

CS540 Introduction to Artificial Intelligence

Lecture 2

Young Wu

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

June 3, 2020

Test

Quiz

$$\frac{2}{3} \frac{1}{n} \sum_{i=1}^n x_i$$

- 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

← 0.85

Socratic Room

CS540C

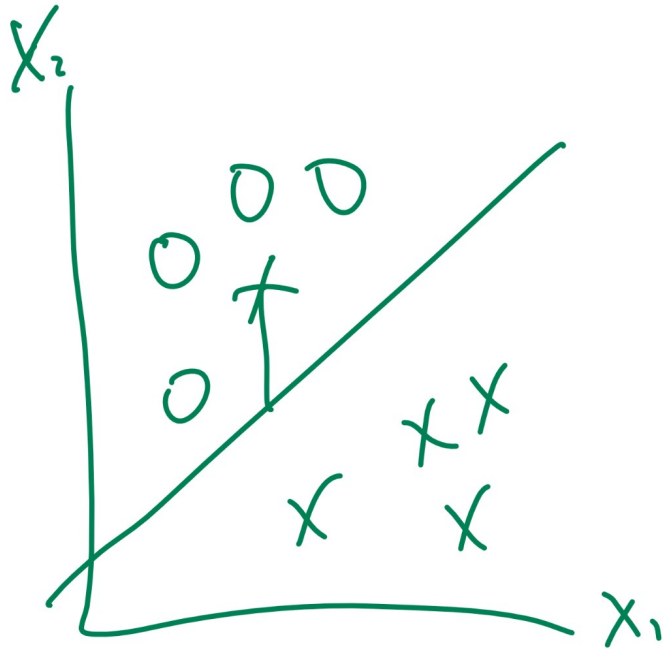
@wisc.edu

↑ ID

Send all questions to "Questions" on public.

Loss Function Diagram

Motivation



min # mistakes.
LTV Perception
min loss / cost

Zero-One Loss Function

Motivation

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

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

1 mistake
 0 mistake

- $\arg \min_f$ 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:

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

prediction label training

Loss Functions Equivalence

Quiz

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

Q2

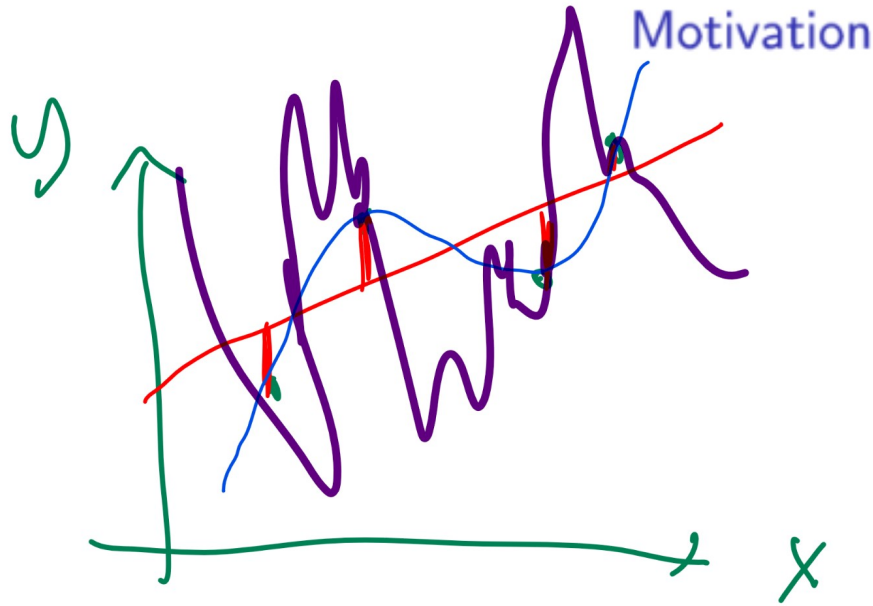
$$C = \sum_{i=1}^n (f(x_i) - y_i)^2, y_i \in \{0, 1\}$$

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

Handwritten analysis and diagrams:

- Diagram 1: A table with two columns. The left column is labeled $(f(x_i) - y_i)^2$ and contains values 0, 1, 1, 0. The right column is labeled $f(x_i)$ and contains values 0, 0, 1, 1. A green box highlights the first two rows, and a blue arrow points from the text $y_i \in \{0, 1\}$ to the first row.
- Diagram 2: A table with two columns. The left column is labeled y_i and contains values 0, 1, 0, 1. The right column is labeled $f(x_i) + y_i$ and contains values 0, 1, 1, 0. A red box highlights the first two rows, and a blue arrow points from the text $y_i \in \{0, 1\}$ to the first row.
- Diagram 3: A table with two columns. The left column is labeled y_i and contains values 0, 1, 0, 1. The right column is labeled $f(x_i) + y_i$ and contains values 0, 1, 1, 0. A red box highlights the first two rows, and a blue arrow points from the text $y_i \in \{0, 1\}$ to the first row.
- Diagram 4: A table with two columns. The left column is labeled y_i and contains values 0, 1, 0, 1. The right column is labeled $f(x_i) + y_i$ and contains values 0, 1, 1, 0. A red box highlights the first two rows, and a blue arrow points from the text $y_i \in \{0, 1\}$ to the first row.

Function Space Diagram



Hypothesis Space

Motivation

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

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

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

bias
 weights
 plane surface
 activation function.

- g is called the activation function.

Linear Threshold Unit

Motivation

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

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

$w^T x + b$

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

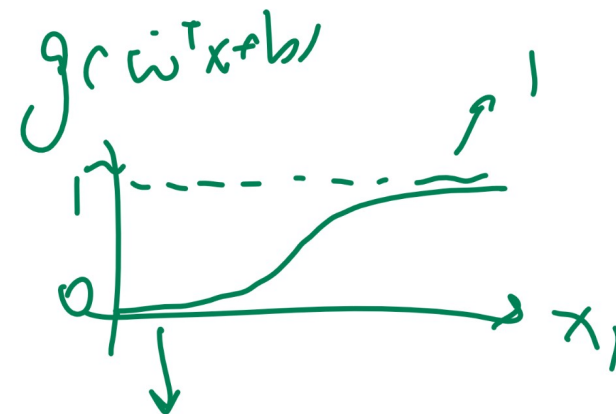
in LI

Sigmoid Activation Function

Motivation

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

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



- This g is also called the logistic function.

Cross Entropy Loss Function

Motivation

$$C = (y_i - f(x_i))^2$$

- 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.

Logistic Regression Objective

Motivation

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

P1

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

Handwritten annotations: '0, 1' with arrows pointing to y_i and $1 - y_i$ respectively.

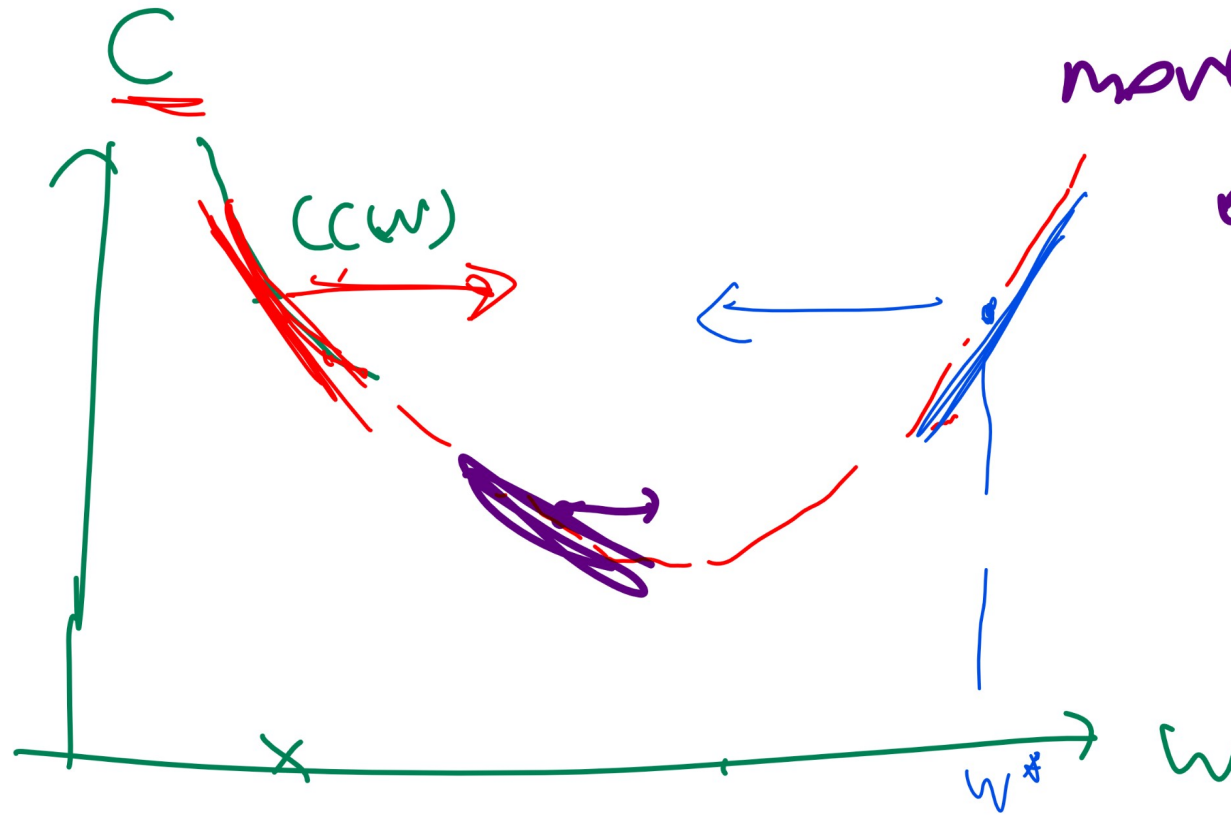
where $a_i = \frac{1}{1 + \exp(-z_i)}$ and $z_i = w^T x_i + b$

Handwritten annotations: A large green bracket underlines the sigmoid function and the linear combination. A green arrow points from the sigmoid function to the a_i term in the objective function above. A green arrow points from the linear combination to the z_i term in the sigmoid function.



Optimization Diagram

Motivation



move in
opposite
direction
of -derivative

Learning Rate Demo

Motivation

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.

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

Logistic Gradient Derivation 1

Definition

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

n → # training

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

w_1, w_2, \dots, w_m

x_1, \dots, x_m

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

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

Logistic Gradient Derivation 2

Definition

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

(Note: In the original image, the terms $y_i - a_i$ and a_i in the denominator are crossed out with green lines, and the entire fraction is circled in blue.)

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

(Note: This equation is enclosed in a blue rounded rectangle.)

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

(Note: In the original image, the fraction $a_i(1-a_i)$ is boxed in red, and the denominator of the second fraction is also boxed in red. A red double-headed arrow indicates the relationship between the two boxes.)

Gradient Descent Step

Definition

- For logistic regression, use chain rule twice.

$P_i \rightarrow$

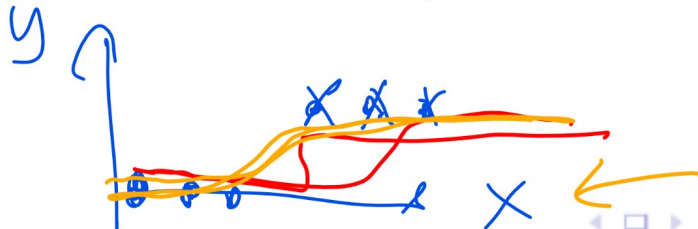
$$w = w - \alpha \sum_{i=1}^n (a_i - y_i) x_i$$

derivative

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

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

- α is the learning rate. It is the step size for each step of gradient descent.



Perceptron Algorithm

Definition

- Update weights using the following rule.

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

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

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

geometric
Trick

if $w^T x + b \geq 0$
otherwise

Gradient Descent

Quiz

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

Q3

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

$$\frac{\partial C}{\partial w}$$

$$= \frac{\partial C}{\partial a} \cdot \frac{\partial a}{\partial w}$$

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

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

cross entropy loss.

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

chain rule

Gradient Descent, Answer Quiz

Gradient Descent, Another One

Quiz

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

Q4

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

linear regression
least squares

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

$$\frac{\partial C}{\partial w}$$

$$\frac{\partial C}{\partial w} = \frac{\partial C}{\partial a} \frac{\partial a}{\partial w}$$

$M1 \rightarrow M12$

top 10

$P1 \rightarrow P6$

top 5

Gradient Descent, Another One, Answer

Quiz

~~$g^{-1}(a) \approx w^T x + b$~~

$a_i = w^T x_i + b$
↓

$w_1 x_1 + w_2 x_2 + \dots + b$

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

$$\frac{\partial C}{\partial w_j} = \sum_{i=1}^n \left[\frac{\partial C}{\partial a_i} \right] \left[\frac{\partial a_i}{\partial w_j} \right]$$

$$= \sum_{i=1}^n \left[\frac{1}{2} 2 (a_i - y_i) \right] \left[x_{ij} \right]$$

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

$$\begin{bmatrix} \frac{\partial C}{\partial w_1} \\ \frac{\partial C}{\partial w_2} \end{bmatrix}$$

$$\frac{\partial C}{\partial w} = \nabla_w C = \sum_{i=1}^n (a_i - y_i) X_i$$

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.