

# CS540 Introduction to Artificial Intelligence

## Lecture 2

Young Wu

Based on lecture slides by Jerry Zhu, Yingyu Liang, and Charles Dyer

May 18, 2020



# Zero-One Loss Function

## Motivation

$$0.5 = y_i \in [0, 1]$$

$$a_i = f(x_i) = 0.4, \quad -1.2$$

- An objective function is needed to select the "best"  $\hat{f}$ . An example is the zero-one loss.

$$\hat{f} = \arg \min_f \sum_{i=1}^n \mathbb{1}_{\{f(x_i) \neq y_i\}}$$

Handwritten annotations: "line" points to the  $f$  in the minimization; "classifier prediction" points to the  $f(x_i)$  term; "# of mistake" points to the indicator function.

- $\arg \min_f$  objective ( $f$ ) outputs the function that minimizes the objective.
- The objective function is called the cost function (or the loss function), and the objective is to minimize the cost.

# Squared Loss Function

## Motivation

- Zero-one loss counts the number of mistakes made by the classifier. The best classifier is the one that makes the fewest mistakes.
- Another example is the squared distance between the predicted and the actual  $y$  value:

$$\hat{f} = \arg \min_f \frac{1}{2} \sum_{i=1}^n (f(x_i) - y_i)^2$$

$a_i$  for perception  
 $\{f(x_i) - y_i\}$   
 sq loss  
 $\left[ \begin{array}{c} f(x_i) \\ 0 \\ 0 \\ \vdots \\ i \end{array} \right] \quad \left[ \begin{array}{c} y_i \\ 0 \\ -1 \\ \vdots \\ 0 \end{array} \right] \quad \leftrightarrow \quad \left[ \begin{array}{c} 0 \\ \vdots \\ 1 \\ 0 \end{array} \right]$



# Loss Functions Equivalence

## Quiz

- Which ones (multiple) of the following functions are equivalent to the squared error for binary classification?

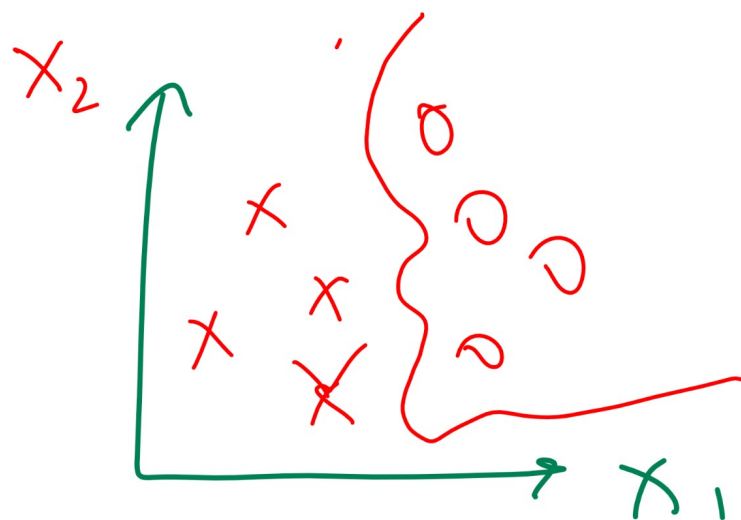
$$C = \sum_{i=1}^n \underbrace{(f(x_i) - y_i)^2}_{\text{sq}}$$

Q9

- |                                                                                                                                                                                                  | sq | abs |
|--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|----|-----|
| <ul style="list-style-type: none"> <li><input checked="" type="radio"/> A: <math>\sum \mathbb{1}_{\{f(x_i) \neq y_i\}}</math> ✓</li> </ul>                                                       | 0  | 0   |
| <ul style="list-style-type: none"> <li><input checked="" type="radio"/> B: <math>\sum \mathbb{1}_{\{f(x_i) = y_i\}}</math></li> </ul>                                                            | 0  | 1   |
| <ul style="list-style-type: none"> <li><input checked="" type="radio"/> C: <math>\sum  f(x_i) - y_i </math> ✓</li> </ul>                                                                         | 1  | 0   |
| <ul style="list-style-type: none"> <li><input checked="" type="radio"/> D: <math>\sum \max\{0, 1 - f(x_i) y_i\}</math></li> </ul>                                                                | 1  | 1   |
| <ul style="list-style-type: none"> <li><input checked="" type="radio"/> E: <math>\frac{1}{2} \sum \max\{0, 1 - (2 \cdot \underline{f(x_i)} - 1)(2 \cdot \underline{y_i} - 1)\}</math></li> </ul> | 0  | 0   |

# Function Space Diagram

## Motivation



# Hypothesis Space

## Motivation

- There are too many functions to choose from.
- There should be a smaller set of functions to choose  $\hat{f}$  from.

$$\hat{f} = \arg \min_{f \in \mathcal{H}} \frac{1}{2} \sum_{i=1}^n (f(x_i) - y_i)^2$$

- The set  $\mathcal{H}$  is called the hypothesis space.

# Activation Function

## Motivation

- Suppose  $\mathcal{H}$  is the set of functions that are compositions between another function  $g$  and linear functions.

$$(\hat{w}, \hat{b}) = \arg \min_{w, b} \frac{1}{2} \sum_{i=1}^n (a_i - y_i)^2$$

where  $a_i = g(w^T x + b)$  ← all functions that can be written like this

LTM perception

- $g$  is called the activation function.

$$a_i = \mathbb{1}(w^T x + b \geq 0)$$

# Linear Threshold Unit

## Motivation

- One simple choice is to use the step function as the activation function:


$$g(\boxed{\cdot}) = \mathbb{1}_{\{\boxed{\cdot} \geq 0\}} = \begin{cases} 1 & \text{if } \boxed{\cdot} \geq 0 \\ 0 & \text{if } \boxed{\cdot} < 0 \end{cases}$$

- This activation function is called linear threshold unit (LTU).

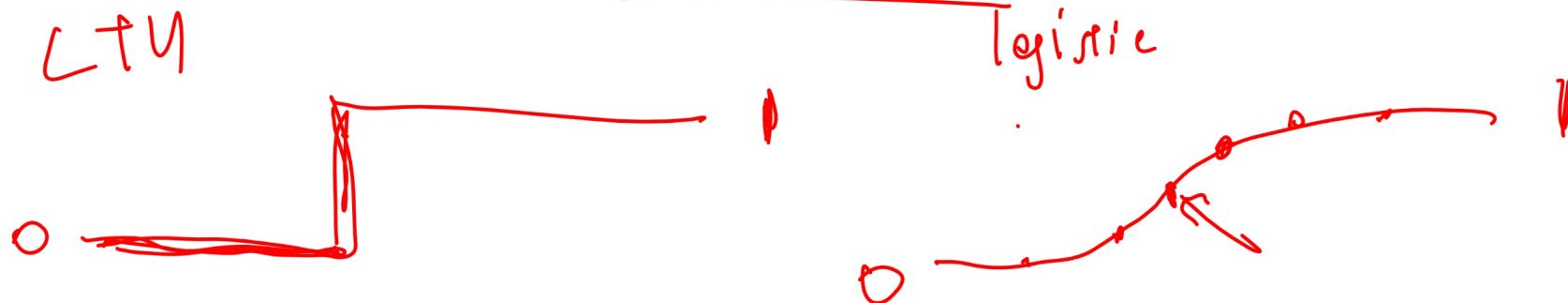
# Sigmoid Activation Function

## Motivation

- When the activation function  $g$  is the sigmoid function, the problem is called logistic regression.

$$g(\square) = \frac{1}{1 + \exp(-\square)}$$


- This  $g$  is also called the logistic function.







# Cross Entropy Loss Function

## Motivation

- The cost function used for logistic regression is usually the log cost function.

$$C(f) = - \sum_{i=1}^n (y_i \log(f(x_i)) + (1 - y_i) \log(1 - f(x_i)))$$

- It is also called the cross-entropy loss function.

like country  
# mistake

# Logistic Regression Objective

## Motivation

- The logistic regression problem can be summarized as the following.

$$(\hat{w}, \hat{b}) = \arg \min_{w, b} - \sum_{i=1}^n (y_i \log(a_i) + (1 - y_i) \log(1 - a_i))$$

where  $a_i = \frac{1}{1 + \exp(-z_i)}$  and  $z_i = w^T x_i + b$

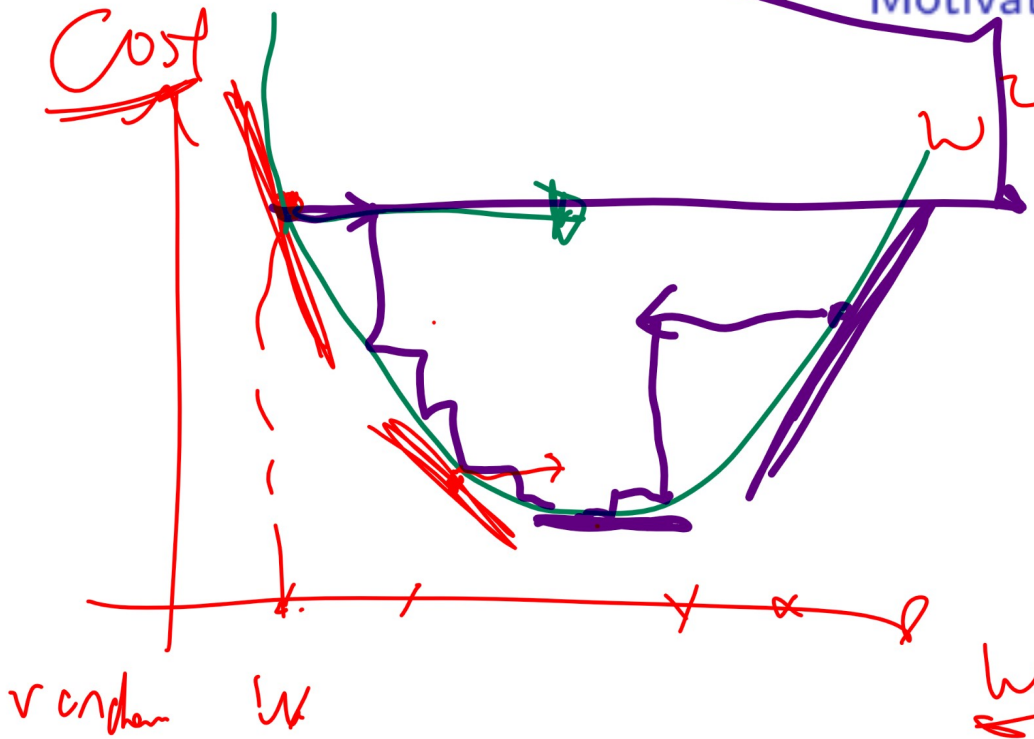
Cost

linear part

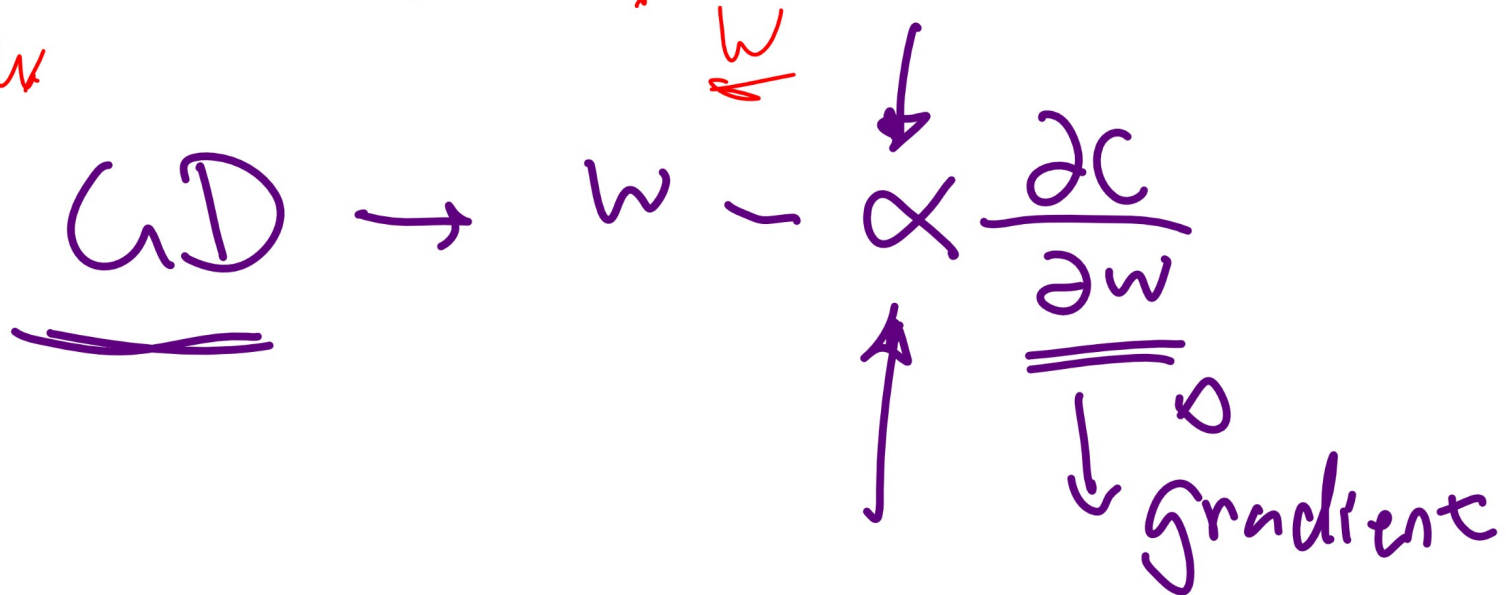
logistic.

# Optimization Diagram

Motivation



$$\frac{\partial C}{\partial w} \leq 0 \Rightarrow w \uparrow$$
$$\frac{\partial C}{\partial w} \geq 0 \Rightarrow w \downarrow$$



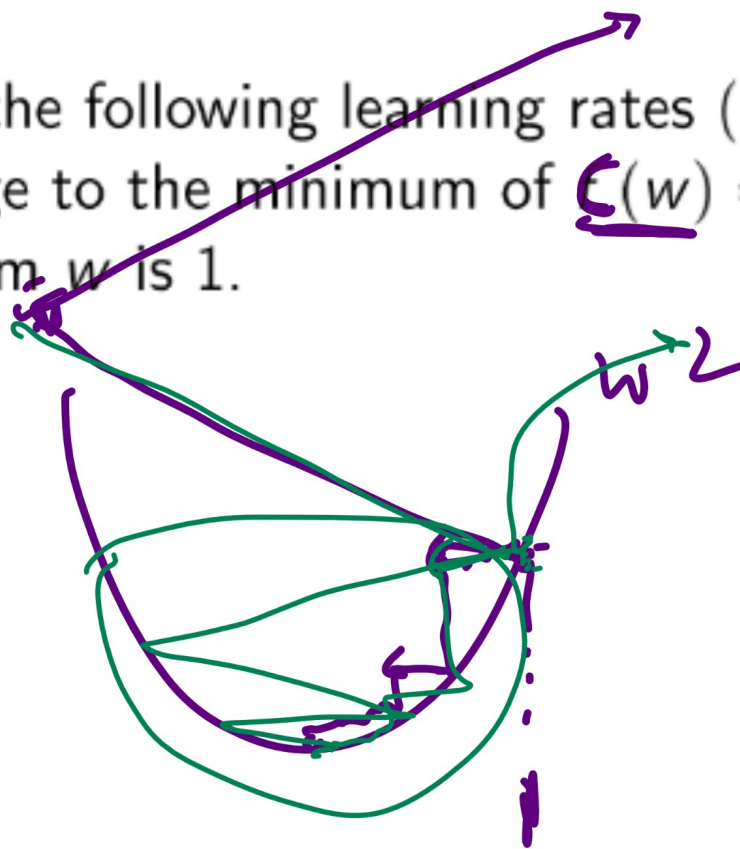


# Simple Gradient Descent

## Quiz

- For which (multiple) of the following learning rates ( $\alpha$ ) would gradient descent converge to the minimum of  $\mathcal{L}(w) = w^2$ . Suppose the initial random  $w$  is 1.

- A: 0.1
- B: 0.2
- C: 0.5
- D: 1
- E: 2



$$1 - 4 = -3$$

$$\begin{aligned} \text{now } w &= w - \alpha \cdot 2w \\ &= |1 - \alpha \cdot 2 \cdot 1| = 0.2 \end{aligned}$$



# Simple Gradient Descent, Another One

## Quiz

*select largest*

- For which (~~multiple~~) of the following learning rates ( $\alpha$ ) would gradient descent converge to the minimum of  $\mathcal{L}(w) = \frac{1}{2}w^2$ .  
Suppose the initial random  $w$  is 1.

- A: 0.1
- B: 0.2
- C: 0.5
- D: 1
- E: 2

$$w = w - \alpha \frac{\partial \mathcal{L}}{\partial w}$$

$$\left| \frac{1 - \alpha}{2} \right| < 1$$







# Logistic Regression

## Description

- Initialize random weights.
- Evaluate the activation function.
- Compute the gradient of the cost function with respect to each weight and bias.
- Update the weights and biases using gradient descent.
- Repeat until convergent.

Use current weight to predict,

# Logistic Gradient Derivation 1

Definition

$$C = - \sum_{i=1}^n y_i \log(a_i) + (1-y_i) \log(1-a_i)$$

$$a_i = \frac{1}{1 + e^{-(w^T x_i + b)}}$$

Chain Rule

$(x_{i1}, x_{i2}, \dots, x_{im})$

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

$$= \sum_{i=1}^n \left( \frac{y_i}{a_i} - \frac{1-y_i}{1-a_i} \right)$$

$$\left( \frac{+ e^{-(w^T x_i + b)} \cdot x_{ij}}{(1 + e^{-(w^T x_i + b)})^2} \right)$$

$$= \sum_{i=1}^n \underbrace{\left( y_i(1-a_i) - (1-y_i)a_i \right)}_{y_i - a_i} x_{ij} \cdot \underbrace{\frac{1}{1 + e^{-(w^T x_i + b)}}}_{a_i} \cdot \underbrace{\frac{e^{-(w^T x_i + b)}}{1 + e^{-(w^T x_i + b)}}}_{1-a_i} x_{ij}$$

# Logistic Gradient Derivation 2

Definition

$$= \sum_{i=1}^n (a_i - y_i) x_{ij}$$

Cross entropy loss

$$w = w - \alpha \sum_{i=1}^n (a_i - y_i) x_i$$

$$\begin{pmatrix} w_1 \\ w_2 \end{pmatrix} = \begin{pmatrix} w_1 \\ w_2 \end{pmatrix} -$$

$$\begin{pmatrix} x_{i1} \\ x_{i2} \end{pmatrix}$$

# Gradient Descent Step

## Definition

- For logistic regression, use chain rule twice.

$$w = w - \alpha \sum_{i=1}^n (a_i - y_i) x_i$$

$$b = b - \alpha \sum_{i=1}^n (a_i - y_i)$$

$$a_i = g(w^T x_i + b), g(\square) = \frac{1}{1 + \exp(-\square)}$$

- $\alpha$  is the learning rate. It is the step size for each step of gradient descent.

# Perceptron Algorithm

## Definition

- Update weights using the following rule.

$$w = w - \alpha (a_i - y_i) x_i$$

$$b = b - \alpha (a_i - y_i)$$

$$a_i = \mathbb{1}_{\{w^T x_i + b \geq 0\}}$$



# Gradient Descent

## Quiz

- What is the gradient descent step for  $w$  if the objective (cost) function is the squared error? Q 11

$$C = \frac{1}{2} \sum_{i=1}^n (a_i - y_i)^2, a_i = g(w^T x_i + b), g'(z) = z \cdot (1 - z)$$

$$a_i = \frac{1}{1 + e^{-w^T x_i + b}}$$

$$\frac{\partial C}{\partial a_i} = (a_i - y_i)$$

$$\frac{\partial a_i}{\partial w} = a_i (1 - a_i) x_i$$

- A:  $w = w - \alpha \sum (a_i - y_i)$
- B:  $w = w - \alpha \sum (a_i - y_i) x_i$  ← cross entropy
- C:  $w = w - \alpha \sum (a_i - y_i) a_i x_i$
- D:  $w = w - \alpha \sum (a_i - y_i) (1 - a_i) x_i$
- E:  $w = w - \alpha \sum (a_i - y_i) a_i (1 - a_i) x_i$

$$\frac{\partial C}{\partial a_i}$$

$$\frac{\partial a_i}{\partial w}$$





# Gradient Descent, Another One

## Quiz

- What is the gradient descent step for  $w$  if the activation function is the identity function?

Linear regression

$$C = \frac{1}{2} \sum_{i=1}^n (a_i - y_i)^2, \quad a_i = w^T x_i + b$$

~~$a_i = y_i$~~   
 $w_1 x_1 + w_2 x_2$

Q12

- A:  $w = w - \alpha \sum (a_i - y_i)$
- B:  $w = w - \alpha \sum (a_i - y_i) x_i$
- C:  $w = w - \alpha \sum (a_i - y_i) a_i x_i$
- D:  $w = w - \alpha \sum (a_i - y_i) (1 - a_i) x_i$
- E:  $w = w - \alpha \sum (a_i - y_i) a_i (1 - a_i) x_i$

$$\frac{\partial C}{\partial w} = \sum \frac{\partial C}{\partial a_i} \cdot \frac{\partial a_i}{\partial w}$$

$(a_i - y_i) x_i$



# Other Non-linear Activation Function

## Discussion

- Activation function:  $g(\square) = \tanh(\square) = \frac{e^{\square} - e^{-\square}}{e^{\square} + e^{-\square}}$
- Activation function:  $g(\square) = \arctan(\square)$
- Activation function (rectified linear unit):  $g(\square) = \square \mathbb{1}_{\{\square \geq 0\}}$
- All these functions lead to objective functions that are convex and differentiable (almost everywhere). Gradient descent can be used.

