

Probability $P(E_1 \cup E_2) = P(E_1) + P(E_2) - P(E_1 \cap E_2)$ $P(X=a) = \sum_b P(X=a, Y=b)$
 Expectation/mean $E[X] = \sum_a a P(X=a)$ Variance $Var[X] = E[(X - E[X])^2]$
 Independence $P(X, Y) = P(X)P(Y)$ Conditional Prob $P(X=a|Y=b) = \frac{P(X=a, Y=b)}{P(Y=b)}$
 Conditional Ind $P(X, Y|Z) = P(X|Z)P(Y|Z)$ like likelihood Chain Rule $P(A_1, A_2, \dots, A_n) = P(A_1)P(A_2|A_1)P(A_3|A_2, A_1) \dots$
 Bayes' Rule $P(F|S) = \frac{P(F, S)}{P(S)} = \frac{P(S|F)P(F)}{P(S)}$ prior $\frac{1}{n-1} X^T X$

PCA $\mu_x = \frac{1}{n} \sum_{i=1}^n \vec{x}_i$, $\vec{x}_i = \vec{x}_i - \mu_x$ so that new $\mu_x = 0 \rightarrow$ covariance matrix $S = \frac{1}{n-1} \sum_{i=1}^n \vec{x}_i \vec{x}_i^T = U \Lambda U^T$
 high dim \rightarrow low $X = X - \bar{X}$ \sum eigenvalues = \sum diagonal elements \leftarrow new representation of \vec{x}_i is $(u_1^T \vec{x}_i, \dots, u_d^T \vec{x}_i)$
 $\vec{x}_i = \sum_{j=1}^m a_{ij} u_j$, $a_{ij} = u_j^T \vec{x}_i$ eigenvectors u_1, \dots, u_d eigenvalues $\lambda_1, \dots, \lambda_d$

Logic Precedence $\neg, \wedge, \vee, \Rightarrow, \Leftrightarrow$ Entailment $A \models B$ $\boxed{B|A}$ if in every interpretation where A true, B also true
 Logical equivalences $(a \wedge b) \equiv (b \wedge a)$ commutativity $(a \vee b) \equiv (b \vee a)$ commutativity $A \oplus B = (A \wedge \neg B) \vee (\neg A \wedge B)$
 $((a \wedge b) \wedge c) \equiv (a \wedge (b \wedge c))$ $((a \vee b) \vee c) \equiv (a \vee (b \vee c))$ associativity $(a \Rightarrow b) \equiv (\neg a \vee b)$ implication
 $\neg(\neg a) \equiv a$ double negation $(a \Rightarrow b) \equiv (\neg b \Rightarrow \neg a)$ contraposition $(a \Leftrightarrow b) \equiv ((a \Rightarrow b) \wedge (b \Rightarrow a))$ biconditional $\neg(a \wedge b) \equiv (\neg a \vee \neg b)$ de Morgan
 $(a \Leftrightarrow b) \equiv ((a \Rightarrow b) \wedge (b \Rightarrow a))$ $(a \vee (b \wedge c)) \equiv ((a \vee b) \wedge (a \vee c))$ distributivity

CNF ① replace \Leftrightarrow ② replace \Rightarrow ③ move negation inward ④ apply distributivity of \vee over \wedge
 NLP $k=0$: Unigram $P(w_1, w_2, \dots, w_n) = P(w_1)P(w_2) \dots P(w_n)$ full independence assumption
 $k=1$: Bigram $P(w_1, w_2, \dots, w_n) = P(w_1)P(w_2|w_1)P(w_3|w_2) \dots P(w_n|w_{n-1})$ Markov assumption
 $k=n-1$: n-gram $P(w_i|w_{i-1}) = \frac{\text{count}(w_{i-1}, w_i)}{\text{count}(w_{i-1})}$ smoothing $\frac{\text{count}(w_{i-1}, w_i) + 1}{\text{count}(w_{i-1}) + V}$ (training count) $P(w_1, \dots, w_n) = P(w_1)P(w_2|w_1)P(w_3|w_2) \dots$

perplexity $PP(W) = P(w_1, w_2, \dots, w_n)^{-\frac{1}{n}}$ lower is better $\frac{1}{n} \sum_{c=1}^V \exp(u_c^T v_c)$
 Word2vec likelihood $L(\theta) = \prod_{t=1}^T \prod_{-a=j=a} P(w_{t+j} | w_t, \theta)$ maximize this $\sum \exp(u_w^T v_c)$
 word vector \uparrow all positions \uparrow window length of $2a$ context word \uparrow center word \uparrow two vectors v_w, u_w for center/context per word

ML Supervised: Classification, regression training set (multiset) = $\{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$
 label discrete label continuous i.i.d. assumption: $(x_i, y_i) \sim p_{XY}$ goal: learn mapping $f: X \rightarrow Y$
 loss function: 0-1 loss for classification $\ell(f, \vec{x}, y) = \mathbb{1}_{f(\vec{x}) \neq y}$ $f(x)$ predicts label y of feature x
 squared loss for regression $\ell(f, \vec{x}, y) = (f(\vec{x}) - y)^2$ empirical risk/training set error $\hat{R}(f) = \frac{1}{n} \sum_{i=1}^n \ell(f, \vec{x}_i, y_i)$
 clustering given: data with no label x_1, \dots, x_n goal: discover patterns and structures in data $\hat{f} = \arg \min_{f \in \mathcal{F}} \hat{R}(f)$ from learning

Unsupervised: dim reduction, clustering output: division into clusters with intra-cluster similarity and inter-cluster dissimilarity
 Unsupervised K-means ① select k cluster centers c_1, \dots, c_k ② for each point x , find closest center in Euclidean distance $y(x) = \arg \min_{i=1..k} \|x - c_i\|$ change x assignments $y(x)$ to centers
 ③ update cluster centers as centroids $c_i = \frac{\sum_{x: y(x)=i} x}{\sum_{x: y(x)=i} 1}$ change centers (cd) ④ update by repeating ②, ③ until convergence
 minimizes distortion $\sum_{d=1..P} \sum_{i=1..k} [x(d) - c_{y(x)}(d)]^2$ hill-climbing
 decide k : $\arg \min_k \text{distortion} + \lambda D k \log N$ \uparrow #dim \uparrow #clusters \uparrow #points

Hierarchical (no need k) input: points output: a hierarchy (binary tree) maximum depth on n points: $n-1$
 agglomerative: bottom-up/divisive: top-down \rightarrow repeat: get a closest pair of clusters and merge
 \rightarrow single-linkage $\min_{x_1 \in A, x_2 \in B} d(x_1, x_2)$
 \rightarrow complete-linkage $\max_{x_1 \in A, x_2 \in B} d(x_1, x_2)$
 \rightarrow average-linkage $\frac{1}{|A||B|} \sum_{x_1 \in A, x_2 \in B} d(x_1, x_2)$

Unsupervised II Part

Density Estimation goal: given samples X_1, \dots, X_n from distribution P , estimate P
 Simplist ideas: histograms define bins; count # of samples in each bin, normalize

Kernel Density Estimation density as combination of kernels
 $f(x) = \frac{1}{nh} \sum_{i=1}^n K\left(\frac{x-X_i}{h}\right)$ ← center at each point
 kernel func often Gaussian width param

Linear Regression

goal: minimize square loss
 Train set: $(\vec{x}_1, y_1), (\vec{x}_2, y_2), \dots, (\vec{x}_n, y_n)$ notation: $f(x) = \theta_0 + x^T \theta$, let $x = [\frac{1}{x}]$, then $f(x) = x^T \theta$
 $X = \begin{bmatrix} x_1^T \\ \vdots \\ x_n^T \end{bmatrix}$ $Y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}$ empirical risk: $\frac{1}{n} \|X\theta - Y\|^2$ error/residual: $|y_i - f(x_i)| = |y_i - \theta^T x_i|$
 solution: $\hat{\theta} = (X^T X)^{-1} X^T Y$ MSE = $\frac{1}{n} \sum_{i=1}^n (\hat{\theta}^T x_i - y_i)^2$
 $n^*(d+1)$ matrix

Classification

logistic regression: $y \in \{0, 1\}$ $\theta^T x \in [0, 1]$ $P(y=1|x) = \frac{1}{1 + \exp(-\theta^T x)}$
 KNN input: train data $(\vec{x}_1, y_1), \dots, (\vec{x}_n, y_n)$; distance func $d(\vec{x}_i, \vec{x}_j)$; num of neighbors k , test data x^*
 ① find k training instances $\vec{x}_{i1}, \dots, \vec{x}_{ik}$ closet to x^* under $d(\vec{x}_i, \vec{x}_j)$
 ② output y^* as the majority class of y_{i1}, \dots, y_{ik} ; break ties randomly (for classification)
 ③ (for regression) output the predicted $y^* = \frac{1}{k} (y_{i1} + \dots + y_{ik})$

pick data: (k)

distance: categorical features - Hamming distance = count (same position, diff items in 2 strs)
 numerical features - p-norm ($p \geq 1$)
 $d(x, x') = \|x - x'\|_p = \left(\sum_{i=1}^d |x_i - x'_i|^p\right)^{\frac{1}{p}}$
 used other situations: 2-norm = Euclidean $d(x, x') = \|x - x'\|_2 = \sqrt{\sum_{i=1}^d (x_i - x'_i)^2}$
 1-norm = Manhattan $d(x, x') = \|x - x'\|_1 = \sum_{i=1}^d |x_i - x'_i|$
 ① training, tuning, test shuffled & splitted randomly
 ② classify tuning with different k
 ③ pick k with least tuning error
 ④ report test error
 Error/Accuracy: error = $\frac{1}{n} \sum_{i=1}^n \mathbb{1}[f(x_i) \neq y_i]$ accuracy = $1 - \text{error}$

MLE

labeled train data $(X_1, y_1), \dots, (X_n, y_n) \rightarrow$ learning algo \rightarrow classifier f
 i.i.d. (fixed underlying distribution)

MLE (Maximum likelihood estimation): best fits data MAP (Maximum a posteriori): best fits data & prior assumption
 $\theta_{MLE} = \arg \max_{\theta} P(X|\theta)$ likelihood func: $L(\theta) = \prod_{i=1}^n P(X_i|\theta)$ $\theta_{MAP} = \arg \max_{\theta} P(X|\theta) P(\theta)$
 $\theta^* = \arg \max_{\theta} \ell(\theta) = \frac{N_H}{N_H + N_T}$

Bernoulli likelihood func $L(\theta) = \theta^{N_H} (1-\theta)^{N_T}$ Log-likelihood $\ell(\theta) = N_H \log \theta + N_T \log (1-\theta)$
 Gaussian likelihood func $L(\mu, \sigma^2 | X) = \prod_{i=1}^n P(X_i; \mu, \sigma^2) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(X_i - \mu)^2}{2\sigma^2}\right)$
 $\mu^*, \sigma^{*2} = \arg \max_{\mu, \sigma^2} \prod_{i=1}^n P(X_i; \mu, \sigma^2)$ $\mu = \frac{1}{n} \sum_{i=1}^n X_i$ $\sigma^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \mu)^2 = E[(X - \mu)^2]$

Naive Bayes Assumption: $P(X_1, \dots, X_k | y) P(y) = \prod_{i=1}^k P(X_i | y) P(y)$

Perceptron

Linear: output $f = \langle w, \vec{x} \rangle + b$ step function activation: $\sigma(x) = \begin{cases} 1 & \text{if } x > 0 \\ 0 & \text{otherwise} \end{cases}$
 with activation: $o = \sigma(\langle w, \vec{x} \rangle + b)$ sigmoid/logistic activation: $\sigma(z) = \text{sigmoid}(z) = \frac{1}{1 + \exp(-z)}$
 multi-layer: $\sigma(z) = \tanh(z) = \frac{1 - \exp(-2z)}{1 + \exp(-2z)}$ tanh activation: $\sigma(z) = \frac{1}{1 + \exp(-z)}$
 Input $n=3$ neurons $h_2 = \sigma\left(\sum_{i=1}^d X_i W_{2i}^{(1)} + b_2\right)$ $h_3 = \sigma\left(\sum_{i=1}^d X_i W_{3i}^{(1)} + b_3\right)$ $W_3^{(2)}$ Output $z = \text{sigmoid}(z) - 1$
 $x \in \mathbb{R}^d$ $W_1^{(2)}$ $W_2^{(2)}$ $W_3^{(2)}$ $W_1^{(1)}$ $W_2^{(1)}$ $W_3^{(1)}$ b_1 b_2 b_3
 MLE $\max_w \sum_i \log \frac{1}{1 + \exp(-y_i w^T \vec{x}_i)}$
 MAP $\min_w \sum_i -\log \frac{1}{1 + \exp(-y_i w^T \vec{x}_i)} + \frac{\lambda}{2} \|w\|_2^2$
 softmax: $p(y|\vec{x}) = \text{softmax}(f) = \frac{\exp(f_y)}{\sum_i \exp(f_i)}$

Distribution

$E(ax+b) = aE[X] + b$ $E[g(X)] = \sum_k g(k) p_X(k) / \int_{-\infty}^{\infty} g(x) f(x) dx$ $\text{Var}(X) = E[(X - E[X])^2] = E[X^2] - (E[X])^2$
 $\text{Var}(aX+b) = a^2 \text{Var}(X)$ $X \sim \text{Geom}(p): p_X(k) = p(1-p)^{k-1}, E = \frac{1}{p}, \text{Var} = \frac{1-p}{p^2}$
 $X \sim \text{Ber}(p): p_X(0) = 1-p, p_X(1) = p, E = p, \text{Var} = p(1-p)$ $X \sim \text{Unif}[a, b]: f_X(t) = \frac{1}{b-a}, E = \frac{a+b}{2}, \text{Var} = \frac{(b-a)^2}{12}$
 $X \sim \text{Bin}(n, p): p_X(k) = \binom{n}{k} p^k (1-p)^{n-k}, E = np, \text{Var} = np(1-p)$
 $X \sim \text{Normal}(\mu, \sigma^2): f_X(t) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(t-\mu)^2}{2\sigma^2}}, E = \mu, \text{Var} = \sigma^2, z = \frac{x-\mu}{\sigma}$

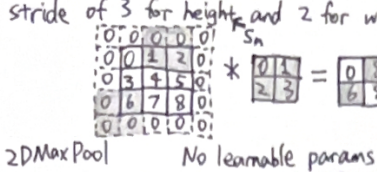
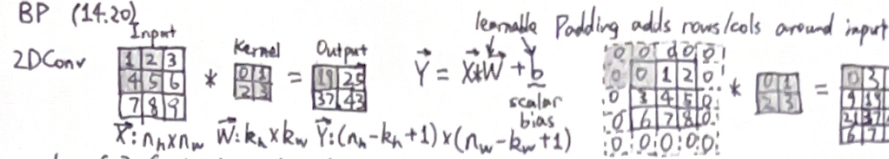
GD

① choose learning rate $\lambda > 0$ loss of NN: $\frac{1}{101} \sum_i \ell(X_i^T, y_i)$ ② init model params w_0 ③ for $t=1, 2, \dots$ update params $w_t = w_{t-1} - \lambda \sum_{i=0}^T \frac{\partial \ell(X_i^T, y_i)}{\partial w_{t-1}}$ until convergence

CS540 SP22 Final Note Ruixuan Tu

BP (14.20)

CNN and DL (15-18)



LeNet < AlexNet < VGG (Lec 16) < ResNet

Avoid Overfitting: Data Augmentation (one regularization) - transform & add new samples, can also in text by thesaurus

Classic regulation $\min_{\theta} \frac{1}{n} \sum_{i=1}^n \ell(f_{\theta}(x_i), y_i) + \lambda R(f_{\theta})$

standard loss regularization

Review: Lec 18

Game (19-20)

Goal-Max reward

Properties: N of players, action space, param

standard loss regularization

regularizer

(turns, payoff) (same time)

Rational Player at DSE if exists

Absolute Best -> Dominant Strategy Equilibrium

Nash Equilibrium: Best response to each other

(zero-sum: win-lose, General-sum: prison's dilemma)

Minimax: max of children max's turn, min at min, root and leaf nodes; heuristics; normal form

$u_i(a_i^*, a_{-i}^*) \geq u_i(b, a_{-i}^*)$

Mixed: $u_i(p, q) = u_i(f_1, q) = u_i(f_2, q)$

Search (21-23)

Problem: State space S, Initial state I ∈ S, Goal state G ∈ S, Successor function Succ(s) ∈ S, Cost(s, s')

BFS, DFS Uniform-cost search: use PQ in BFS Iterative deepening: DFS with depth limit and repeat

Informed A* $0 \leq h(s) \leq h^*(s)$ (PG) open ← S ⇒ open empty, fail ⇒ remove n which $f(n)$ min from open to close

IDA*: with ID

n' not in open/closed, estimate $h(n'), g(n') = g(n) + c(n, n'), f(n') = g(n') + h(n')$

n' in open/closed, $g(n')$ lower, update n' backward ptr to path of $g(n')$

not lower $g(n')$: nothing

Hill-Climbing: Local Optima, always pick max neighbor until descent

Simulated Annealing: $T=1$, for $k=0$ to K : $T \leftarrow T * 0.99$, pick random neighbor $t \leftarrow \text{neighbor}(s)$

(23.20)

if $f(s) \leq f(t)$, then $s \leftarrow t$, else with prob $P(f(s), f(t); T)$ do $s \leftarrow t$

Genetic Algorithm: ① keep a population (a fixed no states) ② selection, cross-over, mutation according to p

$P_i = f(s_i) / \sum_j f(s_j)$ reproduction prob

RL (24-25)

MDP $M = (S, A, P, r, \mu, \gamma)$ state set S, initial s_0 . Action set A. Reward function: $r(s_t, a_t)$

State transition model $P(s_{t+1} | s_t, a_t)$ (Markov assumption). Policy $\pi(s)$. Discount factor $\gamma \in (0, 1)$

$V^{\pi}(s_0) = \sum_{seq} P(seq) V(seq); V(seq) = \sum_{t=0}^{\infty} \gamma^t r_t$

Geometric Series:

$V^*(s) = \max_a r(s, a) + \gamma \sum_{s'} P(s' | s, a) V^*(s')$

$\forall |r| < 1, \sum_{n=1}^{\infty} ar^{n-1} = \frac{a}{1-r}, \frac{1}{1-x} = \sum_{n=0}^{\infty} x^n$

Q-learning (25.21)