

CS540 Introduction to Artificial Intelligence

Lecture 2

Young Wu

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

May 23, 2022

Two-thirds of the Average Game

Quiz

A7

- 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.

Make Up Lectures

Admin

- Next Monday is Memorial Day.
- *A* : No make up lecture
- *B* : Make up lecture next Wednesday
- *C* : Make up lecture next next Wednesday

Supervised Learning

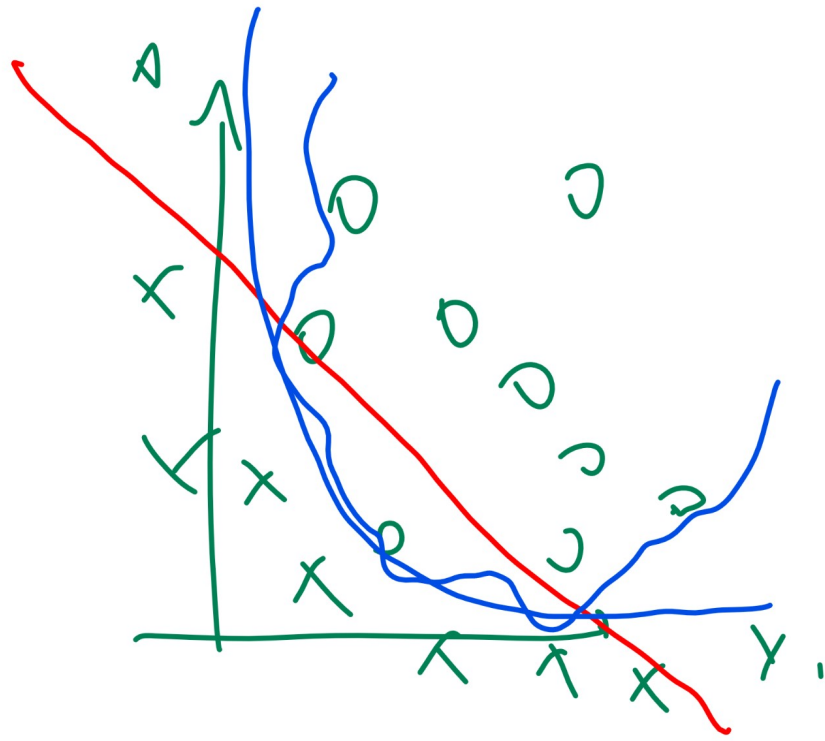
Motivation

$$y \approx w^T x + b$$

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~~



Zero-One Loss Function

Motivation

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

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

Count # incorrect labels by f .

prediction \neq actual

$$\hat{y}_i = f(x_i) \quad y_i$$

loss \sum over items (instances)

- 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

Loss Functions Equivalence

Quiz

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

Q9 $C = \sum_{i=1}^n (f(x_i) - y_i)^2, f(x_i) \in \{0, 1\}, y_i \in \{0, 1\}$

Handwritten annotations: A green bracket underlines the sum. A purple bracket underlines the squared term. A yellow 'C' is written next to the equation. Labels 'f(x_i) predicted' and 'y_i actual' are written below the terms in the equation.

				<i>f(x_i) predicted</i>	<i>y_i actual</i>	<i>square loss</i>	<i>A loss</i>	<i>B loss</i>
<input checked="" type="checkbox"/>	A: $\sum \mathbb{1}_{\{f(x_i) \neq y_i\}}$	*	0	0	0	0	0	1
<input checked="" type="checkbox"/>	B: $\sum \mathbb{1}_{\{f(x_i) = y_i\}}$	~	0	1	1	0	0	1
<input checked="" type="checkbox"/>	C: $\sum f(x_i) - y_i $	~	1	0	1	1	1	0
<input checked="" type="checkbox"/>	D: $\sum \max\{0, 1 - f(x_i) y_i\}$	~	1	0	0	1	1	0
<input checked="" type="checkbox"/>	E: $\sum \frac{1}{2} \max\{0, 1 - (2 \cdot f(x_i) - 1)(2 \cdot y_i - 1)\}$	~	1	1	1	0	0	0

Handwritten annotations: A red 'X' is next to A. A blue 'X' is next to B. A red wavy line is next to C. A yellow box highlights the '1' values in the third column. A red box highlights the '0' values in the sixth column. A blue arrow points from the '0' in the sixth column to the '0' in the eighth column. A purple checkmark is at the bottom right.

Hypothesis Space

Motivation

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

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

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

Activation Function

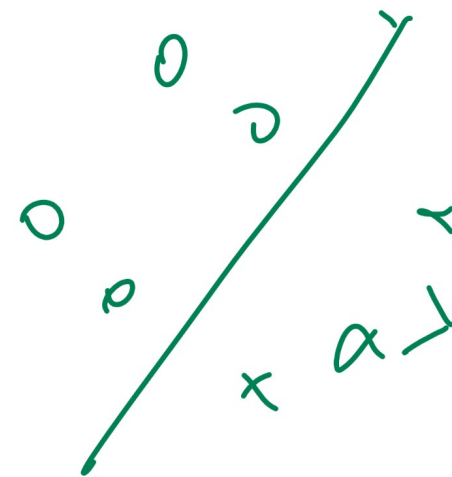
Motivation

- 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~~ ^{indicator} 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

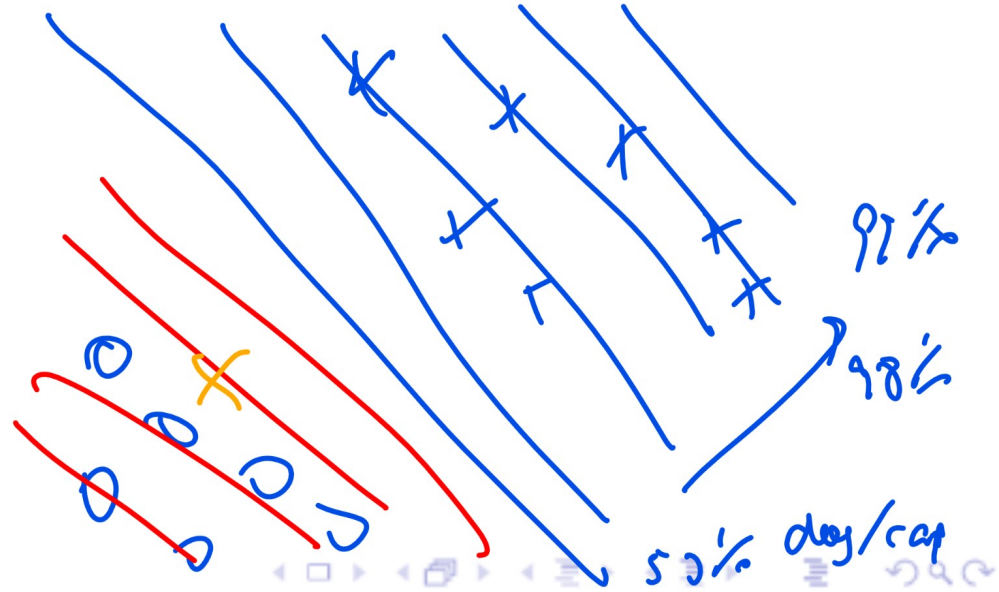
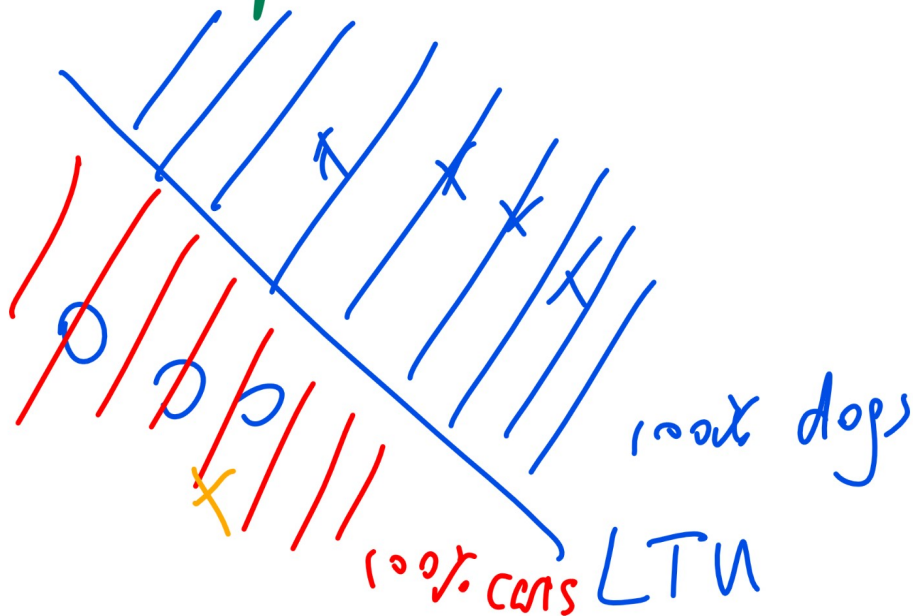
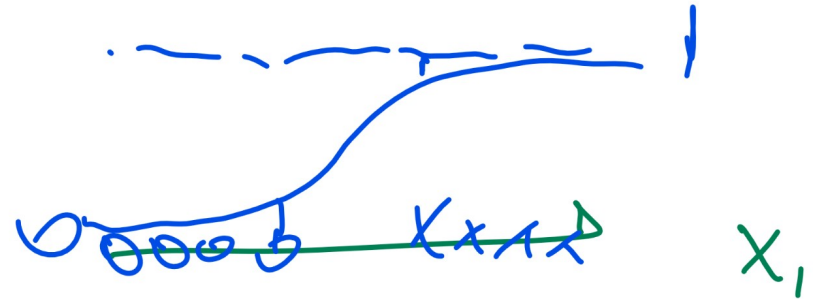
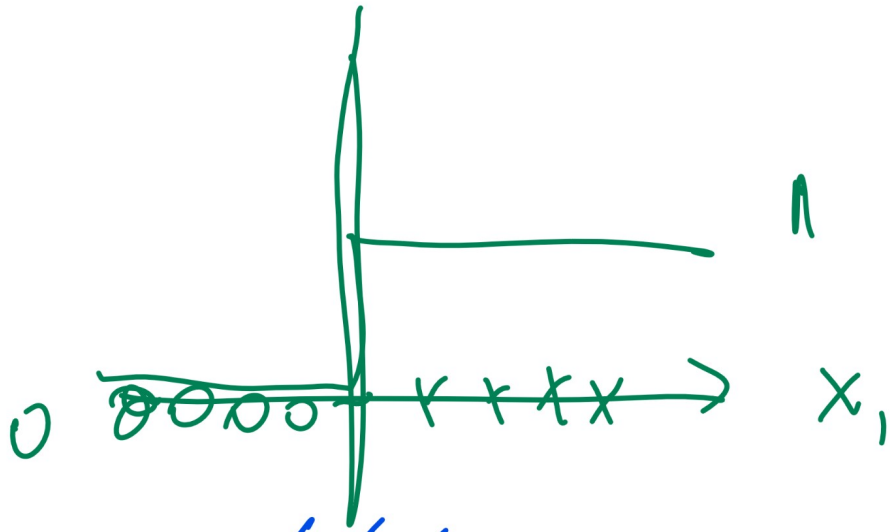
- When the activation function g is the sigmoid function, the problem is called logistic regression.

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

- This g is also called the logistic function.

Sigmoid Function Diagram

Motivation



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.

logistic

Logistic Regression Objective

Motivation

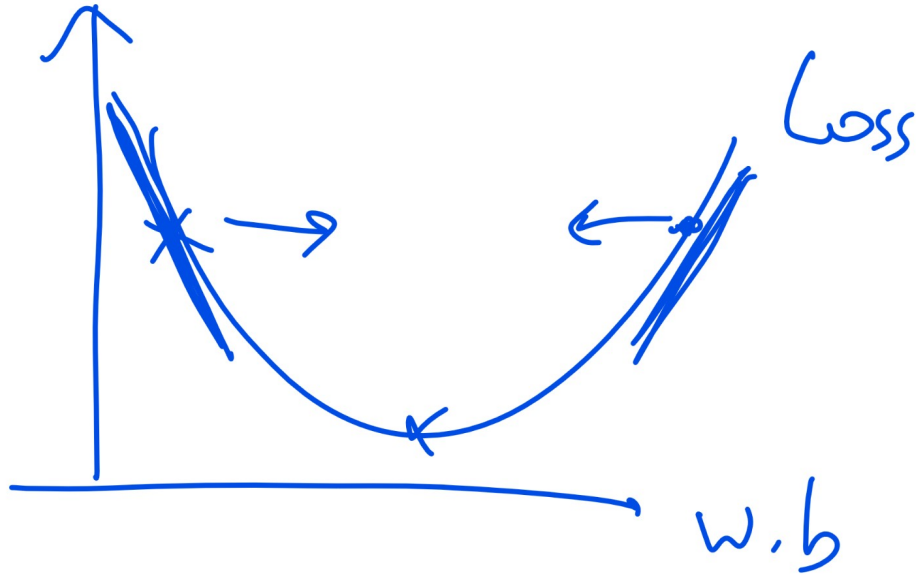
- The logistic regression problem can be summarized as the following.

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

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

Optimization Diagram

Motivation



gradient
↓
derivative

Logistic Regression

Description

- Initialize random weights.
- Evaluate the activation function. $\rightarrow a_i$
- 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.

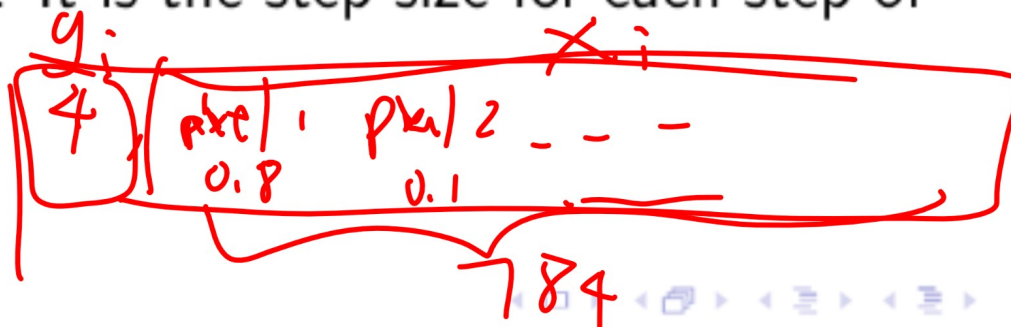
Gradient Descent Step

Definition

- For logistic regression, use chain rule twice.

$$\left. \begin{aligned}
 w &= \underline{w} - \alpha \sum_{i=1}^n (a_i - y_i) x_i \\
 b &= \underline{b} - \alpha \sum_{i=1}^n (a_i - y_i) \\
 a_i &= g(w^T x_i + b), \quad g(\square) = \frac{1}{1 + \exp(-\square)}
 \end{aligned} \right\} \text{gradient descent}$$

- 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.

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

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

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

} Geometry

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.

Gradient Descent

Quiz

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

Cross Entropy

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

Chain Rule

$$w_j = w_j \propto \frac{\partial C}{\partial w_j} = w \sim \sum_{i=1}^n \frac{\partial C}{\partial a_i} \cdot \frac{\partial a_i}{\partial w_j}$$

partial derivative

$$\frac{1}{2} \cdot 2(a_i - y_i) \cdot (1 - a_i) a_i x_j$$

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

Quiz

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

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

Q10.

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

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

$w = w - \alpha \left[\frac{\partial C}{\partial w} \right]$

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

$(a_i - y_i) x_{ij}$

Convexity Diagram

Discussion

Questions?

Admin

10/12

- Missed lectures and quizzes. ✓
- Math used in the course. ✓
- Homework due dates. ✓
- Discussions and sharing solutions. ✓

(x_i, y_i)
pixels $[0, 1]$

