

# CS540 Introduction to Artificial Intelligence

## Lecture 2

Young Wu

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

June 28, 2022

# Two-thirds of the Average Game

Quiz

AI

Socratic room

CS540C

- Pick an integer between 0 and 100 (including 0 and 100) that is the closest to two-thirds of the average of the numbers other people picked.

# Quizzes, Math Homework, Discussions

Admin

- Due dates: Monday, late submission withint a week or so without penalty (regrade requests).
- Share solutions (M2 etc): before due date (one or two days late is okay).
- Share solutions (X1 etc): a week before the exam.
- Group discussions: no due dates.

0.5

# Office Hours, Discussion Sessions

Admin

recorded  
on Zoom

- Answer  $M, P$  homework questions on Saturday evenings?
- A: Yes, I will attend.
- B: Yes, but I will not attend.
- C: No.

# Supervised Learning

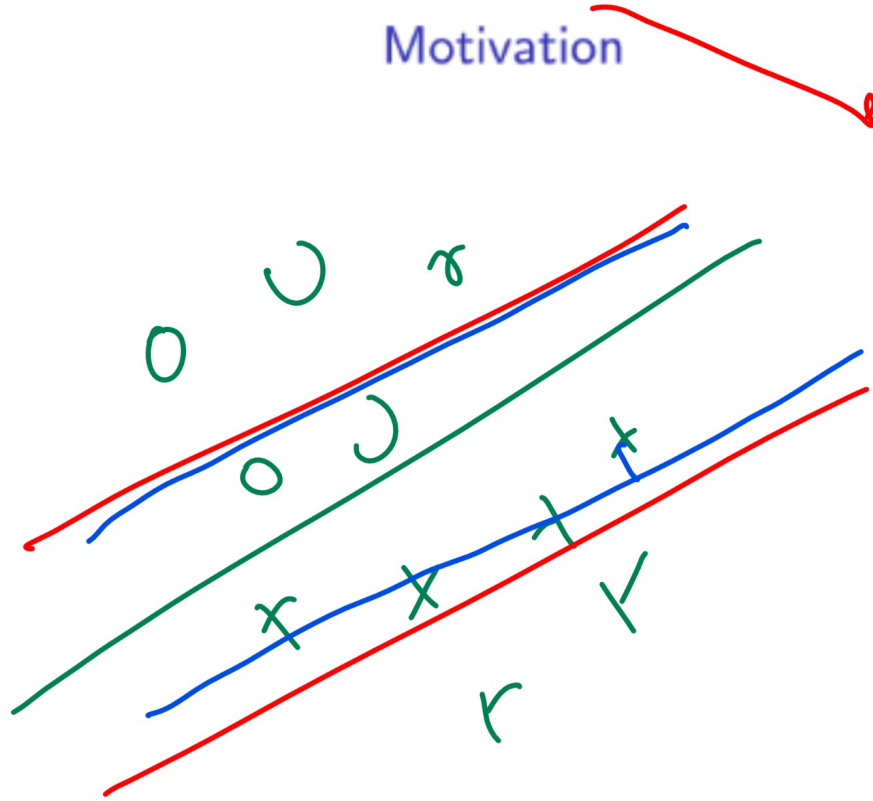
## Motivation

LTU

Data	Features	Labels	-
Training	$\{(x_{i1}, \dots, x_{im})\}_{i=1}^{n'}$	$\{y_i\}_{i=1}^{n'}$	find "best" $\hat{f}$
-	observable	known	-
Test	$(x'_1, \dots, x'_m)$	$y'$	guess $\hat{y} = \hat{f}(x')$
-	observable	unknown	-

# Loss Function Diagram

Motivation



count # mistakes.

# Zero-One Loss Function

## Motivation

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

$$\hat{f} = \operatorname{argmin}_f \sum_{i=1}^n \mathbb{1}_{\{f(x_i) \neq y_i\}}$$


$$\begin{cases} 1 & \text{if } f(x_i) \neq y_i \\ 0 & \text{if } f(x_i) = y_i \end{cases}$$

- $\operatorname{argmin}_f$  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:

$$\hat{f} = \operatorname{argmin}_f \frac{1}{2} \sum_{i=1}^n (f(x_i) - y_i)^2$$




# Loss Functions Equivalence

## Quiz

- Which one of the following functions is not equivalent to the squared error for binary classification?

Q3

$$C = \sum_{i=1}^n (f(x_i) - y_i)^2, f(x_i) \in \{0, 1\}, y_i \in \{0, 1\}$$

*[0, 1]*

- ~~A~~:  $\sum \mathbb{1}_{\{f(x_i) \neq y_i\}}$
- B**:  $\sum \mathbb{1}_{\{f(x_i) = y_i\}}$
- ~~C~~:  $\sum |f(x_i) - y_i|$
- D**:  $\sum \max\{0, 1 - f(x_i) y_i\}$
- ~~E~~:  $\sum \frac{1}{2} \max\{0, 1 - (2 \cdot f(x_i) - 1)(2 \cdot y_i - 1)\}$

	$y_i$	$f(x_i)$
$\mathbb{1}_{\{f(x_i) \neq y_i\}}$	0	0
$\mathbb{1}_{\{f(x_i) = y_i\}}$	0	1
$ f(x_i) - y_i $	1	1
$\max\{0, 1 - f(x_i) y_i\}$	1	0
$\frac{1}{2} \max\{0, 1 - (2 \cdot f(x_i) - 1)(2 \cdot y_i - 1)\}$	1	1

choice

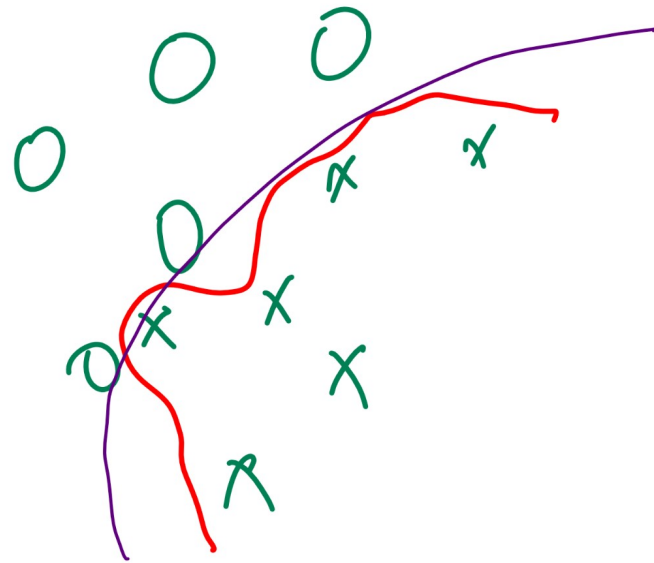
sq	A	B
0	0	1
1	1	0
1	1	0
0	0	1

# Loss Functions Equivalence, Answer

## Quiz

# 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.

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

- 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.

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

where  $a_i = g(w^T x + b)$

- $g$  is called the activation function.

# Linear Threshold Unit

## Motivation

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

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

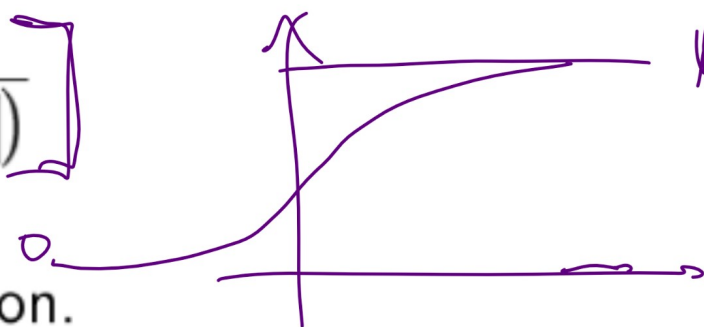
- This activation function is called linear threshold unit (LTU).

# Sigmoid Activation Function

## Motivation

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

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



- 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.

Amount  
Count ⊕  
mistakes

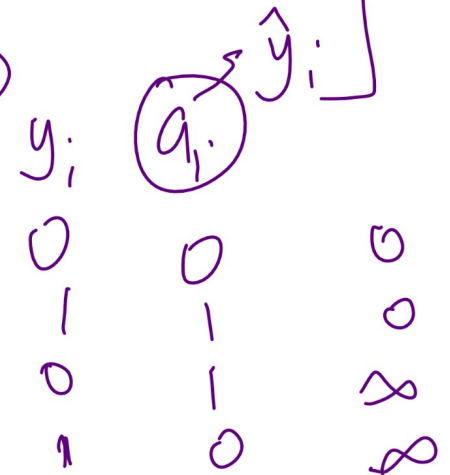
# Logistic Regression Objective

## Motivation

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

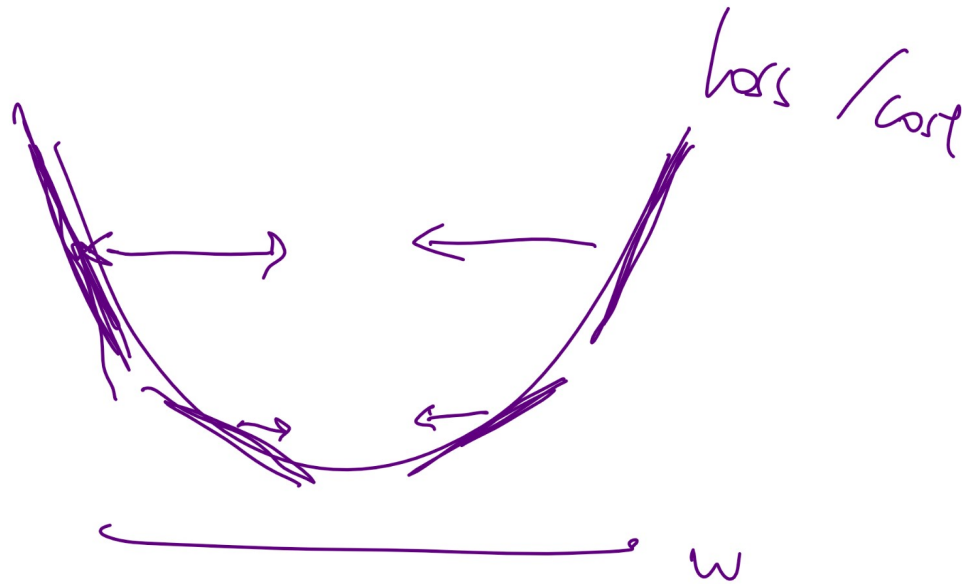
$$(\hat{w}, \hat{b}) = \operatorname{argmin}_{\underline{w, b}} \sum_{i=1}^n (y_i \log(a_i) + (1 - y_i) \log(1 - a_i))$$

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



# Optimization Diagram

## Motivation



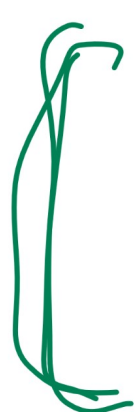
Opposite direction  
of derivative  
(gradient)



# Gradient Descent Step

## Definition

- For logistic regression, use chain rule twice.



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

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

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



CE  
 $y_i \log a_i + (1 - y_i) \log(1 - a_i)$

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

# Perceptron Algorithm

## Definition

- Update weights using the following rule.

$$\left[ \begin{array}{l} w = w - \alpha (a_i - y_i) x_i \\ b = b - \alpha (a_i - y_i) \\ a_i = \mathbb{1}_{\{w^T x_i + b \geq 0\}} \end{array} \right]$$

*geometric*



*NOT*

*gradient  
descent*

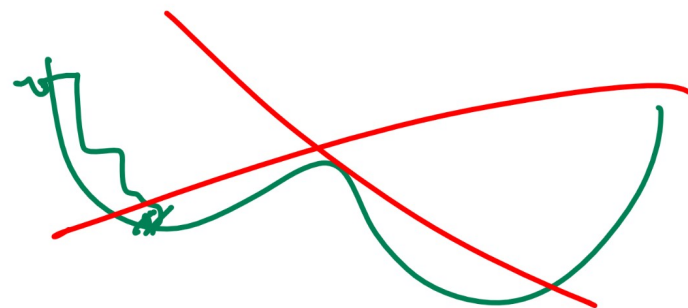
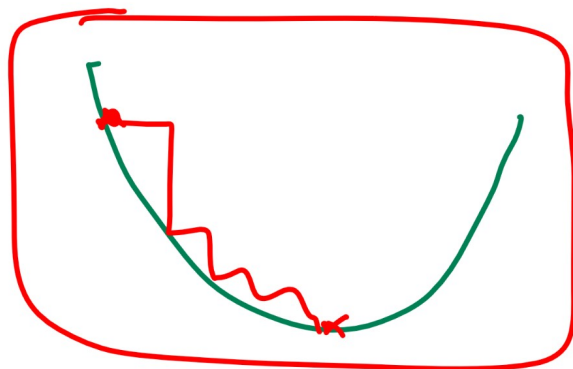
# Learning Rate Diagram

## Definition

# Other Non-linear Activation Function

## Discussion

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





# Gradient Descent

## Quiz

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

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

$w_1 x_{i1} + w_2 x_{i2} + w_3 x_{i3} + \dots$

$$\frac{\partial C}{\partial w_j} = \sum_{i=1}^n \frac{\partial C}{\partial a_i} \frac{\partial a_i}{\partial w_j} = \sum_{i=1}^n (a_i - y_i) \cdot a_i (1 - a_i) \cdot x_{ij}$$

$$w_j = w_j - \frac{\partial C}{\partial w_j}$$

$g'(z)$

# Gradient Descent, Answer

$$\begin{aligned}
 \frac{\partial}{\partial z} \underbrace{\left( \frac{1}{1 + e^{-z}} \right)}_{a_i} &= \frac{+e^{-z} \overset{\text{Quiz}}{}}{(1 + e^{-z})^2} = \frac{e^{-z}}{1 + e^{-z}} \cdot \frac{1}{1 + e^{-z}} \\
 &= \left( 1 - \frac{1}{1 + e^{-z}} \right) \left( \frac{1}{1 + e^{-z}} \right) \\
 &= (1 - a_i) a_i
 \end{aligned}$$



# Gradient Descent

## Quiz

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

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

- 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 :  $w = w - \alpha \sum (a_i - y_i) a_i (1 - a_i) x_i$





# Gradient Descent, Another One Too, Answer

Quiz

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

$$\frac{\partial C}{\partial w_j} = \frac{\partial}{\partial w_j} \sum_{i=1}^n \left( \sum_{j=1}^m w_j x_j + b - y_i \right)^2$$

$$= \sum_{i=1}^n \left( \frac{\partial C}{\partial a_i} \right) \frac{\partial a_i}{\partial w_j} \rightarrow w_1 x_{i1} + w_2 x_{i2} + w_3 x_{i3}$$

$$= \sum_{i=1}^n \frac{1}{2} \cdot 2 (a_i - y_i) x_{ij}$$

i-th item  
 j-th feature.

$$\nabla_w C = \begin{bmatrix} \frac{\partial C}{\partial w_1} \\ \frac{\partial C}{\partial w_2} \\ \frac{\partial C}{\partial w_3} \end{bmatrix} = \sum_{i=1}^n (a_i - y_i) \begin{pmatrix} x_{i1} \\ x_{i2} \\ \vdots \\ x_{i3} \end{pmatrix} = \sum_{i=1}^n (a_i - y_i) X_i$$

# Convexity Diagram

## Discussion