epsilon_2$  on the full distribution. (Hint: Use the results of Chapter 6)*

Does this imply that GANs actually learn the distribution? No, just as in Example 12.2.1: generalization only implies training and test loss are close, not that the distributions are close.

**Problem 13.3.3.** *Under the same conditions as Problem 13.3.2 show that if the distribution  $P_{\text{synth}}$  consists of ?? random iid samples from  $P_{\text{real}}$  then no discriminator can achieve test loss more than  $\epsilon_1 + \epsilon_2$  when comparing the distributions  $P_{\text{synth}}$  and  $P_{\text{real}}$ . (NB: This also implies the test loss is small whp on finite samples.)*

In the previous problem, let's think of  $P_{\text{real}}$  as the distribution on all real-life images. Then  $P_{\text{synth}}$  is quite different from  $P_{\text{real}}$ : it is the uniform distribution on some small set of real-life images. Nevertheless no discriminator (whose size, norm and Lipschitz constant are suitably upper bounded) can distinguish between the two distributions.

### 13.3.1 Experimental verification: Birthday Paradox Test

The theory above suggests that GANs trained using discriminators of a certain “capacity”(in the sense of generalization theory)

<sup>5</sup> Sanjeev Arora, Rong Ge, Yingyu Liang, Tengyu Ma, and Yi Zhang. Generalization and Equilibrium in Generative Adversarial Nets (GANs). *Proc. ICML*, 2017

have solutions with low training and test error where  $P_{synth}$  is supported on a small set of images and is quite different from  $P_{real}$ ). This phenomenon of the GAN ending up with a synthetic distribution consisting of a small set of images is called *mode collapse* and was earlier believed to be a result of either failed training or using a training set of real images that is too small. The results above suggested a different explanation of mode collapse.

This raised a question whether we can detect such model collapse in real-life GANs. In other words, estimate how many “distinct” images it can produce. At first glance, such an estimation seems very difficult. After all, automated/heuristic measures of image similarity can be easily fooled, and we humans surely don’t have enough time to go through millions or billions of images, right?

Luckily, a crude estimate is possible using the simple birthday paradox, a staple of undergrad discrete math.

**Problem 13.3.4** (Birthday paradox). *Consider a uniform distribution on a set of size  $N$ . Show that a random sample of size  $\sqrt{N}$  contains a duplicate probability at least  $1 - 1/e$ . (The name for this paradox comes from its implication that if you put 23 random people in a room, then the odds are good that two of them have the same birthday is significant.)*

Going back to implications for GANs. Imagine for argument’s sake that  $P_{real}$  is the distribution on pictures of faces. What is the number of distinct faces possible? It feels like the diversity of human faces is rather large, because all of us know tens of thousands of faces (including those encountered in the news) and don’t see any unrelated *doppelgangers*; only identical twins. This suggests that the total number of distinct faces possible is large, certainly hundreds of millions, and maybe many billions.

In the GAN setting, the distribution is continuous, not discrete. Thus our proposed birthday paradox test for GANs<sup>6</sup> is as follows.

(a) Pick a sample of size  $s$  from the generated distribution. (b) Use an automated measure of image similarity to flag the 20 (say) most similar pairs in the sample. (c) Visually inspect the flagged pairs and check for images that a human would consider near-duplicates. (d) Repeat.

If this test reveals that samples of size  $s$  have duplicate images with good probability, then suspect that the distribution has support size about  $s^2$ .

**Problem 13.3.5.** *Put in formal calculation for birthday paradox test.*

<sup>6</sup> Sanjeev Arora and Yi Zhang. Theoretical Limitations of Encoder-Decoder GANs Architectures. *Proc. ICLR*, 2018



## Self-supervised Learning

Semantic representations (aka *semantic embeddings*) of complicated data types (e.g. images, text, video) have become central in machine learning, and also crop up in machine translation, language models, GANs, domain transfer, etc. These involve learning a representation function  $f$  such that  $f(x)$  is a compact and “high level” representation of datapoint  $x$  —meaning it retains semantic information while discarding low level details — e.g., the colors of individual pixels in an image. The test of a good representation is that it should greatly simplify solving new classification tasks, by allowing them to be solved via linear classifiers (or other low-complexity classifiers) using small amounts of labeled data.

Researchers are most interested in unsupervised representation learning using unlabeled data. A popular early example is *word embeddings* which used simple linear algebra <sup>1</sup> and proved useful in information retrieval for several decades. More recent word embeddings such as word2vec<sup>2</sup> became the inspiration for semantic embeddings of diverse data types such as molecules, social networks, images, text etc.

In this chapter we encounter *self-supervised* learning, a family of methods for learning good representations. Working with unlabeled data, the learner defines a learning objective for finding a good representation function. An important difference from learning paradigms studied elsewhere in the book is that the training and test tasks are different, and hence the notion of generalization does not capture the final goal of learning. This is an example of *training on task A to later do well on task B*, which one imagines is actually an important aspect of intelligent behavior.

<sup>1</sup> LSI paper

<sup>2</sup> word2vec paper; see wikipedia page of word2vec for links to other similar algorithms



## Adversarial Examples

We provide the definition of adversarial example problem for deep neural network,

**Definition 15.0.1** (Adversarial example, finding the closest adversarial example). *Given a function  $f : \mathbb{R}^d \rightarrow \{0, 1\}$  and an input  $x_0 \in \mathbb{R}^d$ , the goal is to output the smallest  $r > 0$  such that  $\exists x' \in B(x_0, r, \ell_\infty)$  satisfying  $f(x') \neq f(x_0)$ .*

An alternative formulation could be the following

**Definition 15.0.2.** *Adversarial example, finding the largest safe region* Given a function  $f : \mathbb{R}^d \rightarrow \{0, 1\}$  and an input  $x_0 \in \mathbb{R}^d$ , the goal is to output the largest  $r > 0$  such that  $\forall x' \in B(x_0, r, \ell_\infty)$  satisfying  $f(x') = f(x_0)$ .

The above formulation is the exact version, in practice, we might not always need the exact solution. It is natural to define the approximate version of the problem:

**Definition 15.0.3.** *Given a function  $f : \mathbb{R}^d \rightarrow \{0, 1\}$  and an input  $x_0 \in \mathbb{R}^d$ . Let  $r^*$  denote the largest  $r > 0$  such that  $\forall x' \in B(x_0, r, \ell_\infty)$  satisfying  $f(x') = f(x_0)$ . The goal is to output  $r$  such that  $\alpha \cdot r^* \leq r \leq r^*$ , where  $\alpha \in (0, 1)$  is the approximation ratio.*

Katz, Barrett, Dill, Julian, Kochenderfer, CAV'17 [cite] proved that, assuming  $P \neq NP$ , there is no polynomial time that computes  $r$  exactly. Further, Weng et al. [cite] shows that assuming ETH, there is no polynomial that gives an  $1/\log d$  approximation.



## Examples of Theorems, Proofs, Algorithms, Tables, Figures

In this chapter, Zhao provide examples of many things, like Theorems, Lemmas, Algorithms, Tables, and Figures. If anyone has question, feel free to contact Zhao directly.

### 16.1 Example of Theorems and Lemmas

We provide some examples

**Theorem 16.1.1** (*d*-dimension sparse Fourier transform). *There is an algorithm (procedure FOURIERSPARSERECOVERY in Algorithm 2) that runs in ??? times and outputs ??? such that ???.*

Note that, usually, if we provide the algorithm of the Theorem/Lemma. Theorem should try to ref the corresponding Algorithm.

For the name of Theorem/Lemma/Corollary ..., let us only capitalize the first word,

**Lemma 16.1.2** (Upper bound on the gradient). *Blah blah.*

**Problem 16.1.3.** *This is how you put in a problem. It inherits chapter and section numbers.*

**Theorem 16.1.4** (Main result).

### 16.2 Example of Long Equation Proofs

We can rewrite  $\|Ax' - b\|_2^2$  in the following sense,

$$\begin{aligned} \|Ax' - b\|_2^2 &= \|Ax' - Ax^* + AA^\dagger b - b\|_2^2 \\ &= \|Ax^* - Ax'\|_2^2 + \|Ax^* - b\|_2^2 \\ &= \|Ax^* - Ax'\|_2^2 + \text{OPT}^2 \end{aligned}$$

where the first step follows from  $x^* = A^\dagger b$ , the second step follows from Pythagorean Theorem, and the last step follows from  $\text{OPT} := \|Ax^* - b\|_2$ .

### 16.3 Example of Algorithms

Here is an example of algorithm. Usually the algorithm should ref some Theorem/Lemma, and also the corresponding Theorem/Lemma should ref back. This will be easier to verify the correctness.

---

#### Algorithm 2 Fourier Sparse Recovery Algorithm

---

```

1: procedure FOURIERSPARSERECOVERY( $x, n, k, \mu, R^*$ ) ▷
   Theorem 16.1.1
2:   Require that  $\mu = \frac{1}{\sqrt{k}} \|\hat{x}_{-k}\|_2$  and  $R^* \geq \|\hat{x}\|_\infty / \mu$ 
3:    $H \leftarrow 5, v \leftarrow \mu R^* / 2, y \leftarrow \vec{0}$ 
4:   Let  $\mathcal{T} = \{\mathcal{T}^{(1)}, \dots, \mathcal{T}^{(H)}\}$  where each  $\mathcal{T}^{(h)}$  is a list of i.i.d.
   uniform samples in  $[p]^d$ 
5:   while true do
6:      $v' \leftarrow 2^{1-H} v$ 
7:      $z \leftarrow \text{LINFINITYREDUCE}(\{x_t\}_{t \in \mathcal{T}})$ 
8:     if  $v' \leq \mu$  then return  $y + z$  ▷ We found the solution
9:      $y' \leftarrow \vec{0}$ 
10:    for  $f \in \text{supp}(y + z)$  do
11:       $y'_f \leftarrow \Pi_{0.6v}(y_f + z_f)$  ▷ We want  $\|\hat{x} - y'\|_\infty \leq v$  and the
      dependence between  $y'$  and  $\mathcal{T}$  is under control
12:    end for
13:     $y \leftarrow y', v \leftarrow v/2$ 
14:  end while
15: end procedure

```

---

## 16.4 Example of Figures

We should make sure all the pictures are plotted by the same software. Currently, everyone feel free to include their own picture. Zhao will re-plot the picture by tikz finally.

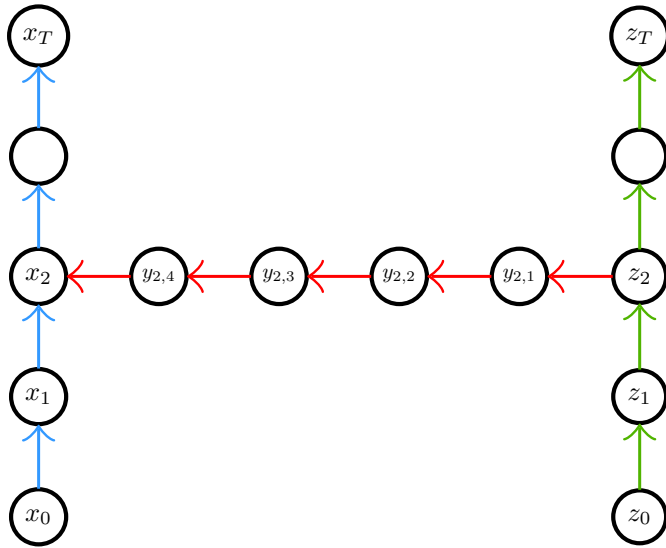


Figure 16.1: A chasing sequence



## 16.5 Example of Tables

Reference	Samples	Time	Filter	RIP
[GMS05]	$k \log^{O(d)} n$	$k \log^{O(d)} n$	Yes	No
[CT06]	$k \log^6 n$	$\text{poly}(n)$	No	Yes
[RV08]	$k \log^2 k \log(k \log n) \log n$	$\tilde{O}(n)$	No	Yes
[HIKP12]	$k \log^d n \log(n/k)$	$k \log^d n \log(n/k)$	Yes	No
[CGV13]	$k \log^3 k \log n$	$\tilde{O}(n)$	No	Yes
[IK14]	$2^{d \log d} k \log n$	$\tilde{O}(n)$	Yes	No
[Bou14]	$k \log k \log^2 n$	$\tilde{O}(n)$	No	Yes
[HR16]	$k \log^2 k \log n$	$\tilde{O}(n)$	No	Yes
[Kap16]	$2^{d^2} k \log n$	$2^{d^2} k \log^{d+O(1)} n$	Yes	No
[KVZ19]	$k^3 \log^2 k \log^2 n$	$k^3 \log^2 k \log^2 n$	Yes	Yes
[NSW19]	$k \log k \log n$	$\tilde{O}(n)$	No	No

Table 16.1: We ignore the  $O$  for simplicity. The  $\ell_\infty/\ell_2$  is the strongest possible guarantee, with  $\ell_2/\ell_2$  coming second,  $\ell_2/\ell_1$  third and exactly  $k$ -sparse being the weaker. We also note that all [RV08, CGV13, Bou14, HR16] obtain improved analyses of the Restricted Isometry property; the algorithm is suggested and analyzed (modulo the RIP property) in [BD08]. The work in [HIKP12] does not explicitly state the extension to the  $d$ -dimensional case, but can easily be inferred from the arguments. [HIKP12, IK14, Kap16, KVZ19] work when the universe size in each dimension are powers of 2.

## 16.6 Exercise

This section provides several examples of exercises.

---

**Exercises**

*Exercise 16.6-1:* Solve the following equation for  $x \in \mathbb{C}$ , with  $\mathbb{C}$  the set of complex numbers:

$$5x^2 - 3x = 5 \quad (16.1)$$

*Exercise 16.6-2:* Solve the following equation for  $x \in \mathbb{C}$ , with  $\mathbb{C}$  the set of complex numbers:

$$7x^3 - 2x = 1 \quad (16.2)$$


---



# Bibliography

- [ADG<sup>+</sup>16] Marcin Andrychowicz, Misha Denil, Sergio Gomez, Matthew W Hoffman, David Pfau, Tom Schaul, and Nando de Freitas. Learning to learn by gradient descent by gradient descent. In *Advances in Neural Information Processing Systems*, 2016.
- [ADH<sup>+</sup>19a] Sanjeev Arora, Simon Du, Wei Hu, Zhiyuan Li, and Ruosong Wang. Fine-grained analysis of optimization and generalization for overparameterized two-layer neural networks. In *International Conference on Machine Learning*, pages 322–332, 2019.
- [ADH<sup>+</sup>19b] Sanjeev Arora, Simon S Du, Wei Hu, Zhiyuan Li, Ruslan Salakhutdinov, and Ruosong Wang. On exact computation with an infinitely wide neural net. *arXiv preprint arXiv:1904.11955*, 2019.
- [ADH<sup>+</sup>19c] Sanjeev Arora, Simon S Du, Wei Hu, Zhiyuan Li, Ruslan Salakhutdinov, and Ruosong Wang. On exact computation with an infinitely wide neural net. *arXiv preprint arXiv:1904.11955*, 2019.
- [ADL<sup>+</sup>19] Sanjeev Arora, Simon S Du, Zhiyuan Li, Ruslan Salakhutdinov, Ruosong Wang, and Dingli Yu. Harnessing the power of infinitely wide deep nets on small-data tasks. *arXiv preprint arXiv:1910.01663*, 2019.
- [AGL<sup>+</sup>17] Sanjeev Arora, Rong Ge, Yingyu Liang, Tengyu Ma, and Yi Zhang. Generalization and Equilibrium in Generative Adversarial Nets (GANs). *Proc. ICML*, 2017.
- [AGNZ18] Sanjeev Arora, Rong Ge, Behnam Neyshabur, and Yi Zhang. Stronger generalization bounds for deep nets via a compression approach. In *Proc. ICML 2018*, pages 254–263, 2018.
- [AZ18] Sanjeev Arora and Yi Zhang. Theoretical Limitations of Encoder-Decoder GANs Architectures. *Proc. ICLR*, 2018.

- [BBV16] Afonso S Bandeira, Nicolas Boumal, and Vladislav Voroninski. On the low-rank approach for semidefinite programs arising in synchronization and community detection. In *Conference on learning theory*, pages 361–382, 2016.
- [BD08] Thomas Blumensath and Mike E Davies. Iterative thresholding for sparse approximations. *Journal of Fourier analysis and Applications*, 14(5-6):629–654, 2008.
- [Ber24] Sergei Bernstein. On a modification of chebyshev’s inequality and of the error formula of laplace. *Ann. Sci. Inst. Sav. Ukraine, Sect. Math*, 1(4):38–49, 1924.
- [BH89] Pierre Baldi and Kurt Hornik. Neural networks and principal component analysis: Learning from examples without local minima. *Neural networks*, 2(1):53–58, 1989.
- [BM02] Peter L Bartlett and Shahar Mendelson. Rademacher and gaussian complexities: Risk bounds and structural results. *Journal of Machine Learning Research*, 3(Nov):463–482, 2002.
- [BM03] P. L. Bartlett and S. Mendelson. Rademacher and Gaussian complexities: Risk bounds and structural results. *Journal of Machine Learning Research*, 2003.
- [BNS16a] Srinadh Bhojanapalli, Behnam Neyshabur, and Nati Srebro. Global optimality of local search for low rank matrix recovery. In *Advances in Neural Information Processing Systems*, pages 3873–3881, 2016.
- [BNS16b] Srinadh Bhojanapalli, Behnam Neyshabur, and Nati Srebro. Global optimality of local search for low rank matrix recovery. In *Advances in Neural Information Processing Systems (NIPS)*, pages 3873–3881, 2016.
- [Bou14] Jean Bourgain. An improved estimate in the restricted isometry problem. In *Geometric Aspects of Functional Analysis*, pages 65–70. Springer, 2014.
- [BR89] Avrim Blum and Ronald L Rivest. Training a 3-node neural network is np-complete. In *Advances in neural information processing systems*, pages 494–501, 1989.
- [Bre67] L. M. Bregman. The relaxation method of finding the common point of convex sets and its application to the solution of problems in convex programming. *USSR computational mathematics and mathematical physics*, 1967.

- [BT03] A. Beck and M. Teboulle. Mirror descent and nonlinear projected subgradient methods for convex optimization. *Operations Research Letters*, 2003.
- [BV04] S. Boyd and L. Vandenberghe. *Convex optimization*. Cambridge university press, 2004.
- [CCS<sup>+</sup>16] Pratik Chaudhari, Anna Choromanska, Stefano Soatto, Yann LeCun, Carlo Baldassi, Christian Borgs, Jennifer Chayes, Levent Sagun, and Riccardo Zecchina. Entropy-sgd: Biasing gradient descent into wide valleys. *arXiv preprint arXiv:1611.01838*, 2016.
- [CGV13] Mahdi Cheraghchi, Venkatesan Guruswami, and Ameya Velingker. Restricted isometry of Fourier matrices and list decodability of random linear codes. *SIAM Journal on Computing*, 42(5):1888–1914, 2013.
- [Che52] Herman Chernoff. A measure of asymptotic efficiency for tests of a hypothesis based on the sum of observations. *The Annals of Mathematical Statistics*, pages 493–507, 1952.
- [CKL<sup>+</sup>21] Jeremy Cohen, Simran Kaur, Yuanzhi Li, J Zico Kolter, and Ameet Talwalkar. Gradient descent on neural networks typically occurs at the edge of stability. *ICLR*, 2021.
- [CLGo1] Rich Caruana, Steve Lawrence, and C Lee Giles. Overfitting in neural nets: Backpropagation, conjugate gradient, and early stopping. In *Advances in Neural Information Processing Systems (NIPS)*, pages 402–408, 2001.
- [CT06] Emmanuel J Candes and Terence Tao. Near-optimal signal recovery from random projections: Universal encoding strategies? *IEEE transactions on information theory*, 52(12):5406–5425, 2006.
- [DHS<sup>+</sup>19] Simon S Du, Kangcheng Hou, Russ R Salakhutdinov, Barnabas Poczos, Ruosong Wang, and Keyulu Xu. Graph neural tangent kernel: Fusing graph neural networks with graph kernels. In *Advances in Neural Information Processing Systems*, pages 5724–5734, 2019.
- [DKB15] Laurent Dinh, David Krueger, and Yoshua Bengio. NICE: Nonlinear Independent Component Analysis. *Proc. ICLR*, 2015.

- [DLL<sup>+</sup>18] Simon S Du, Jason D Lee, Haochuan Li, Liwei Wang, and Xiyu Zhai. Gradient descent finds global minima of deep neural networks. *arXiv preprint arXiv:1811.03804*, 2018.
- [DPBB17] Laurent Dinh, Razvan Pascanu, Samy Bengio, and Yoshua Bengio. Sharp minima can generalize for deep nets. In *International Conference on Machine Learning*, 2017.
- [DSDB17] Laurent Dinh, Jascha Sohl-Dickstein, and Samy Bengio. Density Estimation using Real NVP. *Proc. ICLR*, 2017.
- [DZPS18] Simon S Du, Xiyu Zhai, Barnabas Poczos, and Aarti Singh. Gradient descent provably optimizes over-parameterized neural networks. *arXiv preprint arXiv:1810.02054*, 2018.
- [DZPS19] Simon S Du, Xiyu Zhai, Barnabas Poczos, and Aarti Singh. Gradient descent provably optimizes over-parameterized neural networks. *arXiv preprint arXiv:1810.02054*, 2019.
- [EHJTo4] B. Efron, T. Hastie, I. Johnstone, and R. Tibshirani. Least angle regression. *The Annals of statistics*, 2004.
- [Fri01] Jerome H Friedman. Greedy function approximation: a gradient boosting machine. *Annals of statistics*, 2001.
- [FSS11] Rina Foygel, Ohad Shamir, Nati Srebro, and Ruslan R Salakhutdinov. Learning with the weighted trace-norm under arbitrary sampling distributions. In *Advances in Neural Information Processing Systems*, pages 2133–2141, 2011.
- [GHJY15a] Rong Ge, Furong Huang, Chi Jin, and Yang Yuan. Escaping from saddle points—online stochastic gradient for tensor decomposition. In *Conference on Learning Theory*, pages 797–842, 2015.
- [GHJY15b] Rong Ge, Furong Huang, Chi Jin, and Yang Yuan. Escaping from saddle points—online stochastic gradient for tensor decomposition. In *Conf. Learning Theory (COLT)*, 2015.
- [GJZ17a] Rong Ge, Chi Jin, and Yi Zheng. No spurious local minima in nonconvex low rank problems: A unified geometric analysis. In *Proceedings of the 34th International Conference on Machine Learning—Volume 70*, pages 1233–1242. JMLR. org, 2017.

- [GJZ17b] Rong Ge, Chi Jin, and Yi Zheng. No spurious local minima in nonconvex low rank problems: A unified geometric analysis. *arXiv preprint arXiv:1704.00708*, 2017.
- [GLM16a] Rong Ge, Jason D Lee, and Tengyu Ma. Matrix completion has no spurious local minimum. In *Advances in Neural Information Processing Systems*, pages 2973–2981, 2016.
- [GLM16b] Rong Ge, Jason D Lee, and Tengyu Ma. Matrix completion has no spurious local minimum. In *Advances in Neural Information Processing Systems (NIPS)*, 2016.
- [GLM18] Rong Ge, Jason D Lee, and Tengyu Ma. Learning one-hidden-layer neural networks with landscape design. In *ICLR*. *arXiv preprint arXiv:1711.00501*, 2018.
- [GLSS18a] Suriya Gunasekar, Jason Lee, Daniel Soudry, and Nathan Srebro. Characterizing implicit bias in terms of optimization geometry. *arXiv preprint arXiv:1802.08246*, 2018.
- [GLSS18b] Suriya Gunasekar, Jason D Lee, Daniel Soudry, and Nati Srebro. Implicit bias of gradient descent on linear convolutional networks. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc., 2018.
- [GMS05] Anna C Gilbert, S Muthukrishnan, and Martin Strauss. Improved time bounds for near-optimal sparse Fourier representations. In *Optics & Photonics 2005*, pages 59141A–59141A. International Society for Optics and Photonics, 2005.
- [GWB<sup>+</sup>17] Suriya Gunasekar, Blake E Woodworth, Srinadh Bhojanapalli, Behnam Neyshabur, and Nati Srebro. Implicit regularization in matrix factorization. In *Advances in Neural Information Processing Systems*, pages 6151–6159, 2017.
- [HHS17] Elad Hoffer, Itay Hubara, and Daniel Soudry. Train longer, generalize better: closing the generalization gap in large batch training of neural networks. In *Advances in Neural Information Processing Systems*, 2017.

- [HIKP12] Haitham Hassanieh, Piotr Indyk, Dina Katabi, and Eric Price. Nearly optimal sparse Fourier transform. In *Proceedings of the forty-fourth annual ACM symposium on Theory of computing*, pages 563–578. ACM, 2012.
- [HMR18] Moritz Hardt, Tengyu Ma, and Benjamin Recht. Gradient descent learns linear dynamical systems. In *JLMR*. arXiv preprint arXiv:1609.05191, 2018.
- [Hoe63] Wassily Hoeffding. Probability inequalities for sums of bounded random variables. *Journal of the American Statistical Association*, 58(301):13–30, 1963.
- [HR16] Ishay Haviv and Oded Regev. The restricted isometry property of subsampled Fourier matrices. In *SODA*, pages 288–297. arXiv preprint arXiv:1507.01768, 2016.
- [HS97] Sepp Hochreiter and Jürgen Schmidhuber. Flat minima. *Neural Computation*, 1997.
- [IK14] Piotr Indyk and Michael Kapralov. Sample-optimal Fourier sampling in any constant dimension. In *Foundations of Computer Science (FOCS), 2014 IEEE 55th Annual Symposium on*, pages 514–523. IEEE, 2014.
- [IS15] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International Conference on Machine Learning (ICML)*, pages 448–456, 2015.
- [JGH18] Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and generalization in neural networks. In *Advances in neural information processing systems*, pages 8571–8580, 2018.
- [JGN<sup>+</sup>17] Chi Jin, Rong Ge, Praneeth Netrapalli, Sham M Kakade, and Michael I Jordan. How to escape saddle points efficiently. *arXiv preprint arXiv:1703.00887*, 2017.
- [JNG<sup>+</sup>19] Chi Jin, Praneeth Netrapalli, Rong Ge, Sham M Kakade, and Michael I Jordan. On nonconvex optimization for machine learning: Gradients, stochasticity, and saddle points. *arXiv preprint arXiv:1902.04811*, 2019.
- [Kap16] Michael Kapralov. Sparse Fourier transform in any constant dimension with nearly-optimal sample complexity in sublinear time. In *Symposium on Theory of Computing Conference, STOC’16, Cambridge, MA, USA, June 19–21, 2016*, 2016.



- [Kaw16] Kenji Kawaguchi. Deep learning without poor local minima. In *Adv in Neural Information Proc. Systems (NIPS)*, 2016.
- [KD19] Diederik P. Kingma and Prafulla Dhariwal. GLOW: Generative Flow with Invertible  $1 \times 1$  convolutions. *Proc. Neurips*, 2019.
- [KGC17] Jan Kukavcka, Vladimir Golkov, and Daniel Cremers. Regularization for deep learning: A taxonomy. *arXiv preprint arXiv:1710.10686*, 2017.
- [KSK11] Sham M Kakade, Varun Kanade, Ohad Shamir, and Adam Kalai. Efficient learning of generalized linear and single index models with isotonic regression. In *Advances in Neural Information Processing Systems*, pages 927–935, 2011.
- [KMN<sup>+</sup>16] Nitish Shirish Keskar, Dheevatsa Mudigere, Jorge Nocedal, Mikhail Smelyanskiy, and Ping Tak Peter Tang. On large-batch training for deep learning: Generalization gap and sharp minima. In *International Conference on Learning Representations*, 2016.
- [KS09] Adam Tauman Kalai and Ravi Sastry. The isotron algorithm: High-dimensional isotonic regression. In *COLT*. Citeseer, 2009.
- [KST09] Sham M Kakade, Karthik Sridharan, and Ambuj Tewari. On the complexity of linear prediction: Risk bounds, margin bounds, and regularization. In *Advances in neural information processing systems*, 2009.
- [KVZ19] Michael Kapralov, Ameya Velingker, and Amir Zandieh. Dimension-independent sparse Fourier transform. In *Proceedings of the Thirtieth Annual ACM-SIAM Symposium on Discrete Algorithms*, pages 2709–2728. SIAM, 2019.
- [Lano2] John Langford. *Quantitatively tight sample complexity bounds*. PhD Thesis CMU, 2002.
- [LMAPH19] Stéphane Lathuilière, Pablo Mesejo, Xavier Alameda-Pineda, and Radu Horaud. A comprehensive analysis of deep regression. *IEEE transactions on pattern analysis and machine intelligence*, 2019.
- [LMZ18] Yuanzhi Li, Tengyu Ma, and Hongyang Zhang. Algorithmic regularization in over-parameterized matrix

- sensing and neural networks with quadratic activations. In *Conference On Learning Theory*, pages 2–47, 2018.
- [LSJR16] Jason D Lee, Max Simchowitz, Michael I Jordan, and Benjamin Recht. Gradient descent converges to minimizers. *arXiv preprint arXiv:1602.04915*, 2016.
- [MA19] Poorya Mianjy and Raman Arora. On dropout and nuclear norm regularization. In *International Conference on Machine Learning*, 2019.
- [MAV18] Poorya Mianjy, Raman Arora, and Rene Vidal. On the implicit bias of dropout. In *International Conference on Machine Learning*, pages 3537–3545, 2018.
- [McA99] David A McAllester. Some pac-bayesian theorems. *Machine Learning*, 37(3):355–363, 1999.
- [MRT18] Mehryar Mohri, Afshin Rostamizadeh, and Ameet Talwalkar. *Foundations of machine learning*. MIT press, 2018.
- [NBMS18] Behnam Neyshabur, Srinadh Bhojanapalli, David McAllester, and Nathan Srebro. A pac-bayesian approach to spectrally-normalized margin bounds for neural networks. *ICLR*, 2018.
- [Nes98] Yurii Nesterov. *Introductory Lectures on Convex Programming Volume I: Basic Course*. Springer, 1998.
- [Nes00] Yurii Nesterov. Squared functional systems and optimization problems. In *High performance optimization*, pages 405–440. Springer, 2000.
- [Ney17] Behnam Neyshabur. Implicit regularization in deep learning. *arXiv preprint arXiv:1709.01953*, 2017.
- [NH92] Steven J Nowlan and Geoffrey E Hinton. Simplifying neural networks by soft weight-sharing. *Neural computation*, 4(4):473–493, 1992.
- [NK19] V Nagarajan and Zico Kolter. Uniform convergence may be unable to explain generalization in deep learning. *NeurIPS*, 2019.
- [NP06] Yurii Nesterov and Boris T Polyak. Cubic regularization of Newton method and its global performance. *Mathematical Programming*, 108(1):177–205, 2006.

- [NSS15] Behnam Neyshabur, Ruslan R Salakhutdinov, and Nati Srebro. Path-sgd: Path-normalized optimization in deep neural networks. In *Advances in Neural Information Processing Systems*, pages 2422–2430, 2015.
- [NSW19] Vasileios Nakos, Zhao Song, and Zhengyu Wang. (Nearly) Sample-optimal sparse Fourier transform in any dimension; RIPless and Filterless. In *FOCS*. <https://arxiv.org/pdf/1909.11123.pdf>, 2019.
- [NTS15a] Behnam Neyshabur, Ryota Tomioka, and Nathan Srebro. In search of the real inductive bias: On the role of implicit regularization in deep learning. In *International Conference on Learning Representations*, 2015.
- [NTS15b] Behnam Neyshabur, Ryota Tomioka, and Nathan Srebro. Norm-based capacity control in neural networks. In *Conference on Learning Theory*, pages 1376–1401, 2015.
- [NY83] A. Nemirovskii and D. Yudin. *Problem complexity and method efficiency in optimization*. Wiley, 1983.
- [PKCS17] Dohyung Park, Anastasios Kyrillidis, Constantine Caramanis, and Sujay Sanghavi. Non-square matrix sensing without spurious local minima via the burer-monteiro approach. In *AISTATS*. arXiv preprint arXiv:1609.03240, 2017.
- [RDS04] Cynthia Rudin, Ingrid Daubechies, and Robert E Schapire. The dynamics of adaboost: Cyclic behavior and convergence of margins. *Journal of Machine Learning Research*, 5(Dec):1557–1595, 2004.
- [RV08] Mark Rudelson and Roman Vershynin. On sparse reconstruction from Fourier and Gaussian measurements. *Communications on Pure and Applied Mathematics: A Journal Issued by the Courant Institute of Mathematical Sciences*, 61(8):1025–1045, 2008.
- [SF12] Robert E Schapire and Yoav Freund. *Boosting: Foundations and algorithms*. MIT press, 2012.
- [SHK<sup>+</sup>14] Nitish Srivastava, Geoffrey E Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: a simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research (JMLR)*, 15(1), 2014.

- [SHS17] Daniel Soudry, Elad Hoffer, and Nathan Srebro. The implicit bias of gradient descent on separable data. *arXiv preprint arXiv:1710.10345*, 2017.
- [Smi18] Le Smith, Kindermans. Don't Decay the Learning Rate, Increase the Batch Size. In *ICLR*, 2018.
- [SQW16a] Ju Sun, Qing Qu, and John Wright. Complete dictionary recovery over the sphere i: Overview and the geometric picture. *IEEE Transactions on Information Theory*, 63(2):853–884, 2016.
- [SQW16b] Ju Sun, Qing Qu, and John Wright. A geometric analysis of phase retrieval. In *IEEE International Symposium on Information Theory (ISIT)*, pages 2379–2383, 2016.
- [SQW18] Ju Sun, Qing Qu, and John Wright. A geometric analysis of phase retrieval. *Foundations of Computational Mathematics*, 18(5):1131–1198, 2018.
- [SS10] Nathan Srebro and Ruslan R Salakhutdinov. Collaborative filtering in a non-uniform world: Learning with the weighted trace norm. In *Advances in Neural Information Processing Systems*, pages 2056–2064, 2010.
- [SSS10] Shai Shalev-Shwartz and Yoram Singer. On the equivalence of weak learnability and linear separability: New relaxations and efficient boosting algorithms. *Machine learning*, 2010.
- [SY19] Zhao Song and Xin Yang. Quadratic suffices for overparametrization via matrix chernoff bound. *arXiv preprint arXiv:1906.03593*, 2019.
- [Tel13] Matus Telgarsky. Margins, shrinkage, and boosting. *arXiv preprint arXiv:1303.4172*, 2013.
- [Tro15] Joel A Tropp. An introduction to matrix concentration inequalities. *Foundations and Trends® in Machine Learning*, 8(1-2):1–230, 2015.
- [WRS<sup>+</sup>17] Ashia C Wilson, Rebecca Roelofs, Mitchell Stern, Nati Srebro, and Benjamin Recht. The marginal value of adaptive gradient methods in machine learning. In *Advances in Neural Information Processing Systems*, 2017.
- [ZBH<sup>+</sup>16a] Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep learning requires rethinking generalization. *arXiv preprint arXiv:1611.03530*, 2016.

- [ZBH<sup>+</sup>16b] Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep learning requires rethinking generalization. *arXiv preprint arXiv:1611.03530*, 2016.
- [ZY<sup>+</sup>05] Tong Zhang, Bin Yu, et al. Boosting with early stopping: Convergence and consistency. *The Annals of Statistics*, 2005.