Stochastic Gradient
OOOOOOOOO

Regularization
OOOOOOOOOOOOOOOOOO

Multi–Class Classification
OOOOOOOOOO

# CS540 Introduction to Artificial Intelligence
# Lecture 4

Young Wu
Based on lecture slides by Jerry Zhu and Yingyu Liang

May 29, 2019

**Stochastic Gradient**
●○○○○○○○○

Regularization
○○○○○○○○○○○○○○○○○○

Multi-Class Classification
○○○○○○○○○○

# Test

## Quiz (Graded)

- A:
- B:
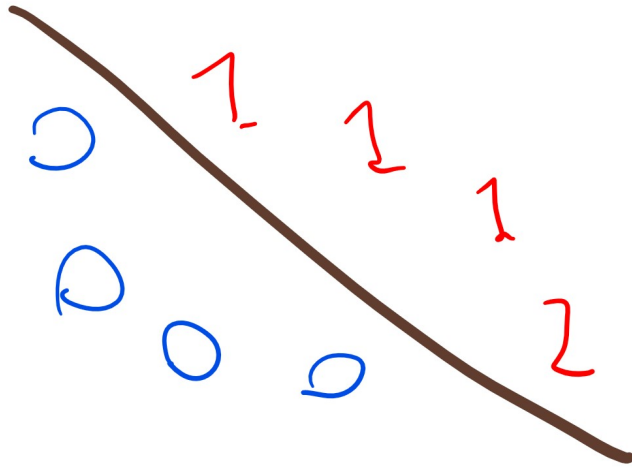- C:
- D: Choose this.
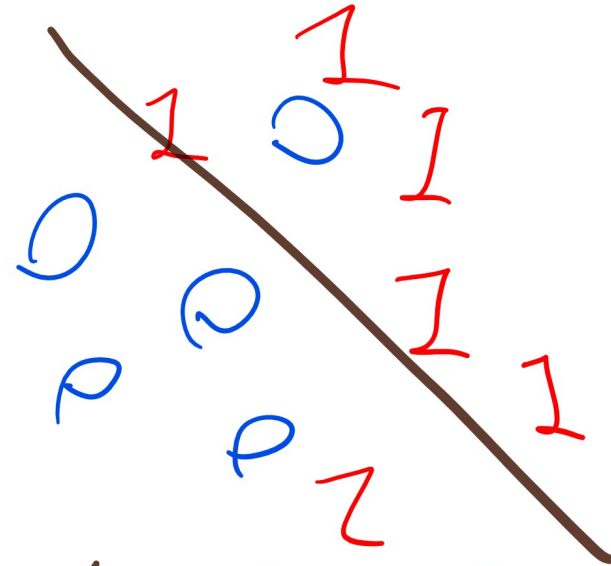- E:

# Homework

### Quiz (Participation)

- Have you finished homework 1
- A: Waiting for solution.
- B: Will start soon.
- C: Started.
- D: Does not work due to bugs.
- E: Finished: 90+ percent accuracy.

**Stochastic Gradient**
OO●OOOOOO

Regularization
OOOOOOOOOOOOOOOOO

Multi-Class Classification
OOOOOOOOOO

# Neural Network Diagram

## Review



Perceptron Algorithm

Logistic Regression

$\downarrow$

Gradient Descent

$$w = w - \alpha \nabla_w C$$

$\uparrow$ opposite direction

Neural Network

**Stochastic Gradient**
○○○●○○○○○

Regularization
○○○○○○○○○○○○○○○○○

Multi-Class Classification
○○○○○○○○○○

# Multi-Layer Neural Network Diagram

Review

$$\nabla_w C = \sum_{i=1}^{n} \frac{\partial C}{\partial a_i} \cdot \nabla_w a_i$$

slow.

**Stochastic Gradient**
○○○○●○○○○

Regularization
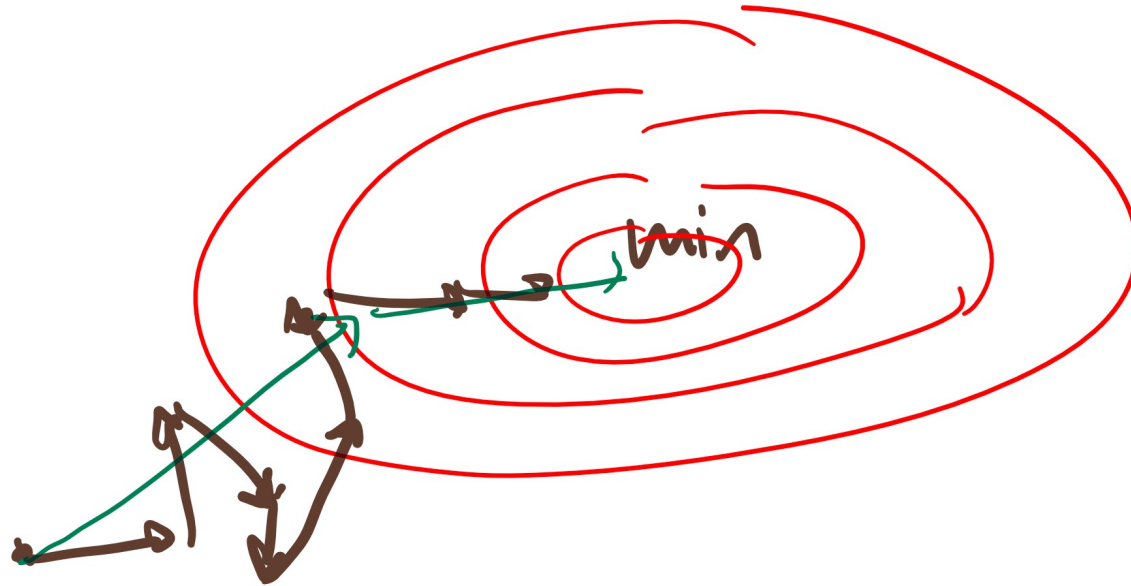○○○○○○○○○○○○○○○○○○

Multi-Class Classification
○○○○○○○○○○○

# Stochastic Gradient Descent
## Motivation

- Each gradient descent step requires the computation of gradients for all training instances $i = 1, 2, ..., n$. It is very costly.

- Stochastic gradient descent picks one instance $x_i$ randomly, compute the gradient, and update the weights and biases.

- When a batch of instances is selected randomly each time, it is called batch gradient descent.

**Stochastic Gradient**
oooooo●ooo

Regularization
ooooooooooooooooooo

Multi-Class Classification
oooooooooo

# Stochastic Gradient Descent Diagram

## Motivation

**Stochastic Gradient**
○○○○○○○●○○

Regularization
○○○○○○○○○○○○○○○○○○

Multi-Class Classification
○○○○○○○○○○

# Stochastic Gradient Descent, Part 1

## Algorithm

- Inputs, Outputs: same as backpropagation.

- Initialize the weights.

- Randomly permute (shuffle) the training set. Evaluate the activation functions at one instance at a time.

- Compute the gradient using the chain rule.

weight
at layer
↳

no Sum over $i$

$$\frac{\partial C}{\partial w_{j'j}^{(l)}} = \delta_{ij}^{(l)} a_{ij'}^{(l-1)}$$

$$\frac{\partial C}{\partial b_j^{(l)}} = \delta_{ij}^{(l)}$$

**Stochastic Gradient**
○○○○○○○●○

**Regularization**
○○○○○○○○○○○○○○○○○○

**Multi-Class Classification**
○○○○○○○○○○○

# Stochastic Gradient Descent, Part 2
## Algorithm

- Update the weights and biases using gradient descent.

For $l = 1, 2, ..., L$

$$w_{j'j}^{(l)} \leftarrow w_{j'j}^{(l)} - \alpha \frac{\partial C}{\partial w_{j'j}^{(l)}}, j' = 1, 2, ...., m^{(l-1)}, j = 1, 2, ...., m^{(l)}$$

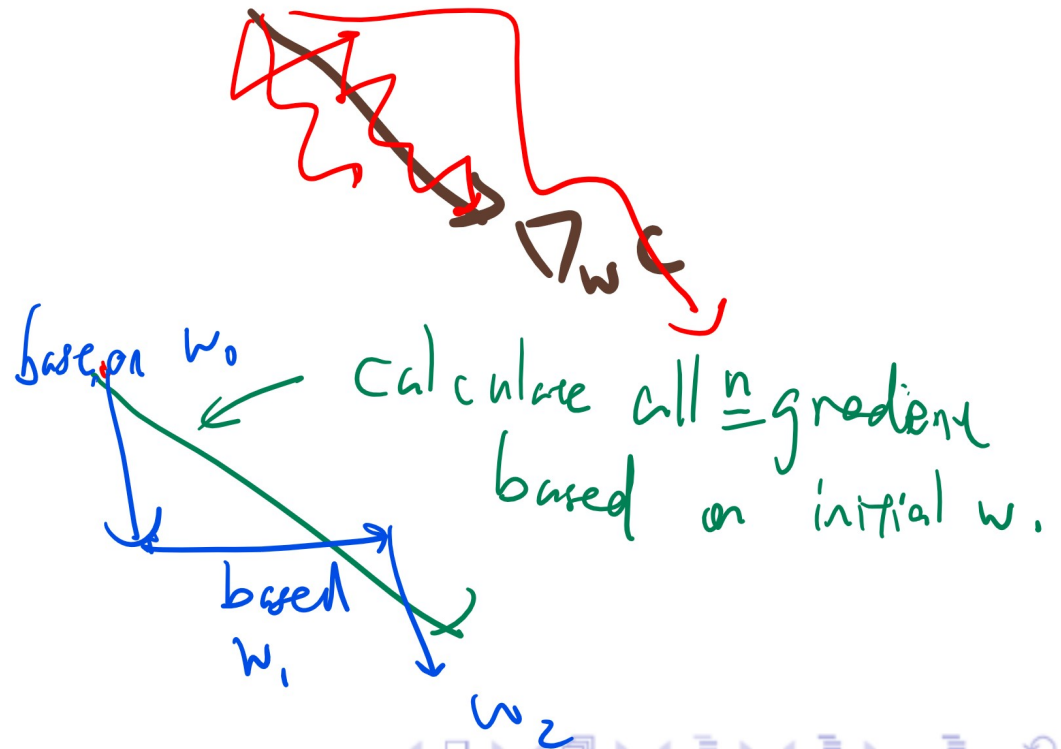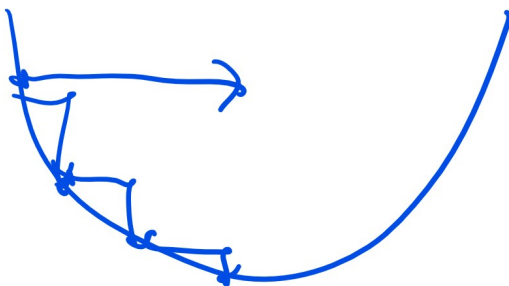$$b_j^{(l)} \leftarrow b_j^{(l)} - \alpha \frac{\partial C}{\partial b_j^{(l)}}, j = 1, 2, ...., m^{(l)}$$

- Repeat the process until convergent.

$$|C - C^{\text{prev}}| < \varepsilon$$

**Stochastic Gradient**
○○○○○○○○●

Regularization
○○○○○○○○○○○○○○○○○

Multi-Class Classification
○○○○○○○○○○

# Stochastic vs Full Gradient Descent
## Quiz (Participation)

- Given the same initial weights and biases, stochastic gradient descent with instances picked randomly without replacement and full gradient descent lead to the same updated weights.

- A: Do not choose this.

- B: True.

- C: Do not choose this.

- D: False.

- E: Do not choose this.

$\nabla_w$

based on $w_0$

Calculate all $\underline{n}$ gradient based on initial $w$.

based $w_1$

$w_2$

Stochastic Gradient
ooooooooo

Regularization
●oooooooooooooooooo

Multi–Class Classification
ooooooooooo

# Generalization Error
## Motivation

- With a large number of hidden units and small enough learning rate $\alpha$, a multi-layer neural network can fit every finite training set perfectly.

- It does not imply the performance on the test set will be good.

- This problem is called overfitting.

Stochastic Gradient
○○○○○○○○○

Regularization
○●○○○○○○○○○○○○○○○○

Multi-Class Classification
○○○○○○○○○○

# Generalization Error Diagram

## Motivation

Stochastic Gradient
○○○○○○○○○

Regularization
○○●○○○○○○○○○○○○○○○○○

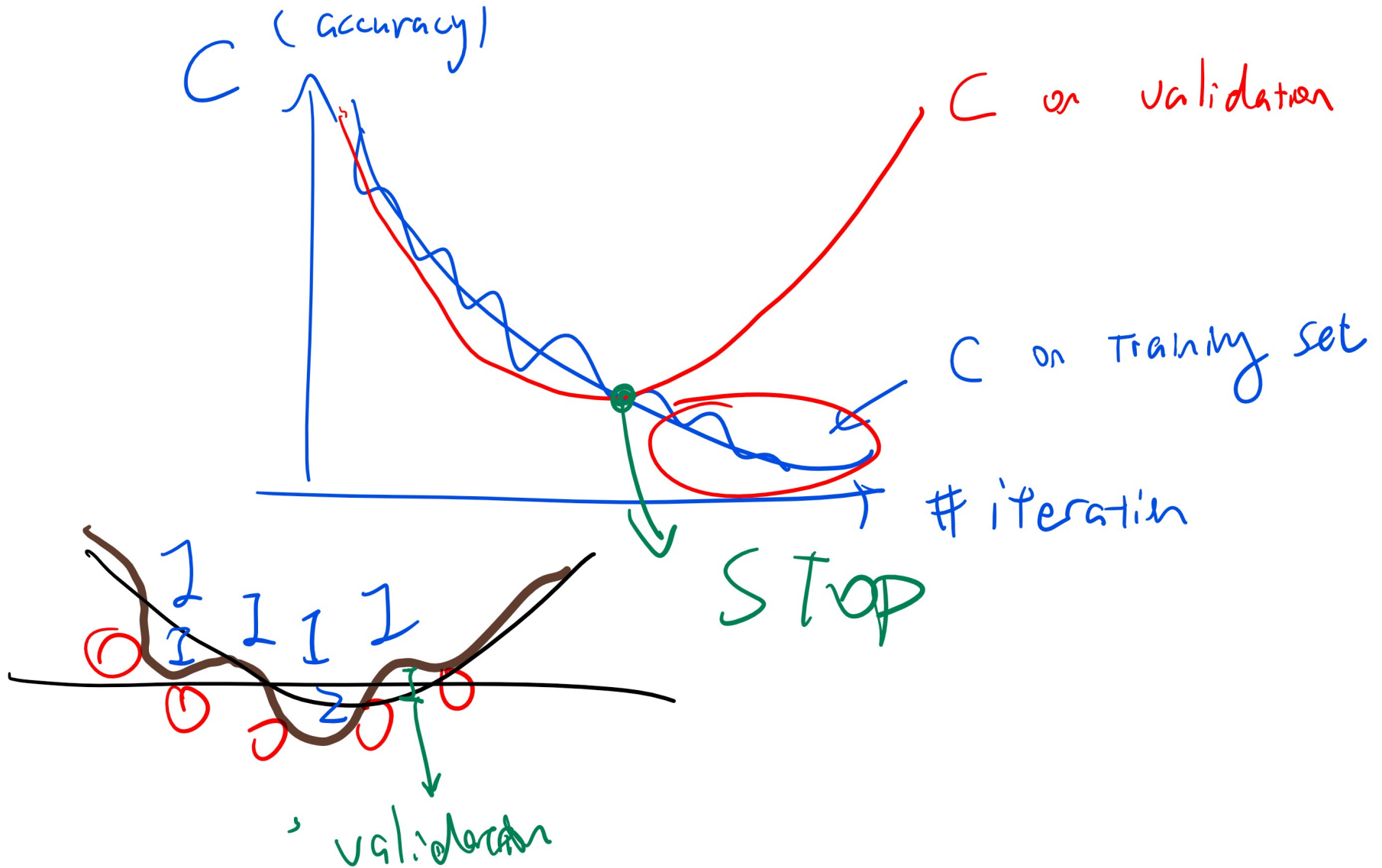Multi-Class Classification
○○○○○○○○○○

# Method 1, Validation Set
## Discussion

- Set aside a subset of the training set as the validation set.

- During training, the cost (or accuracy) on the training set will always be decreasing until it hits 0.

- Train the network until the cost (or accuracy) on the validation set begins to increase.

Stochastic Gradient
○○○○○○○○

**Regularization**
○○○●○○○○○○○○○○○○○

Multi-Class Classification
○○○○○○○○○○

# Validation Set Diagram

## Discussion

Stochastic Gradient
○○○○○○○○○

Regularization
○○○○●○○○○○○○○○○○○○○

Multi-Class Classification
○○○○○○○○○○

# Method 2, Drop Out

## Discussion

*only NN*

- At each hidden layer, a random set of units from that layer is set to 0.

- For example, each unit is retained with probability $p = 0.5$. During the test, the activations are reduced by $p = 0.5$ (or 50 percent).

- The intuition is that if a hidden unit works well with different combinations of other units, it does not rely on other units and it is likely to be individually useful.

*→ set to 0*

*1) set to 0*

Stochastic Gradient

oooooooo

Regularization

ooooo●ooooooooooooo

Multi-Class Classification

ooooooooooo

# Drop Out Diagram

## Discussion

Stochastic Gradient
OOOOOOOOO

Regularization
OOOOOO●OOOOOOOOOO

Multi-Class Classification
OOOOOOOOOO
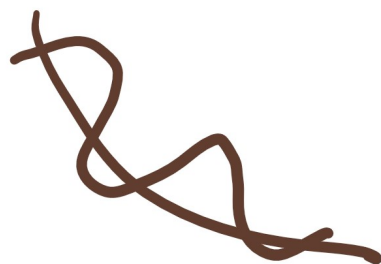
# Method 3, L1 and L2 Regularization
## Discussion

- The idea is to include an additional cost for non-zero weights.

- The models are simpler if many weights are zero.

- For example, if logistic regression has only a few non-zero weights, it means only a few features are relevant, so only these features are used for prediction.

Stochastic Gradient
○○○○○○○○○

Regularization
○○○○○○○●○○○○○○○○

Multi-Class Classification
○○○○○○○○○○

# Method 3, L1 Regularization

## Discussion

- For L1 regularization, add the 1-norm of the weights to the cost.

$$C = \sum_{i=1}^{n} (a_i - y_i)^2 + \lambda \left\| \begin{bmatrix} w \\ b \end{bmatrix} \right\|_1$$

$$= \sum_{i=1}^{n} (a_i - y_i)^2 + \lambda \left( \sum_{i=1}^{m} |w_i| + |b| \right)$$

Cost of having non-zero weights.

min C ⟹ try to make w ≈ 0

- Linear regression with L1 regularization is called LASSO (least absolute shrinkage and selection operator).

Stochastic Gradient
○○○○○○○○○

Regularization
○○○○○○○○●○○○○○○○○

Multi-Class Classification
○○○○○○○○○○

# Method 3, L2 Regularization

## Discussion

- For L2 regularization, add the 2-norm of the weights to the cost.

$$C = \sum_{i=1}^{n} (a_i - y_i)^2 + \frac{\lambda}{2} \left\| \begin{bmatrix} w \\ b \end{bmatrix} \right\|_2^2$$

$$\underbrace{\text{learning rate}} = \sum_{i=1}^{n} (a_i - y_i)^2 + \lambda \left( \frac{1}{2} \sum_{i=1}^{m} w_i^2 + b^2 \right)$$

$$w = w - \alpha \nabla_w C - \lambda w$$

regularization parameter.

easy for gradient descent.

Stochastic Gradient
○○○○○○○○○

Regularization
○○○○○○○○○○●○○○○○○○

Multi-Class Classification
○○○○○○○○○○

# L1 and L2 Regularization Comparison
## Discussion

- L1 regularization leads to more weights that are exactly 0. It is useful for feature selection.

- L2 regularization leads to more weights that are close to 0. It is easier to do gradient descent because 1-norm is not differentiable.

Stochastic Gradient
○○○○○○○○

Regularization
○○○○○○○○○●○○○○○○

Multi-Class Classification
○○○○○○○○○○

# L1 and L2 Regularization Diagram

## Discussion



$w_2$

min C
best $w$, such that $\|w\|_1 = 1$

fix
$|w_1| + |w_2| = 1$

decrease C

min C
best $w$ s.t. $\|w\|_2 = 1$

C = 3
C = 3.5
C = 4
C = 5

$w_1$

fix
$w_1^2 + w_2^2 = 1$

L1 $\rightarrow$ many $w = 0$ $\leftarrow$ feature selection

L2 $\rightarrow$ no $w = 0$ many $\approx 0$

Stochastic Gradient
○○○○○○○○○

Regularization
○○○○○○○○○○○○●○○○○○○

Multi–Class Classification
○○○○○○○○○○

# Method 4, Data Augmentation
## Discussion

- More training data can be created from the existing ones, for example, by translating or rotating the handwritten digits.

Stochastic Gradient
OOOOOOOOO

Regularization
OOOOOOOOOOOOO●OOOOO

Multi–Class Classification
OOOOOOOOOO

# Hyperparameters

## Discussion

- It is not clear how to choose the learning rate $\alpha$, the stopping criterion $\varepsilon$, and the regularization parameters. $\lambda, \beta \ldots$

- For neural networks, it is also not clear how to choose the number of hidden layers and the number of hidden units in each layer.

- The parameters that are not parameters of the functions in the hypothesis space are called hyperparameters.

Stochastic Gradient
○○○○○○○○○

Regularization
○○○○○○○○○○○○○●○○○○

Multi-Class Classification
○○○○○○○○○○

# K Fold Cross Validation

## Discussion

train on training set to find $v, b$

test on validation to compare performance

$C$, accuracy.

- Partition the training set into $K$ groups.

- Pick one group as the validation set.

- Train the model on the remaining training set.

- Repeat the process for each of the $K$ groups.

- Compare accuracy (or cost) for models with different hyperparameters and select the best one.

Stochastic Gradient
○○○○○○○○○

Regularization
○○○○○○○○○○○○○○●○○○

Multi-Class Classification
○○○○○○○○○○

# 5 Fold Cross Validation Example

## Discussion

- Partition the training set $S$ into 5 subsets $S_1, S_2, S_3, S_4, S_5$

$$S_i \cap S_j = \emptyset \text{ and } \bigcup_{i=1}^{5} S_i = S$$

| Iteration | Training | Validation |
|-----------|----------|------------|
| 1 | $S_2 \cup S_3 \cup S_4 \cup S_5$ | $S_1$ |
| 2 | $S_1 \cup S_3 \cup S_4 \cup S_5$ | $S_2$ |
| 3 | $S_1 \cup S_2 \cup S_4 \cup S_5$ | $S_3$ |
| 4 | $S_1 \cup S_2 \cup S_3 \cup S_5$ | $S_4$ |
| 5 | $S_1 \cup S_2 \cup S_3 \cup S_4$ | $S_5$ |

get C on all training instances.

Stochastic Gradient
○○○○○○○○○

Regularization
○○○○○○○○○○○○○○○●○○

Multi-Class Classification
○○○○○○○○○○

# Leave One Out Cross Validation
## Discussion

- If $K = n$, each time exactly one training instance is left out as the validation set. This special case is called Leave One Out Cross Validation (LOOCV).

Stochastic Gradient
000000000

Regularization
00000000000000000●0

Multi-Class Classification
0000000000

# Cross Validation, Part II
## Quiz (Graded)

5

- March 2018 Midterm Q9   *will repeat*

  $0\ 1\ 1\ 1\ 0 \neq 1$

- Consider the majority classifier that predict $\hat{y} =$ mode of the training data labels. What is the 2-fold cross validation accuracy (percentage of correct classification) on the following training set.

  $2/10 = 20\%$

  $S_1$      $S_2$

| x | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
|---|---|---|---|---|---|---|---|---|---|----|
| y | 1 | 1 | 0 | 1 | 1 | 0 | 0 | 1 | 0 | 0  |

- A: 0 percent, B: 10 percent, C. 20 percent

  *correct.*

- D: 50 percent, E: 100 percent

  $\hat{y}_{S_1} = 1$

Stochastic Gradient
○○○○○○○○○

Regularization
○○○○○○○○○○○○○○○○●

Multi-Class Classification
○○○○○○○○○○

# Cross Validation, Part I

## Quiz (Graded)

- March 2018 Midterm Q9

- Consider the majority classifier that predict $\hat{y}$ = mode of the training data labels. What is the LOOCV accuracy (percentage of correct classification) on the following training set.

*K = 10-fold CV*

*pnt on midterm!*

| x | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
|---|---|---|---|---|---|---|---|---|---|----|
| y | 1 | 1 | 0 | 1 | 1 | 0 | 0 | 1 | 0 | 0  |

$\hat{y} = 0$

- A: 0 percent, B: 10 percent, C: 20 percent
- D: 50 percent, E: 100 percent

Stochastic Gradient
○○○○○○○○○

Regularization
○○○○○○○○○○○○○○○○○○○

Multi-Class Classification
●○○○○○○○○○○

# Multi-Class Classification
## Discussion

- When there are $K$ categories to classify, the labels can take $K$ different values, $y_i \in \{1, 2, ..., K\}$.

- Logistic regression and neural network cannot be directly applied to these problems.

Stochastic Gradient
○○○●○○○○○

Regularization
○○○○○○○○○○○○○○○○○

Multi-Class Classification
○●○○○○○○○○

# Method 1, One VS All
## Discussion

$0$ vs not $0$     $1$  vs  not $1$   . . . .

- Train a binary classification model with labels $y'_i = \mathbb{1}_{\{y_i=j\}}$ for each $j = 1, 2, ..., K$.

- Given a new test instance $x_i$, evaluate the activation function $a_i^{(j)}$ from model $j$.

$$\hat{y}_i = \arg\max_j a_i^{(j)}$$

- One problem is that the scale of $a_i^{(j)}$ may be different for different $j$.

Stochastic Gradient
OOOOOOOOO

Regularization
OOOOOOOOOOOOOOOOOO

Multi-Class Classification
OO●OOOOOOOO

# Method 2, One ~~VS~~ One

## Discussion

$\partial$ vs 1       $0$ vs 2       $0$ vs 3       $\sim$ -

- Train a binary classification model with for each of the $\dfrac{K(K-1)}{2}$ pairs of labels.

- Given a new test instance $x_i$, apply all $\dfrac{K(K-1)}{2}$ models and output the class that receives the largest number of votes.

$$\hat{y}_i = \arg\max_j \sum_{j' \neq j} \hat{y}_i^{(j \text{ vs } j')}$$

- One problem is that it is not clear what to do if multiple classes receive the same number of votes.

Stochastic Gradient
ooooooooo

Regularization
oooooooooooooooooo

Multi-Class Classification
oooo●oooooo

# One Hot Encoding

$$y = 1, 2, 3, 4.$$

← no order

- If $y$ is not binary, use one-hot encoding for $y$.
- For example, if $y$ has three categories, then

$$y_i \in \left\{ \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} \right\}$$

$y=1 \qquad y=2 \qquad y=3$

Stochastic Gradient
oooooooooo

Regularization
oooooooooooooooooooo

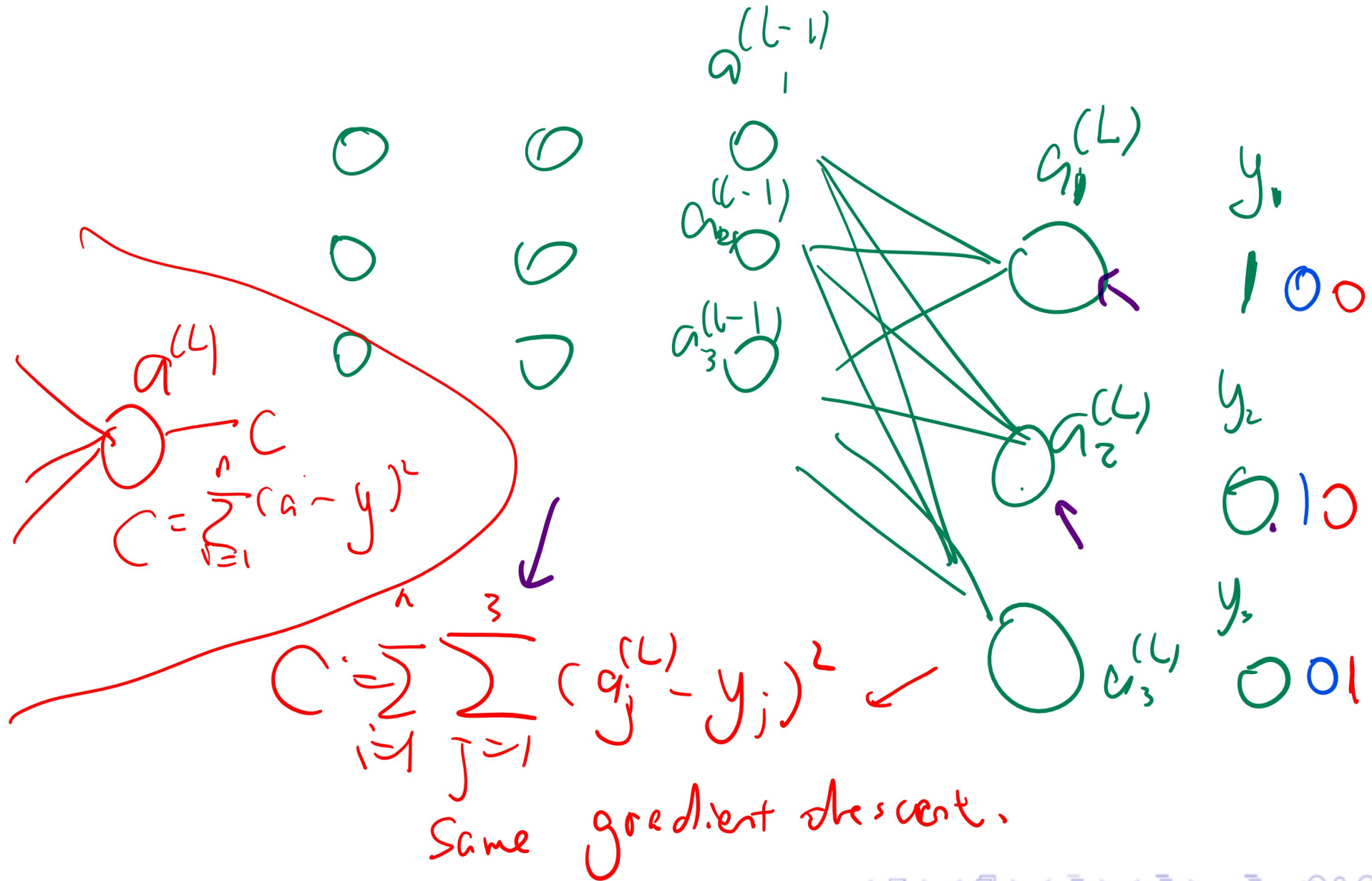Multi-Class Classification
ooooo●ooooo

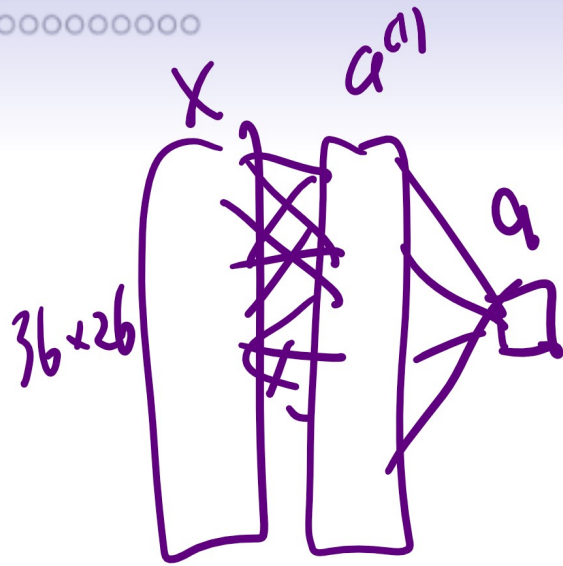# Method 3, Softmax Function

Discussion

- For both logistic regression and neural network, the last layer will have $K$ units, $a_{ij}$, for $j = 1, 2, ..., K$ and the softmax function is used instead of the sigmoid function.

$$a_{ij} = g\left(w_j^T x_i + b_j\right) = \frac{\exp\left(w_j^T x_i + b_j\right)}{\sum_{j'=1}^{K} \exp\left(w_{j'}^T x_i + b_{j'}\right)}, j = 1, 2, ..., K$$

Stochastic Gradient
○○○○○○○○○

Regularization
○○○○○○○○○○○○○○○○○

Multi-Class Classification
○○○○○●○○○○

# Softmax Function Diagram

## Discussion



$a_1^{(l-1)}$

$a_2^{(l-1)}$

$a_3^{(l-1)}$

$a_1^{(L)}$

$y_1$

$\boxed{1} \ O \ O$

$a^{(L)}$

$C$

$$C = \sum_{i=1}^{n}(a - y)^2$$

$a_2^{(L)}$

$y_2$

$O.10$

$a_3^{(L)}$

$y_3$

$O \ O1$

$$C = \sum_{i=1}^{n}\sum_{j=1}^{3}(a_j^{(L)} - y_j)^2$$

Same gradient descent.

Stochastic Gradient
○○○○○○○○○

Regularization
○○○○○○○○○○○○○○○○

Multi-Class Classification
○○○○○○○●○○○
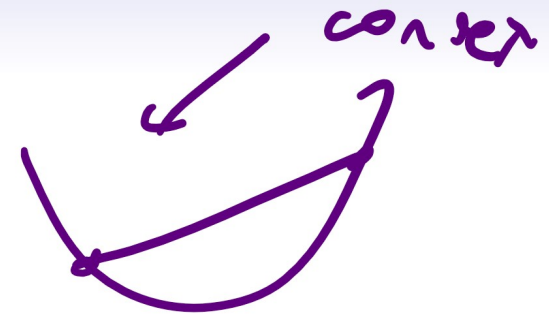
# Autoencoder

## Discussion

- A multi-layer neural network with the same input and output $y_i = x_i$ is called an autoencoder.

- The hidden layers have fewer units than the dimension of the input $m$.

- The hidden units form an encoding of the input with reduced dimensionality.

Stochastic Gradient
○○○○○○○○○

Regularization
○○○○○○○○○○○○○○○○○○

**Multi-Class Classification**
○○○○○○○○●○○

# Autoencode Diagram

## Discussion

Stochastic Gradient
ooooooooo

Regularization
oooooooooooooooooooo

Multi-Class Classification
oooooooooo●o

# Generative Adversarial Network
## Discussion

- Two competitive neural networks.

1. Generative network input random noise and output fake images.

2. Discriminative network input real and fake images and output label real or fake.

Stochastic Gradient
○○○○○○○○○

Regularization
○○○○○○○○○○○○○○○○○○

Multi-Class Classification
○○○○○○○○○○●

# Generative Adversarial Network Diagram

## Discussion