# CS540 Introduction to Artificial Intelligence
# Lecture 2

Young Wu

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

June 22, 2021

**Generalized Linear Models**
●○○○○○○○○○○○○

Logistic Regression
○○○○○○○○○○○

Gradient Descent
○○○○○○○○

# Two-thirds of the Average Game

## Quiz

*Socrative*

*CS 540*

*Wisc 2D*

*Q1*

- Pick the number that is the closest to two-thirds of the average of the numbers other people picked.

  - A: 0
  - B: 1
  - C: 2
  - D: 3
  - E: 4

$$avg = 1.2$$

$$\frac{2}{3}\, avg = 0.8$$

**Generalized Linear Models**
○●○○○○○○○○○○○

Logistic Regression
○○○○○○○○○○○

Gradient Descent
○○○○○●○○

# Remind Me to Start Recording
## Admin

- The annotated slides are posted on Q1, Q2, etc.
- The lecture recordings are shared on Canvas, are they visible?
- The messages you send in chat will be recorded: you can change your Zoom name now before I start recording.

**Generalized Linear Models**
○○●○○○○○○○○○○

Logistic Regression
○○○○○○○○○○○

Gradient Descent
○○○○○○●○

# Two-thirds of the Average Game
## Quiz

A1

- 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.

**Generalized Linear Models**

○○○○●○○○○○○○

Logistic Regression

○○○○○○○○○○○

Gradient Descent

○○○○○○○●

# Supervised Learning Example
## Motivation

| Data | images of cats and dogs |
|---|---|
| Features (Input) | height, length, eye color, ... |
| Output | cat or dog |

| Data | emails |
|---|---|
| Features (Input) | word count, capitalization, ... |
| Output | spam or ham |

**Generalized Linear Models**
○○○○●○○○○○○

Logistic Regression
○○○○○○○○○○○

Gradient Descent
○○○○○○○○

# Supervised Learning
## Motivation

$x_i$

0   ham
1   spam

| Data | Features (Input) | Output | - |
|---|---|---|---|
| Training | $\{(x_{i1}, ..., x_{im})\}_{i=1}^{n'}$ | $\{y_i\}_{i=1}^{n'}$ | find "best" $\hat{f}$ |
| - | observable | known | - |
| Test | $(x_1', ..., x_m')$ | $y'$ | guess $\hat{y} = \hat{f}(x')$ |
| - | observable | unknown | - |

**Generalized Linear Models**
○○○○○●○○○○○○

Logistic Regression
○○○○○○○○○○○

Gradient Descent
○○○○○○○○

# Loss Function Diagram

## Motivation



best → objective

find $\hat{f}$   minimize   #l mistake   or training

= 2

cost

loss

$\hat{f}$

**Generalized Linear Models**
○○○○○○●○○○○○

**Logistic Regression**
○○○○○○○○○○○

**Gradient Descent**
○○○○○○○○

# Zero-One Loss Function

## Motivation

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

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

$$\mathbb{1}_E = \begin{cases} 1 & E \text{ true} \\ 0 & E \text{ false} \end{cases}$$

Loss

$$\begin{cases} 1 & \text{mistake} \\ 0 & \text{correct prediction} \end{cases}$$

- $\arg\min_{f}$ 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.

**Generalized Linear Models**
ooooooooo●oooo

Logistic Regression
ooooooooooo

Gradient Descent
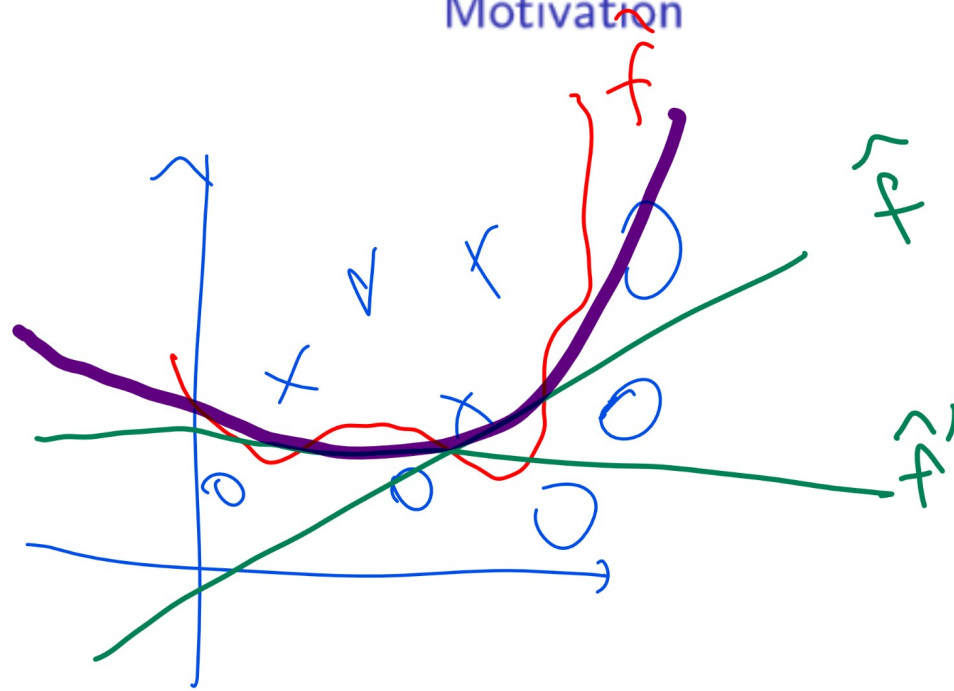ooooooooo

# 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:

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

**Generalized Linear Models**
OOOOOOOOO●OOO

Logistic Regression
OOOOOOOOOOO

Gradient Descent
OOOOOOOO

# Function Space Diagram

## Motivation

**Generalized Linear Models**
○○○○○○○○○●○○

Logistic Regression
○○○○○○○○○○○

Gradient Descent
○○○○○○○○

# Hypothesis Space
## Motivation

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

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

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

**Generalized Linear Models**
OOOOOOOOOOO●O

Logistic Regression
OOOOOOOOOOO

Gradient Descent
OOOOOOOO

# Activation Function
## Motivation

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

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

$$f(x_i)$$

$$\text{where } a_i = g\left(w^T x_i + b\right)$$

- $g$ is called the activation function.

**Generalized Linear Models**
○○○○○○○○○○○○●

Logistic Regression
○○○○○○○○○○○

Gradient Descent
○○○○○○○○

# Linear Threshold Unit

## Motivation

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

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

$$w^T x + b$$

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

Generalized Linear Models
○○○○○○○○○○○○○○

Logistic Regression
●○○○○○○○○○○○

Gradient Descent
○○○○○○○○

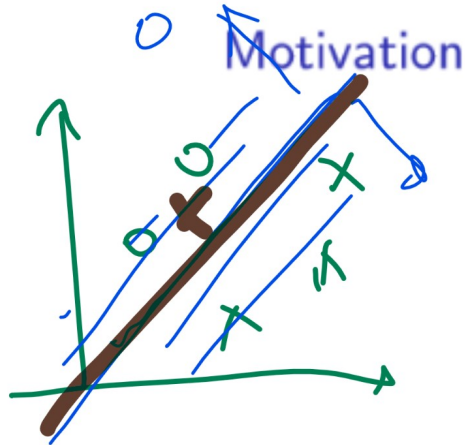# Sigmoid Activation Function

## Motivation

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

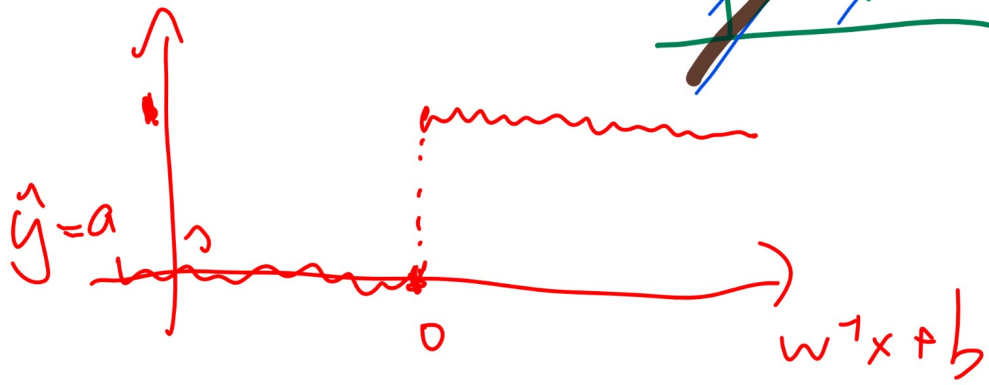$$g\left(\boxed{\cdot}\right) = \frac{1}{1 + \exp\left(-\boxed{\cdot}\right)}$$

$$w^T x + b$$

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

Generalized Linear Models
○○○○○○○○○○○○○

Logistic Regression
○●○○○○○○○○○

Gradient Descent
○○○○○○○○

# Sigmoid Function Diagram

Motivation

$LTU$

$0$

$\hat{y} = a$

$0$

$w^T x + b$

logistic

$a_i = 0.99$
$\hat{y} = 1$

$1$

$0$

$0$

$w^T x + b$

$a \neq \hat{y}$

Prob that label is 1

Perceptron Algorithm

Gradient Descent

Generalized Linear Models
○○○○○○○○○○○○○

Logistic Regression
○○●○○○○○○○○

Gradient Descent
○○○○○○○○

# 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.

Convex

Generalized Linear Models
○○○○○○○○○○○○○

Logistic Regression
○○○●○○○○○○

Gradient Descent
○○○○○○○○

# Logistic Regression Objective
## Motivation

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

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

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

Generalized Linear Models
ooooooooooooo

Logistic Regression
ooooo●ooooo

Gradient Descent
oooooooo

# Optimization Diagram

## Motivation



Loss / cost

loss

slope

w, b

random w

want w*, b*

Convex

**Generalized Linear Models**
ooooooooooooo

**Logistic Regression**
ooooo●ooooo

**Gradient Descent**
oooooooo

# Learning Rate Demo

## Motivation

Cost

negative
direction of gradient

slope

$$w = w - \boxed{\alpha} (a_i - y_i) x_i$$

fraction

Generalized Linear Models
○○○○○○○○○○○○○

Logistic Regression
○○○○○○●○○○○

Gradient Descent
○○○○○○○○

# 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.

Generalized Linear Models
○○○○○○○○○○○○○

Logistic Regression
○○○○○○○●○○○

Gradient Descent
○○○○○○○○

# Logistic Gradient Derivation 1

CE
↙ Lose

## Definition

→ natural $\lg$ ← ln

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

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

logistic activation

$j = 1, 2 \dots m$

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

— Chain rule.

Same $w$ for all $i$

$$= -\sum_{i=1}^{n} \left( \frac{y_i}{a_i} - \frac{1-y_i}{1-a_i} \right) \left( \frac{e^{-(w^T x_i + b)}}{\left(1 + e^{-(w^T x_i + b)}\right)^2} \cdot x_{ij} \right)$$

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

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

Generalized Linear Models
○○○○○○○○○○○○○

Logistic Regression
○○○○○○○○●○○

Gradient Descent
○○○○○○○○

# Logistic Gradient Derivation 2

## Definition

$$w^T x_i = w_1 \boxed{x_{i1}} + w_2 x_{i2} + \cdots$$

$$\boxed{\frac{\partial w^T x_i}{\partial w_j} = x_{ij}}$$

$$\nabla_w C = \begin{bmatrix} \frac{\partial C}{\partial w_1} \\ \frac{\partial C}{\partial w_2} \\ \vdots \\ \frac{\partial C}{\partial w_m} \end{bmatrix}$$

$$= \sum_{i=1}^{n} (a_i - y_i) \begin{bmatrix} x_{i1} \\ x_{i2} \\ \vdots \\ x_{im} \end{bmatrix}$$

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

$$a_i (1 - a_i)$$

$$\frac{1}{1 + e^{-(w^T x + b)}} \cdot \left( 1 - \frac{1}{1 + e^{-(w^T x + b)}} \right)$$

$$\frac{e^{-(w^T x + b)}}{1 + e^{-(w^T x + b)}}$$

$$\frac{e^{-(w^T x + b)}}{1 + e^{-(w^T x + b)}}$$

Generalized Linear Models
○○○○○○○○○○○○○○

Logistic Regression
○○○○○○○○○●○

Gradient Descent
○○○○○○○○

# Gradient Descent Step

## Definition

- For logistic regression, use chain rule twice.

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

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

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

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

Generalized Linear Models
○○○○○○○○○○○○○

Logistic Regression
○○○○○○○○○○●

Gradient Descent
○○○○○○○○

# Perceptron Algorithm
## Definition

- Update weights using the following rule.

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

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

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

Generalized Linear Models
○○○○○○○○○○○○○

Logistic Regression
○○○○○○○○○○○

Gradient Descent
●○○○○○○○

# Other Non-linear Activation Function
## Discussion

- Activation function: $g\left(\boxed{\cdot}\right) = \tanh\left(\boxed{\cdot}\right) = \dfrac{e^{\boxed{\cdot}} - e^{-\boxed{\cdot}}}{e^{\boxed{\cdot}} + e^{-\boxed{\cdot}}}$

- Activation function: $g\left(\boxed{\cdot}\right) = \arctan\left(\boxed{\cdot}\right)$

- Activation function (rectified linear unit): $g\left(\boxed{\cdot}\right) = \boxed{\cdot}\,\mathbb{1}_{\left\{\boxed{\cdot} \geq 0\right\}}$

  *ReLU*

- All these functions lead to objective functions that are convex and differentiable (almost everywhere). Gradient descent can be used.

# Loss Functions Equivalence

## Quiz

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

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

- 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: $\sum \frac{1}{2} \max\{0, 1 - (2 \cdot f(x_i) - 1)(2 \cdot y_i - 1)\}$

Generalized Linear Models
OOOOOOOOOOOOO

Logistic Regression
OOOOOOOOOOOO

**Gradient Descent**
OO●OOOOO

# Loss Functions Equivalence, Answer

## Quiz

Generalized Linear Models
○○○○○○○○○○○○○○

Logistic Regression
○○○○○○○○○○○

**Gradient Descent**
○○○●○○○○○

# Gradient Descent
## Quiz

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

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

*(handwritten: square loss)*

*(handwritten: Q.3)*

$$w = w - \alpha \nabla_w C = \begin{pmatrix} \frac{\partial C}{\partial w_1} \\ \vdots \\ \frac{\partial C}{\partial w_n} \end{pmatrix}$$

- A: $w = w - \alpha \sum (a_i - y_i)$
- B: $w = w - \alpha \sum (a_i - y_i) x_i$ $\longleftarrow$ logistic with CE *(handwritten: 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: $w = w - \alpha \sum (a_i - y_i) a_i (1 - a_i) x_i$

Generalized Linear Models

○○○○○○○○○○○○

Logistic Regression

○○○○○○○○○○○

**Gradient Descent**

○○○○○●○○○○

# Gradient Descent, Answer

## Quiz

Generalized Linear Models
ооооооооооооо

Logistic Regression
ооооооооооо

Gradient Descent
ооооооооо

# Gradient Descent, Another One
## Quiz

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

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

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

Generalized Linear Models
○○○○○○○○○○○○○

Logistic Regression
○○○○○○○○○○○○

Gradient Descent
○○○○○○○●○

# Gradient Descent, Another One, Answer

## Quiz

Generalized Linear Models
ooooooooooooo

Logistic Regression
ooooooooooo

Gradient Descent
ooooooo●

# Remind Me to Stop Recording
## Admin

- If you accidentally selected an obviously incorrect answer earlier, you can enter the question name and the correct answer here.