

# CS540 Introduction to Artificial Intelligence

## Lecture 2

Young Wu

Based on lecture slides by Jerry Zhu and Yingyu Liang

May 22, 2019

# Admin

## Admin

- Math and Stat Review posted under W1.
- Complete slides (with diagrams and quiz questions etc) will be posted Thursday or Friday.
- Homework will be posted on Friday (due in 9 days, not 2 days).
- Exact due dates are on Canvas: programming homework can be submitted two weeks late (except for the last two homework (one week late)).

• Lecture this Friday. no lecture Monday

# Activation Function

## Review

- The supervised learning problem with activation function is the following.

$$\left( \hat{w}_0, \hat{w}_1, \dots, \hat{w}_m, \hat{b} \right) = \arg \min_{w_1, \dots, w_m, b} C$$

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

$$\text{and } a_i = g(w_1 x_{i1} + w_2 x_{i2} + \dots + w_m x_{im} + b)$$

$$g = \mathbb{1}_{\{w^T x \geq 0\}} \quad \leftarrow \text{Perceptron Algorithm}$$

# Sigmoid Activation Function

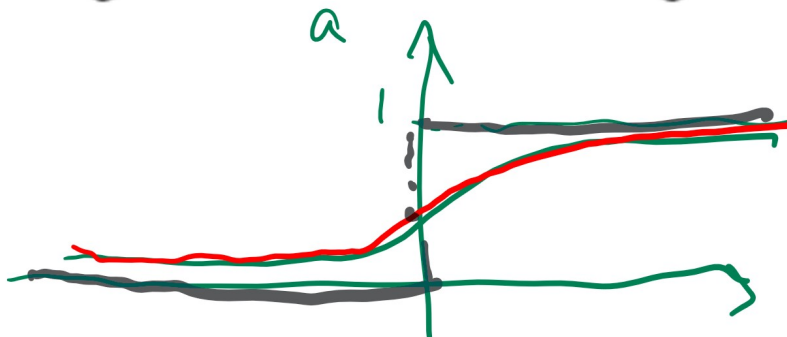
## Motivation

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

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

$\mathbb{1}_{\{\boxed{\cdot} \geq 0\}}$

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

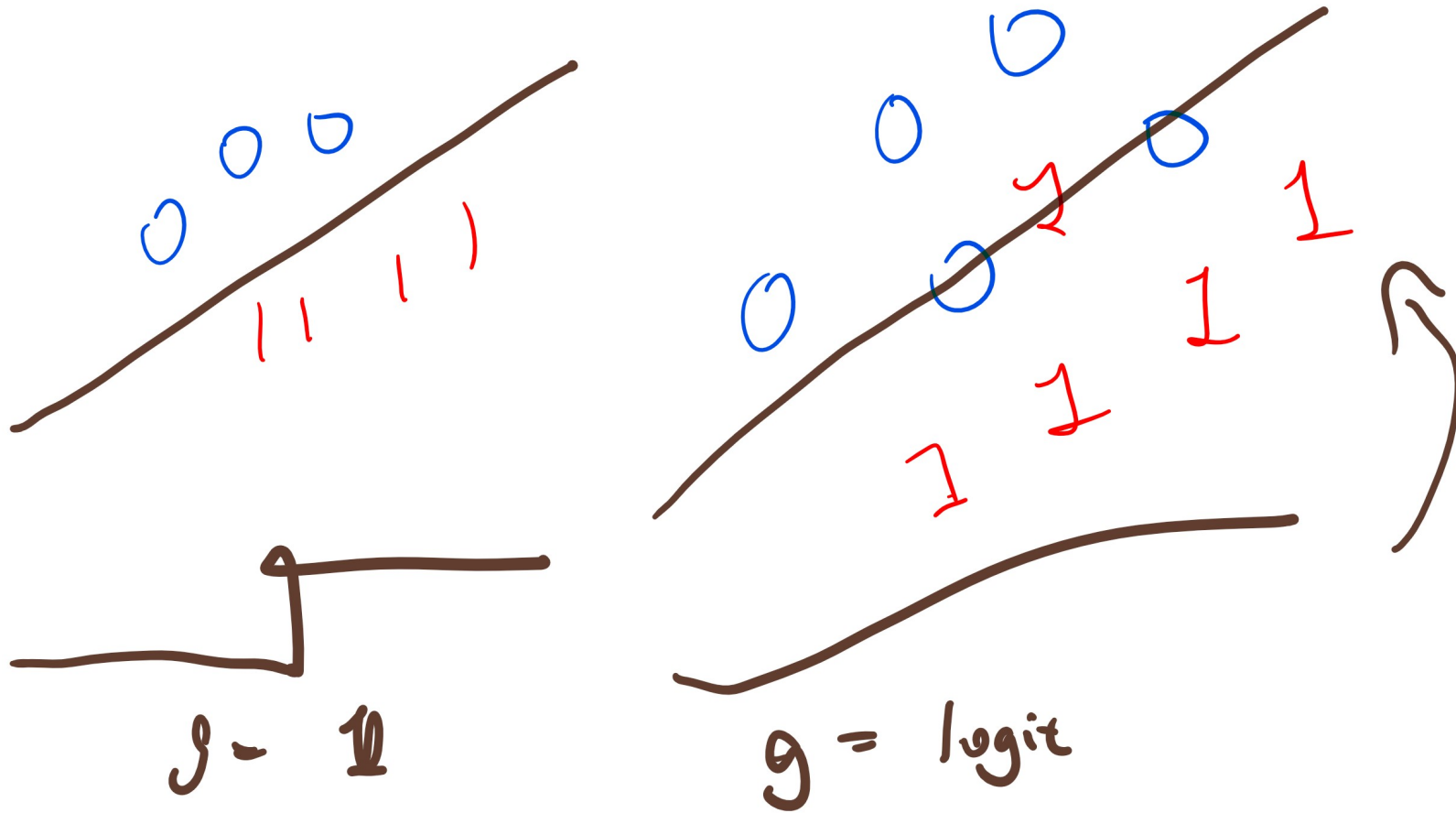


$\boxed{\cdot} \rightarrow w^T x + b$



# Sigmoid Function Diagram

Motivation



# Cross Entropy Loss Function

## Motivation

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

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

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

# Logistic Regression

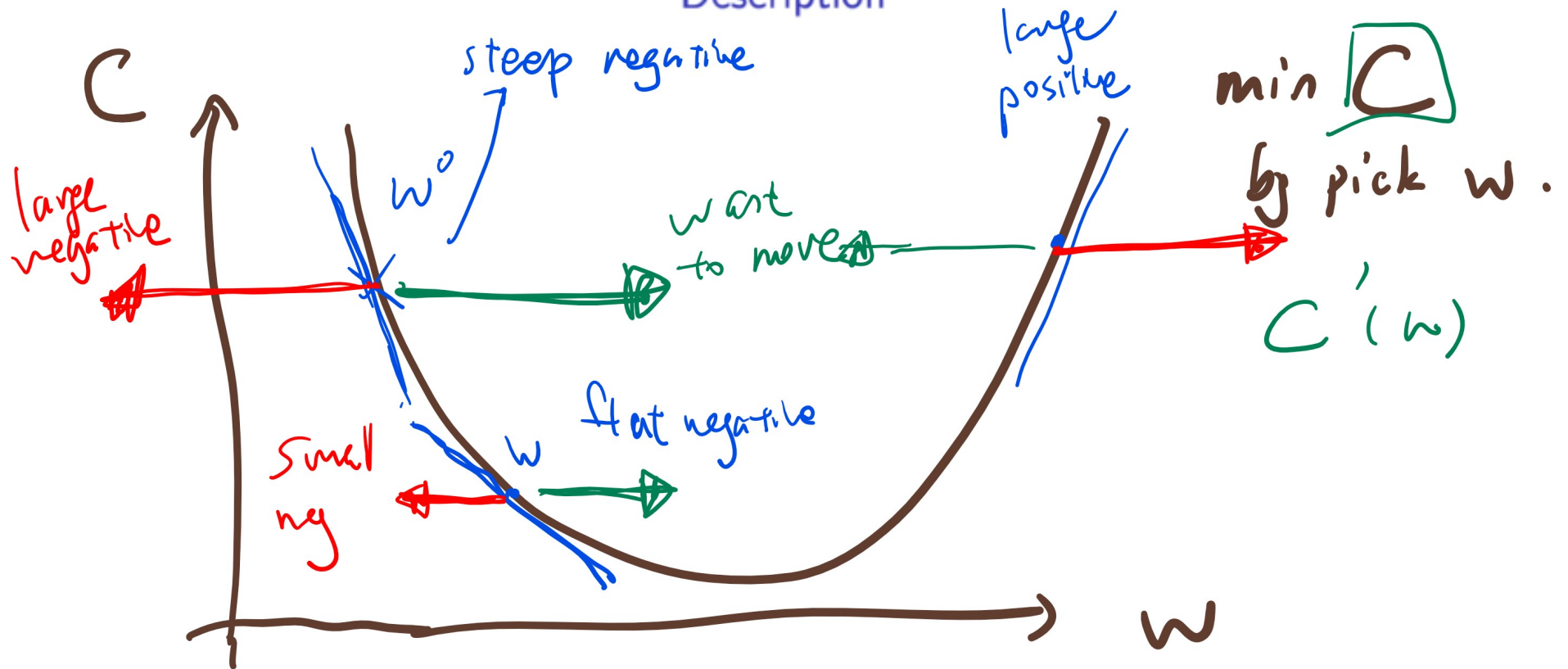
## Description

- 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.

last line  $\rightarrow + \alpha (a_i - y_i) x_i$

# Optimization Diagram

Description



# Gradient Descent Intuition

## Definition

- If a small increase in  $w_1$  causes the distances from the points to the regression line to decrease: increase  $w_1$ .
- If a small increase in  $w_1$  causes the distances from the points to the regression line to increase: decrease  $w_1$ .
- The change in distance due to change in  $w_1$  is the derivative.
- The change in distance due to change in  $\begin{bmatrix} w \\ b \end{bmatrix}$  is the gradient.

# Gradient

## Definition

- The gradient is the vector of derivatives.

- The gradient of

$f(x_i) = w^T x_i + b = w_1 x_{i1} + w_2 x_{i2} + \dots + w_m x_{im} + b$  is:

$\nabla_w f$  = 
$$\begin{bmatrix} \frac{\partial f}{\partial w_1} \\ \frac{\partial f}{\partial w_2} \\ \dots \\ \frac{\partial f}{\partial w_m} \end{bmatrix} = \begin{bmatrix} x_{i1} \\ x_{i2} \\ \dots \\ x_{im} \end{bmatrix} = x_i$$

*partial derivative*

$\nabla_b f = 1$

$\nabla_w w^T x_i = x_i$

# Chain Rule

## Definition

- The gradient of  $f(x_i) = g(w^T x_i + b) = g(w_1 x_{i1} + w_2 x_{i2} + \dots + w_m x_{im} + b)$  can be found using the chain rule.

$$\nabla_w f = g'(w^T x_i + b) x_i$$

$$\nabla_b f = g'(w^T x_i + b) \cdot 1$$

$$\nabla_w (w^T x_i + b)$$

- In particular, for the logistic function  $g$  :

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

$$g'(\square) = g(\square) (1 - g(\square))$$



# Logistic Gradient Derivation

Definition

$$\begin{aligned}
 g(x) &= \frac{1}{1 + e^{-x}} = (1 + e^{-x})^{-1} \\
 g'(x) &= -1 \cdot (1 + e^{-x})^{-2} \cdot \boxed{e^{-x} \cdot (-1)} \\
 &= \frac{1}{1 + e^{-x}} \cdot \frac{e^{-x}}{1 + e^{-x}} \\
 &= g(x) (1 - g(x))
 \end{aligned}$$

$\frac{d}{dx} (1 + e^{-x})$



# Gradient Descent Step

## Definition

- For logistic regression, use chain rule twice.

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

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

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

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

# Gradient Descent Derivation

## Definition

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

Chain Rule

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

from prev slide

$$C = - \sum_i y_i \log a_i + (1-y_i) \log(1-a_i)$$

$$\frac{dC}{da_i} = - \frac{y_i}{a_i} - \frac{1-y_i}{1-a_i} (-1) = - \frac{y_i}{a_i} + \frac{1-y_i}{1-a_i}$$

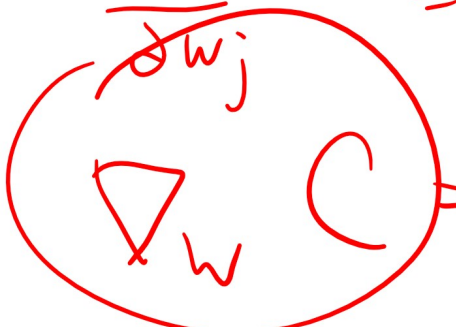
$$\frac{\partial C}{\partial w_j} = \left( - \frac{y_i}{a_i} + \frac{1-y_i}{1-a_i} \right) a_i (1-a_i) x_{ij}$$

# Learning Rate Diagram

Definition

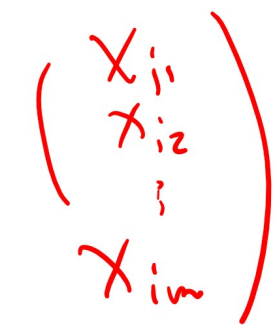
$$\left( - (1 - a_i) y_i + a_i (1 - y_i) \right) X_{ij}$$

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



Gradient vector

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



$$w = w - \alpha \nabla_w C$$

3

# Gradient Descent

Quiz (Graded)

ignore ①

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

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

$$\nabla_w a_i = a_i(1 - a_i)x_i$$

- A:  $w = w - \alpha \sum (a_i - y_i) x_i$
- B:  $w = w - \alpha \sum (a_i - y_i) a_i x_i$
- C:  $w = w - \alpha \sum (a_i - y_i) (1 - a_i) x_i$
- ✓ • D:  $w = w - \alpha \sum (a_i - y_i) a_i (1 - a_i) x_i$
- E: None of the above

$$\begin{aligned} \nabla_w C &= \frac{\partial C}{\partial a_i} \nabla_w a_i \\ &= (a_i - y_i) \end{aligned}$$

$$|x| = \begin{cases} x & \text{if } x > 0 \\ -x & \text{if } x < 0 \end{cases}$$

# Gradient Descent, Another One

Quiz (Graded)



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

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

$$C = \sum_{i=1}^n |a_i - y_i|$$

$\nabla_w a_i$  for logit  
 $= \sum_{i=1}^n a_i (1 - a_i) x_i$

- A:  $w = w - \alpha \sum (a_i - y_i) a_i (1 - a_i) x_i$
- B:  $w = w - \alpha \sum |a_i - y_i| a_i (1 - a_i) x_i$
- C:  $w = w - \alpha \sum \mathbb{1}_{\{a_i - y_i > 0\}} a_i (1 - a_i) x_i$
- D:  $w = w - \alpha \sum \text{sign}(a_i - y_i) a_i (1 - a_i) x_i$
- E: None of the above

$\begin{cases} 1 & a_i - y_i > 0 \\ 0 & \text{not} \end{cases}$

$\begin{cases} +1 & a_i - y_i > 0 \\ 0 & a_i = y_i \\ -1 & a_i - y_i < 0 \end{cases}$



# Logistic Regression, Part 1

## Algorithm

- Inputs: instances:  $\{x_i\}_{i=1}^n$  and  $\{y_i\}_{i=1}^n$
- Outputs: weights and biases:  $w_1, w_2, \dots, w_m$  and  $b$
- Initialize the weights.

$w_1, \dots, w_m, b \sim \text{Unif} [0, 1]$

$w = (X^T X)^{-1} X^T y$

- Evaluate the activation function.

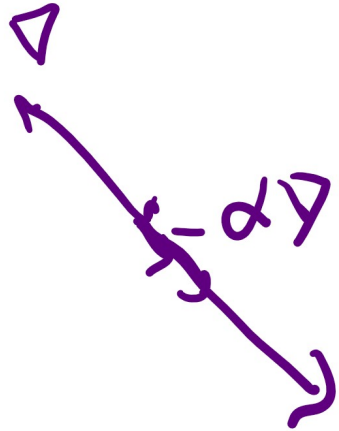
$g(x) = \frac{1}{1 + \exp(-x)}$

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

# Logistic Regression, Part 2

## Algorithm

- Update the weights and bias using gradient descent.



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

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

- Repeat the process until convergent.

$$|C - C^{\text{prev}}| < \epsilon \quad \approx 0.0001$$

# Other Non-linear Activation Function

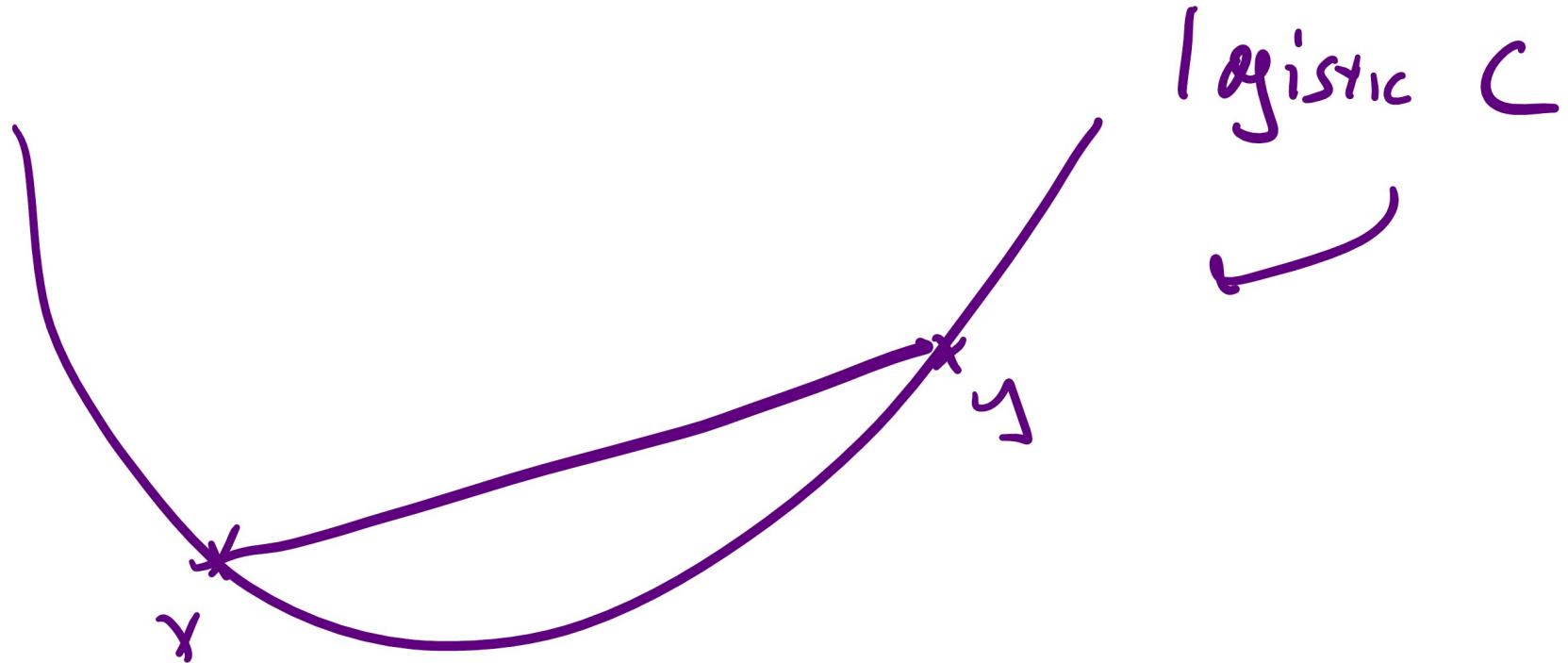
## Discussion

- 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. Gradient descent can be used.



# Convexity Diagram

## Discussion



# Convexity

## Discussion

- If a function is convex, gradient descent with any initialization will converge to the global minimum.
- If a function is not convex, gradient descent with different initializations may converge to different local minima.
- A twice differentiable function is convex if and only if its second derivative is non-negative.
- In the multivariate case, it means the Hessian matrix is positive semidefinite.

# Positive Semidefinite

## Discussion

- Hessian matrix is the matrix of second derivatives:

$$H : H_{ij} = \frac{\partial^2 f}{\partial x_i \partial x_j}$$

$$H = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} \\ \frac{\partial^2 f}{\partial x_1 \partial x_2} & \frac{\partial^2 f}{\partial x_2^2} \end{bmatrix}$$

$f(x_1, x_2)$

- A matrix  $H$  is positive semidefinite if  $x^T H x \geq 0 \forall x \in \mathbb{R}^n$ .

- A symmetric matrix is positive semidefinite if and only if all of its eigenvalues are non-negative.

$$Hx = \lambda x$$

↑  
eigenvalue

↖ eigen vector

# Convex Functions

## Quiz (Participation)

- What is the Hessian (second derivative) of

$$f(x) = x_1^2 + 4x_1x_2 + x_2^2 = [x_1 \ x_2] \begin{bmatrix} 1 & 2 \\ 2 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

$2x_1 + 4x_2$

2

$2x_1 + 4x_2$

4

$$\frac{\partial^2 f}{\partial x_1^2}$$

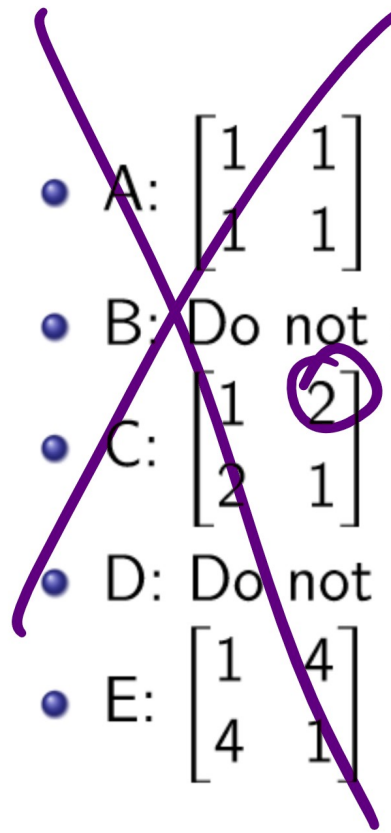
$$\frac{\partial^2 f}{\partial x_1 \partial x_2}$$

$$\frac{\partial^2 f}{\partial x_2^2}$$



2 4  
4 2

- A:  $\begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$
- B: Do not choose this.
- C:  $\begin{bmatrix} 1 & 2 \\ 2 & 1 \end{bmatrix}$
- D: Do not choose this.
- E:  $\begin{bmatrix} 1 & 4 \\ 4 & 1 \end{bmatrix}$



# Definiteness

Quiz (Participation) put on midterm.

$$\lambda = -1$$

• Which ones (two) of the following are the eigenvalues of  $\begin{bmatrix} 1 & 2 \\ 2 & 1 \end{bmatrix}$ ? Two eigenvectors are  $\begin{bmatrix} 1 \\ 1 \end{bmatrix}$  and  $\begin{bmatrix} -1 \\ 1 \end{bmatrix}$ .

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

$$Hx = \lambda x$$

$$\begin{matrix} \uparrow \\ \text{e.v.} \end{matrix} \begin{bmatrix} 1 & 2 \\ 2 & 1 \end{bmatrix} \begin{bmatrix} -1 \\ 1 \end{bmatrix} = \begin{bmatrix} 1 \\ -1 \end{bmatrix}$$

$$\begin{bmatrix} ? \\ ? \end{bmatrix} = \begin{bmatrix} 1 & 2 \\ 2 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} = ? \begin{bmatrix} ? \\ ? \end{bmatrix}$$

# Image as Input

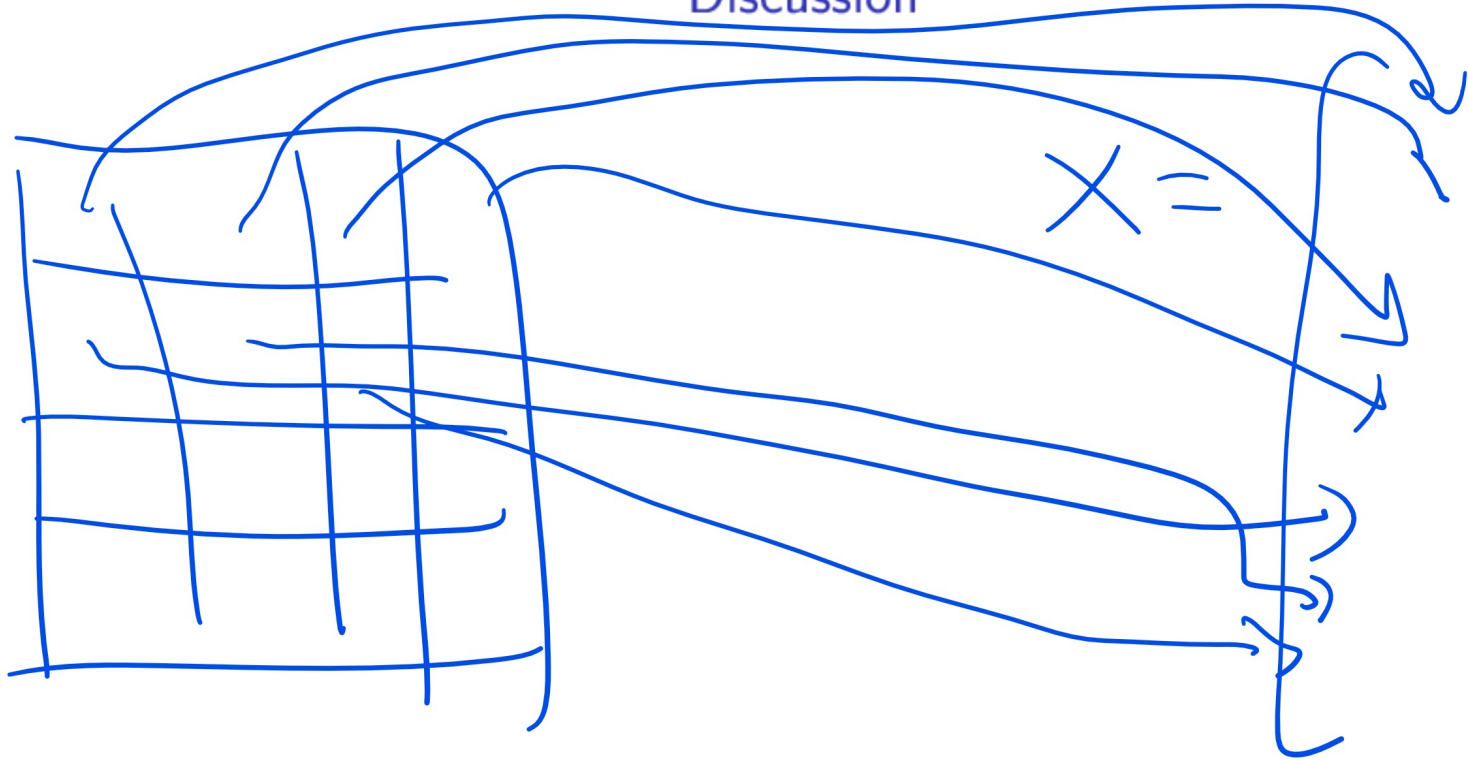
## Discussion

- Simplest feature vector for an image is the flattened pixel intensities.
- One way to compute pixel intensity is to use the average of the RGB values divided by 255.
- Pixel intensity of each pixel is between 0 and 1.
- An  $n_w$  pixel by  $n_h$  pixel image then can be flattened into a  $m = n_w n_h$  dimensional input feature vector  $x$ .



# Flattened Feature Vector Diagram

Discussion



$$g(w^T x)$$

# AND Operator Data

## Quiz (Participation)

- Sample data for AND

$x_1$	$x_2$	$y$
0	0	0
0	1	0
1	0	0
1	1	1



# Learning AND Operator

## Quiz (Participation)

- Which one of the following is AND?
- A:  $\hat{y} = \mathbb{1}_{\{1x_1+1x_2-1.5 \geq 0\}}$
- B:  $\hat{y} = \mathbb{1}_{\{1x_1+1x_2-0.5 \geq 0\}}$
- C:  $\hat{y} = \mathbb{1}_{\{-1x_1+0.5 \geq 0\}}$
- D:  $\hat{y} = \mathbb{1}_{\{-1x_1-1x_2+0.5 \geq 0\}}$
- E: None of the above

# OR Operator Data

## Quiz (Graded)

- Sample data for OR

$x_1$	$x_2$	$y$
0	0	0
0	1	1
1	0	1
1	1	1

# Learning OR Operator

## Quiz (Graded)

- Which one of the following is OR?
- A:  $\hat{y} = \mathbb{1}_{\{1x_1+1x_2-1.5 \geq 0\}}$
- B:  $\hat{y} = \mathbb{1}_{\{1x_1+1x_2-0.5 \geq 0\}}$
- C:  $\hat{y} = \mathbb{1}_{\{-1x_1+0.5 \geq 0\}}$
- D:  $\hat{y} = \mathbb{1}_{\{-1x_1-1x_2+0.5 \geq 0\}}$
- E: None of the above

# XOR Data

## Quiz (Graded)

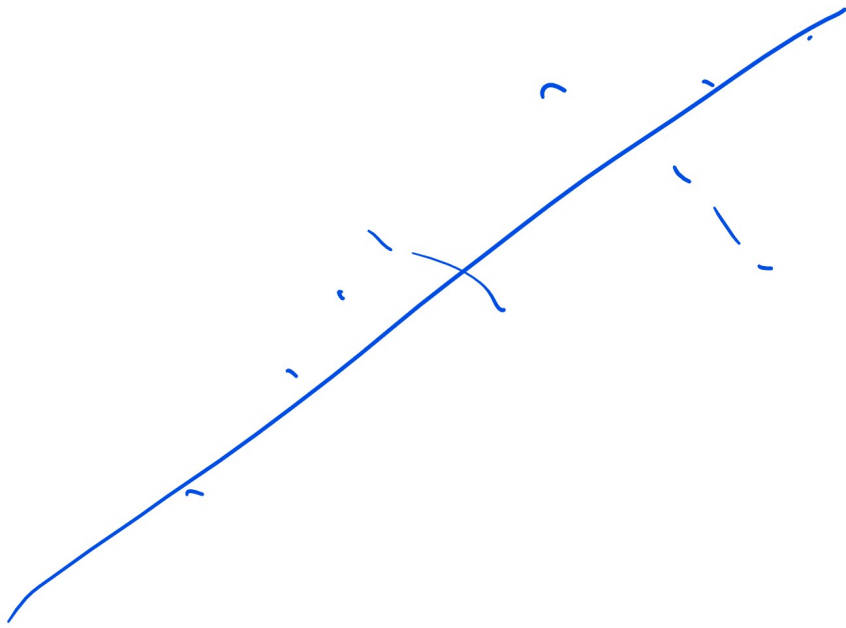
- Sample data for XOR

$x_1$	$x_2$	$y$
0	0	0
0	1	1
1	0	1
1	1	0

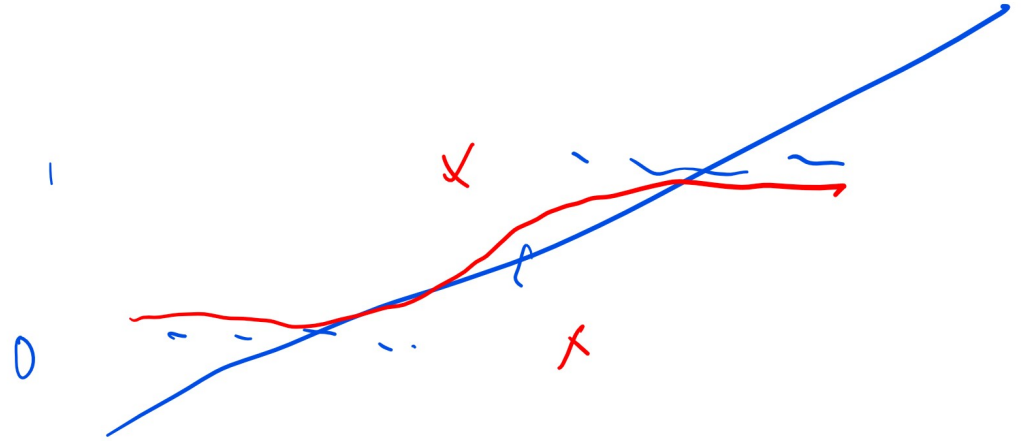
# Learning XOR Operator

## Quiz (Graded)

- Which one of the following is XOR?
- A:  $\hat{y} = \mathbb{1}_{\{1x_1 + 1x_2 - 1.5 \geq 0\}}$
- B:  $\hat{y} = \mathbb{1}_{\{1x_1 + 1x_2 - 0.5 \geq 0\}}$
- C:  $\hat{y} = \mathbb{1}_{\{-1x_1 + 0.5 \geq 0\}}$
- D:  $\hat{y} = \mathbb{1}_{\{-1x_1 - 1x_2 + 0.5 \geq 0\}}$
- E: None of the above



$$y = wx + b$$



$$\hat{y} = \begin{cases} 1 & \text{if } y = \underline{w}^T x \geq 0.5 \\ 0 & \text{if } y < 0.5 \end{cases}$$

