

# CSE 4820/5819

# Introduction to

# Logistic Regression

Suining He

Department of Computer Science and Engineering

University of Connecticut

[suining.he@uconn.edu](mailto:suining.he@uconn.edu)



# Outline

- Background
- Gradient Descent
- Comparison with Linear Regression



# Classification

- **Data:** A set of data records (also called examples, instances or cases) described by
  - Each example is labelled with a pre-defined class.
- **Goal:** To learn a classification model from the data that can be used to predict the classes of new (future, or test) cases/instances.



# Example: Loan Application Binary Classification

ID	Age	Has_Job	Own_House	Credit_Rating	Class
1	young	false	false	fair	No
2	young	false	false	good	No
3	young	true	false	good	Yes
4	young	true	true	fair	Yes
5	young	false	false	fair	No
6	middle	false	false	fair	No
7	middle	false	false	good	No
8	middle	true	true	good	Yes
9	middle	false	true	excellent	Yes
10	middle	false	true	excellent	Yes
11	old	false	true	excellent	Yes
12	old	false	true	good	Yes
13	old	true	false	good	Yes
14	old	true	false	excellent	Yes

old	false	false	fair	?
-----	-------	-------	------	---



# Why Use Logistic Regression

- Why?
  - Estimation by maximum likelihood
  - Interpreting coefficients
  - Hypothesis testing
  - Evaluating the performance of the model
- There are many important research topics for which the dependent variable is "limited."
- For example: voting, morbidity or mortality, and participation data is not continuous or distributed normally.



# Background

- Binary **logistic** regression is a type of regression analysis where the dependent variable is binary
  - 0/1, True/False, Yes/No
  - For example: loan application.
    - Coded No (not approved) or Yes (approved)
  - $\underbrace{y \in \{0, 1\}}$

$P$  → probability a data pt.  
belongs to 1



# Logistic Regression Model: An Overview

- The logistic regression model solves these problems:

$$\ln[p/(1-p)] = \mathbf{x}^T \mathbf{w} + b$$

*(logarithm  
natural base  $\ln \frac{P}{1-P}$ )*

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{bmatrix} \quad \mathbf{w} = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_d \end{bmatrix}$$

- $p$  is the probability that an event occurs
  - $p/(1 - p)$  is the **odds ratio**
  - $\ln[p/(1 - p)]$  is the **log odds ratio, or logit**



# Logistic Regression Model: An Overview

- The logistic distribution constrains the estimated probabilities to lie between 0 and 1.
- The estimated probability (using  $\hat{y}$  in the place of  $p$ ) is

$$P = \hat{y} = \frac{1}{1 + \exp(-z)}$$

- where  $z = \mathbf{x}^T \mathbf{w} + b$

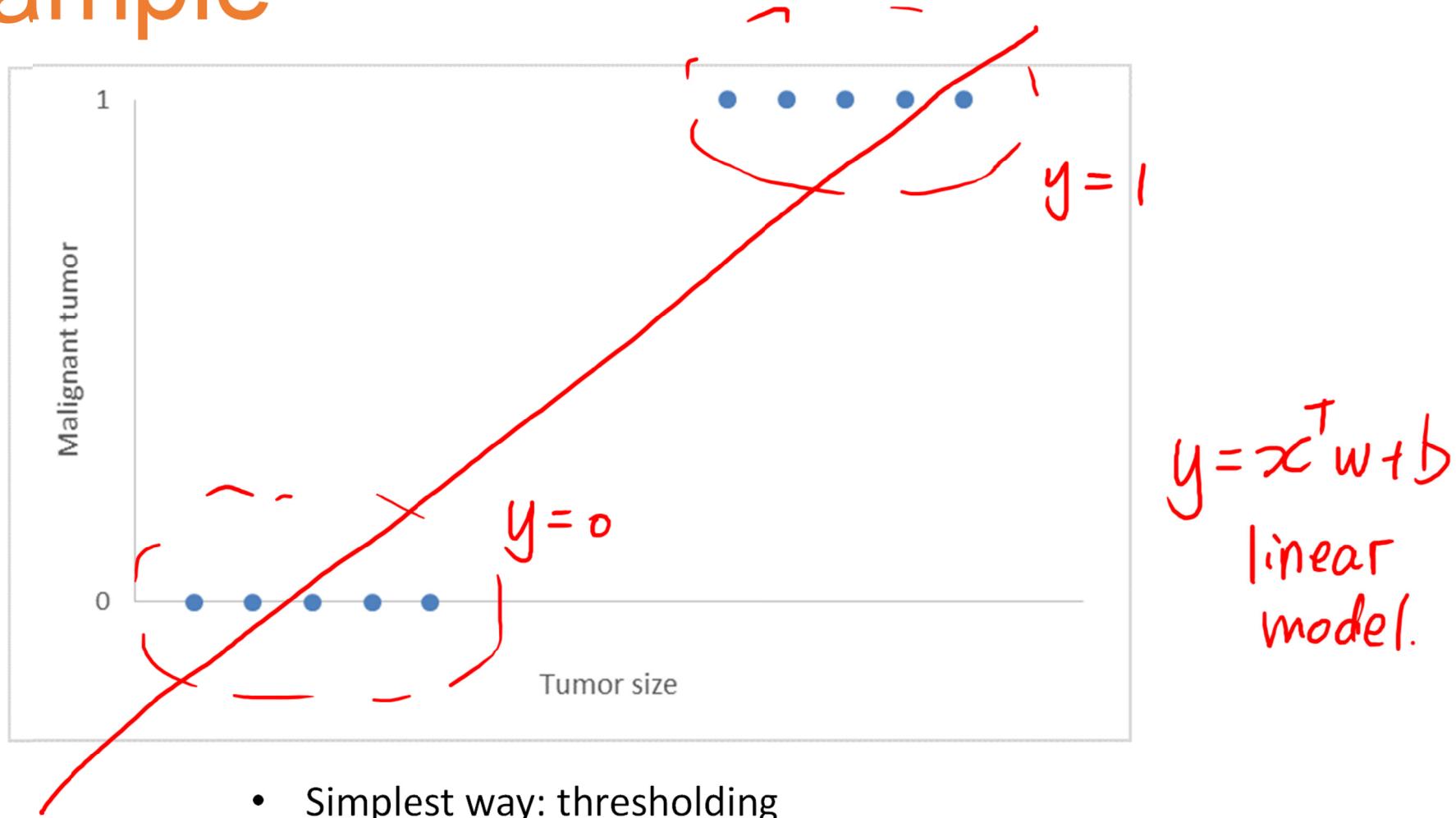
- If you let  $b + \mathbf{x}^T \mathbf{w} = 0$ , then  $p = .50$

- As  $b + \mathbf{x}^T \mathbf{w}$  gets really big, p approaches 1
- As  $b + \mathbf{x}^T \mathbf{w}$  gets really small, p approaches 0

$$\begin{aligned} z = 0 & \quad P = \frac{1}{1 + \exp(0)} \\ z \rightarrow +\infty & \quad P \rightarrow \frac{1}{1 + 0} = \frac{1}{2} = 0.5 \\ z \rightarrow -\infty & \quad P \rightarrow 0 \\ P \rightarrow 0 & \end{aligned}$$

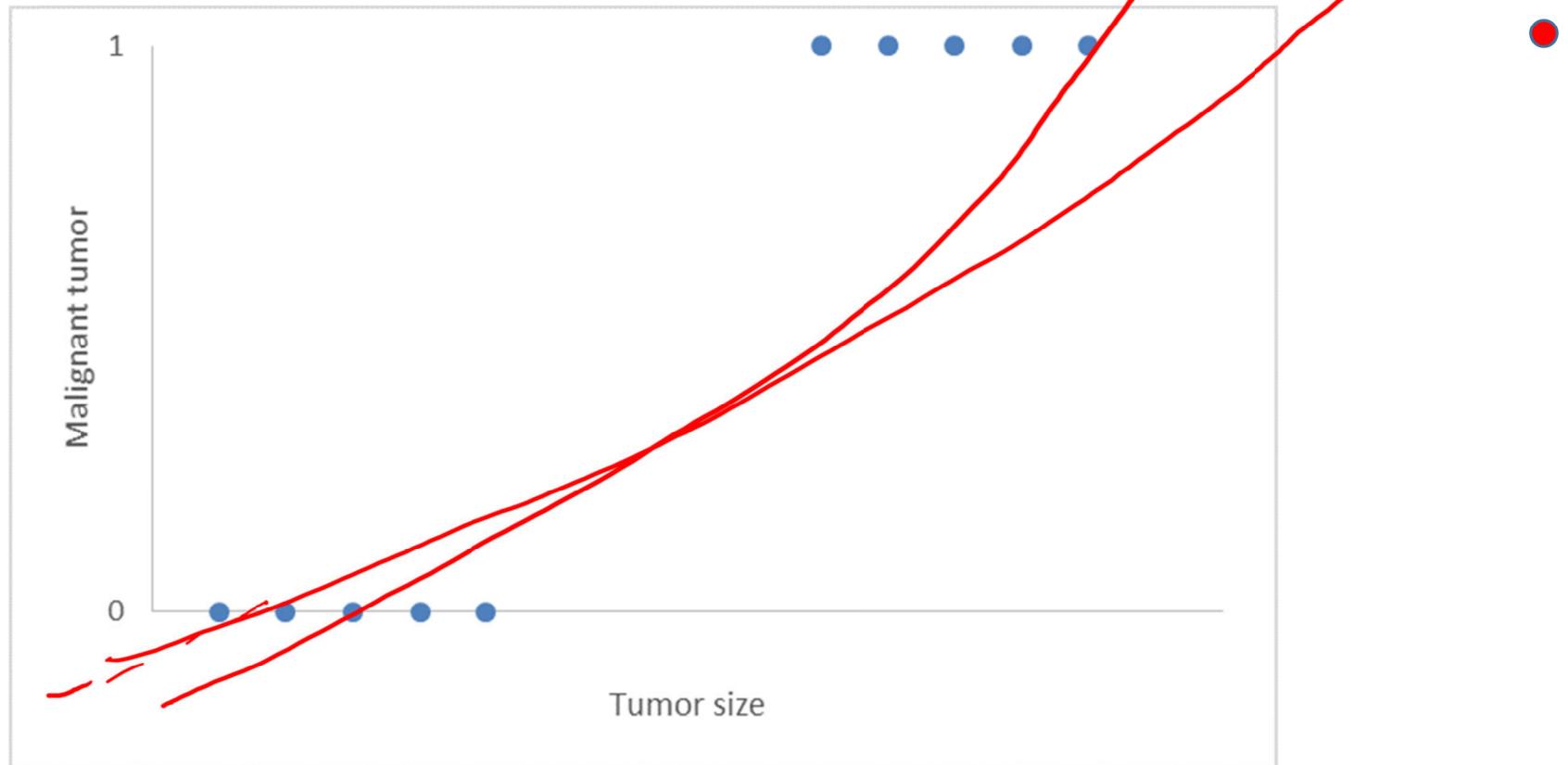


# Example





# Example



- What about the linear model learned in linear regression?



# Origin and Connection with Linear Regression

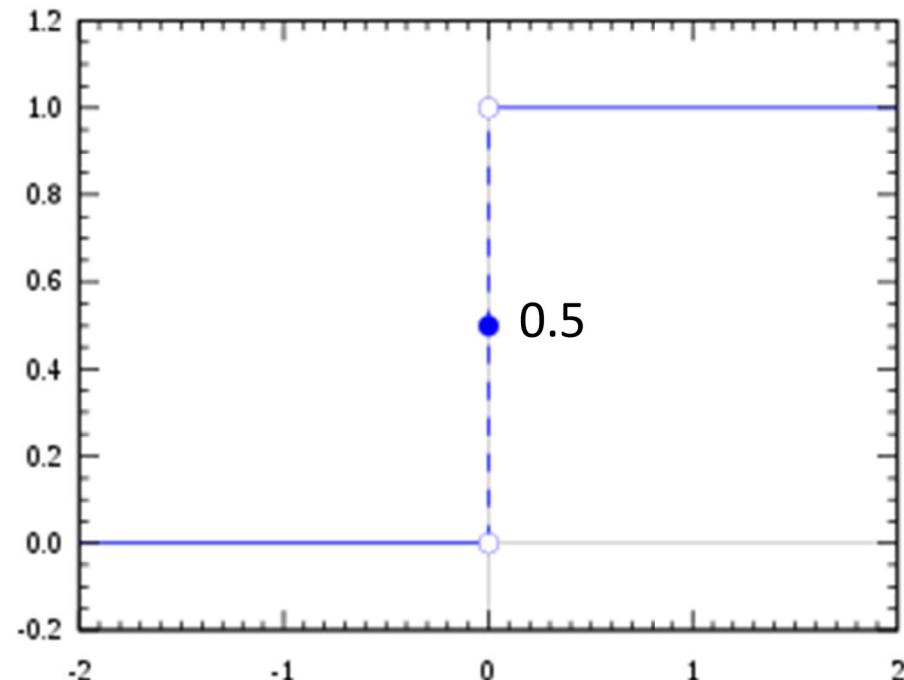
- Recall  $y = \mathbf{x}^T \mathbf{w} + b$  returns  $y \in R$ 
  - It can be any real value
  - We need to turn it into the values of 0 and 1 for the *logistic* use
- One simple way:

- Unit-step function

$$z = \mathbf{x}^T \mathbf{w} + b$$

$$\hat{y} = \begin{cases} 0, & \text{if } z < 0 \\ 0.5, & \text{if } z = 0 \\ 1, & \text{if } z > 0 \end{cases}$$

probability P





2024/9/19

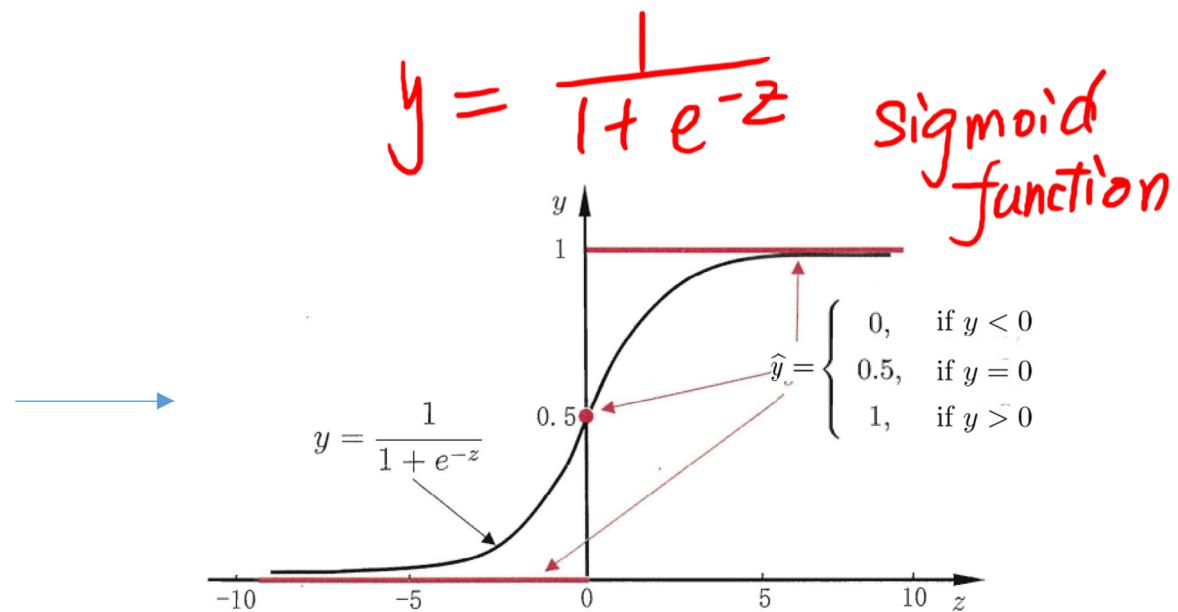
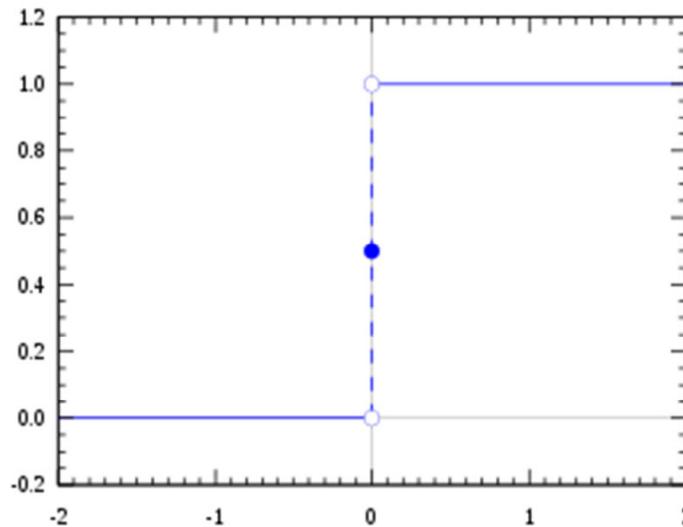
CSE 4820/5819  
Introduction to Machine Learning  
University of Connecticut

12



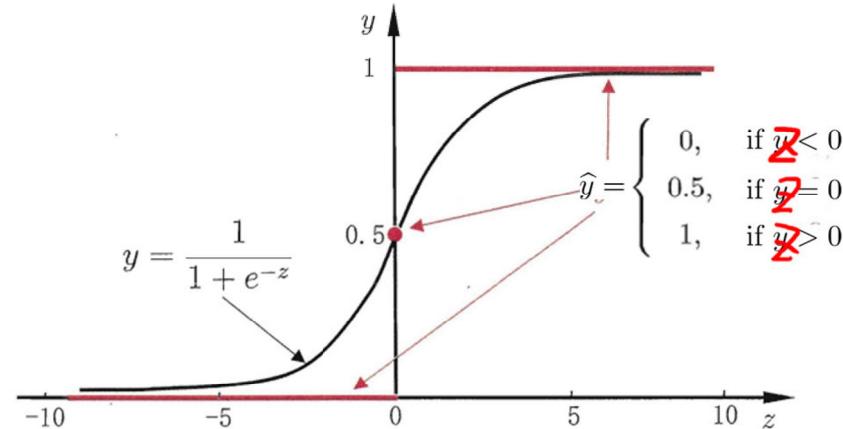
# Normally, we consider Unit Step Function

- Unit function is not **continuous**
  - We need a *surrogate* function
  - A surrogate function is a function to approximate another





# Surrogate Function



- Logistic function

$$\hat{y} = \frac{1}{1 + \exp(-z)}$$

- With  $z = \mathbf{x}^T \mathbf{w} + b$ , then we can have

probability  $\Leftarrow \hat{y} = \frac{1}{1 + \exp(-(x^T w + b))}$



# Logistic Regression

- Find a function  $f$  that best predicts the conditional probability directly?

$p(y_i = \text{"yes" or "no"} | \mathbf{x}_i)$  (binary classification)

- Still find a linear function of  $\mathbf{x}_i$  parameterized with  $\mathbf{w}$ , i.e.,  $\mathbf{x}_i^T \mathbf{w} + b$
- But  $\mathbf{x}_i^T \mathbf{w} + b$  is used to compute the probability ( $p$ )
- $P \in [0,1]$ , however,  $\mathbf{x}_i^T \mathbf{w} + b \in (-\infty, +\infty)$

*So, how to link  $\mathbf{x}_i^T \mathbf{w} + b$  to  $p$ ?*



# Further Transformation

probability.

- We can further convert  $\hat{y} = \frac{1}{1+\exp(-(\mathbf{x}^T \mathbf{w} + b))}$  into

$$\ln \frac{\hat{y}}{1-\hat{y}} = \mathbf{x}^T \mathbf{w} + b$$

*(linear model).*

- Recall that we are expecting  $y$  has value of either 1 or 0 (e.g.,  $y = 1$  for class 1 and  $y = 0$  for class 2)
- The **odds** is defined as  $\frac{\hat{y}}{1-\hat{y}}$ 
  - Log odds (logit):  $\ln \frac{\hat{y}}{1-\hat{y}}$



⇒ output prob of data.

linear model.

$$\hat{y} = \frac{1}{1 + \exp(-(x^T w + b))} \longrightarrow \ln \frac{\hat{y}}{1 - \hat{y}} = x^T w + b$$

binary classification

toss a coin

P → probability head. ( $\hat{y}$ )

1-P → prob tail. ( $1 - \hat{y}$ ).

logit  $\ln \frac{\hat{y}}{1 - \hat{y}}$   $1 - \hat{y} = \frac{1 + \exp(-(x^T w + b))}{1 + \exp(-(x^T w + b))} - \frac{1}{1 + \exp(-(x^T w + b))}$

$$\frac{\hat{y}}{1 - \hat{y}} = \frac{\frac{1}{1 + \exp(-(x^T w + b))}}{\frac{\exp(-(x^T w + b))}{1 + \exp(-(x^T w + b))}} = \frac{1}{1 + \exp(-(x^T w + b))} = \exp(x^T w + b)$$

$$\ln \frac{\hat{y}}{1 - \hat{y}} = \ln(\exp(x^T w + b)) = x^T w + b.$$

CSE 4820/5819



$$\vec{x} = \begin{bmatrix} x_w \\ x_h \end{bmatrix}$$

$$\vec{w} = \begin{bmatrix} w_1 \\ w_2 \end{bmatrix}$$

weight  
height

# Quick Example for Settings



$y$

0 Not Evolved

1 Evolved

$P \rightarrow$  probability that pokemon evolved ( $\text{prob}(y=1)$ )

9/19/2024

$$\ln \frac{\hat{y}}{1-\hat{y}} = \vec{x}^T \vec{w} + b \Leftrightarrow \hat{y} = \frac{1}{1+\exp(-(\vec{x}^T \vec{w} + b))}$$

1st Pika.  $\vec{x}_1 = \begin{bmatrix} 4.79 \\ 0.38 \end{bmatrix}$

$$\hat{y}_1 = \frac{1}{1+\exp(-(\vec{x}_1^T \vec{w} + b))}$$

$$= \frac{1}{1+\exp\left(-\left(\begin{bmatrix} 4.79 \\ 0.38 \end{bmatrix}^T \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} + b\right)\right)}$$

$\Rightarrow 0 \equiv$  ground truth (not evolved)

$$\hat{y}_2 = \frac{1}{1+\exp\left(-\left(\begin{bmatrix} 31.42 \\ 0.82 \end{bmatrix}^T \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} + b\right)\right)}$$

$\Rightarrow 1 \equiv$  ground truth (evolved)



# Connection Between Linear Regression and Logistic Regression

$$\text{Odds}(P) = \frac{P(y_i=1|\mathbf{x}_i)}{P(y_i=0|\mathbf{x}_i)} = \frac{P(y_i=1|\mathbf{x}_i)}{1-P(y_i=1|\mathbf{x}_i)} = \frac{P}{1-P}$$

$P \rightarrow$  probab belongs to 1

$1-P \rightarrow$  prob belongs to 0



# Connection Between Linear Regression and Logistic Regression

$$\ln \frac{P}{1-P} = \mathbf{x}_i^T \mathbf{w} + b$$

$$\text{Odds}(P) = \frac{P(y_i=1|\mathbf{x}_i)}{P(y_i=0|\mathbf{x}_i)} = \frac{P(y_i=1|\mathbf{x}_i)}{1-P(y_i=1|\mathbf{x}_i)}$$

Probability $P(y_i = 1 \mathbf{x}_i)$	Odds
1.0	$+\infty$
0.99	99
0.75	3
0.5	1
0.25	0.3333
0.01	0.0101
0	0

$$\frac{P}{1-P}$$

Probability:  $P \in [0,1]$

Odd value:  $\text{Odds}(P) \in (0, +\infty)$



$\mathbf{x}^T \mathbf{w} + b \in (-\infty, +\infty)?$

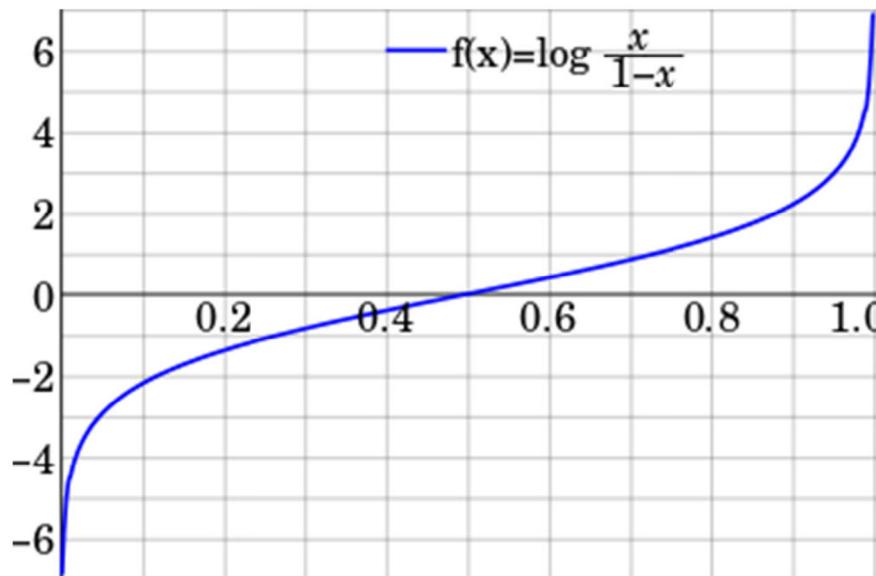


# Logit Function

- Take the logarithm of odds

logit.

$$\ln(Odds) = \ln\left(\frac{P}{1 - P}\right)$$



We call this function the **logit function** of  $P$ , i.e.,

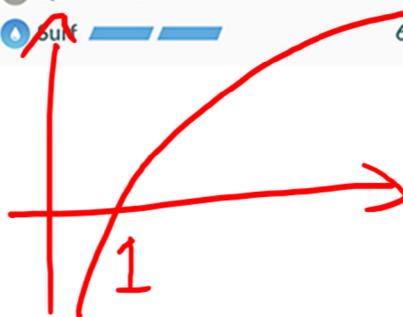
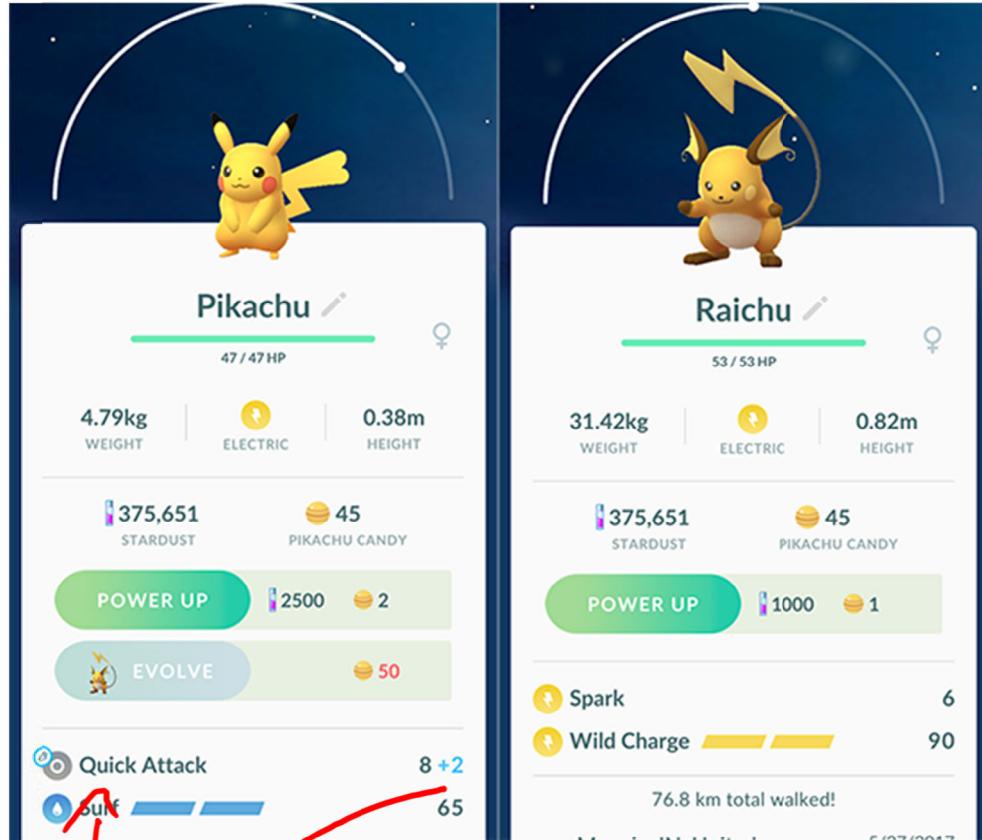
$$\text{logit}(P) = \ln\left(\frac{P}{1 - P}\right)$$

$$\text{logit}(P) \in (-\infty, +\infty)$$

$\underbrace{\mathbf{x}^T \mathbf{w} + b}_{\text{attribute of data pt.}} \in (-\infty, +\infty)$

$$\log\left(\frac{y}{1-y}\right) = \boxed{z^T w + b}$$

# Quick Example (Again!)



$$f(x) = \log x$$

②

$$y_2 = 1 \quad \hat{y}_2 \rightarrow 1$$

$$\hat{y}_1 \rightarrow 0$$

output

① Pikachu

$$y_1 = 0.$$

true label.

odds ratio.

$$\frac{\hat{y}_1}{1-\hat{y}_1} = 0$$

logit:

$$\log\left(\frac{\hat{y}_1}{1-\hat{y}_1}\right) \rightarrow -\infty$$

$$\frac{\hat{y}_1}{1-\hat{y}_1} \rightarrow 0$$

$$\log\left(\frac{\hat{y}_1}{1-\hat{y}_1}\right) \rightarrow -\infty$$

sign

$$\log\left(\frac{\hat{y}_2}{1-\hat{y}_2}\right) \rightarrow +\infty$$

sign



# Interpretation of Model

$$f_{\mathbf{w}, b}(x) = \hat{y} = \frac{1}{1 + \exp(-(\mathbf{x}^T \mathbf{w} + b))}$$

- Interpreted as:
  - Estimated probability that  $y = 1$  on the input  $x$
- Tumor example  $x = \text{size of tumor}$ 
  - $\hat{y} = 0.7$
  - Tell patient that 70% chance of tumor being malignant

$$f_{\mathbf{w}, b}(x) = P(\hat{y} = 1 | x; \mathbf{w}, b)$$

- “Probability that  $\hat{y} = 1$ , given  $x$ , parameterized by  $\mathbf{w}, b$ ”
- $P(\hat{y} = 1 | x; \hat{\mathbf{w}}) + P(\hat{y} = 0 | x; \hat{\mathbf{w}}) = 1$
- $P(\hat{y} = 1 | x; \hat{\mathbf{w}}) = 1 - P(\hat{y} = 0 | x; \hat{\mathbf{w}})$

CSE 4820/5819



$$\hat{\mathbf{w}} = [\mathbf{w}; b]$$

$$\hat{\mathbf{w}} = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_d \\ b \end{bmatrix}$$

# Logistic Regression

- Assume log odds is a linear function of  $\mathbf{x}$

logit  $\ln\left(\frac{P(y_i = 1|\mathbf{x}_i)}{1 - P(y_i = 1|\mathbf{x}_i)}\right) = \mathbf{x}_i^T \hat{\mathbf{w}}$  linear model.  
if  $\mathbf{x}_i^T \hat{\mathbf{w}} > 0 \quad P > 0.5$

When  $\mathbf{x}_i^T \hat{\mathbf{w}} > 0, P(y_i = 1|\mathbf{x}_i) > 0.5$

If  $\mathbf{x}_i^T \hat{\mathbf{w}} < 0 \quad P < 0.5$

$$P(y_i = 1|\mathbf{x}_i) = \frac{\exp(\mathbf{x}_i^T \hat{\mathbf{w}})}{1 + \exp(\mathbf{x}_i^T \hat{\mathbf{w}})}$$

- Linear decision** boundary, the line  $\mathbf{x}_i^T \hat{\mathbf{w}} = 0 \quad P = 0.5$
- Decision Rule:** For binary classification, if  $P(y=1|X) > P(y=0|X)$  (or  $P(y=1|X) > 0.5$ ), then  $Y=1$ ; otherwise,  $Y=0$



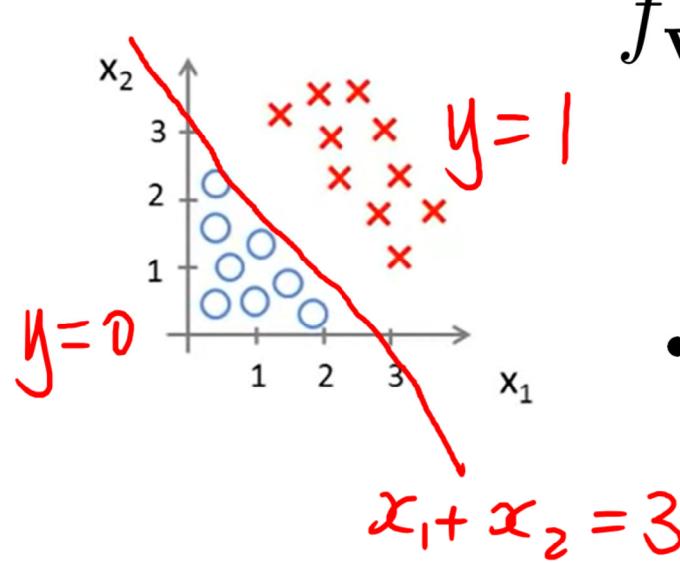
# Quick Deduction

$$\hat{w} = \begin{bmatrix} w_1 \\ \vdots \\ w_d \\ b \end{bmatrix}$$

$$\begin{aligned} p_i &= \frac{1}{1 + \exp(-(x_i^T w + b))} \\ &= \frac{1}{1 + \exp(-x_i^T \hat{w})} \\ &= \frac{1 \cdot \exp(x_i^T \hat{w})}{1 \cdot \exp(x_i^T \hat{w}) + \exp(-x_i^T \hat{w}) \cdot \exp(-x_i^T \hat{w})} \\ &= \frac{\exp(x_i^T \hat{w})}{\exp(x_i^T \hat{w}) + 1} \end{aligned}$$



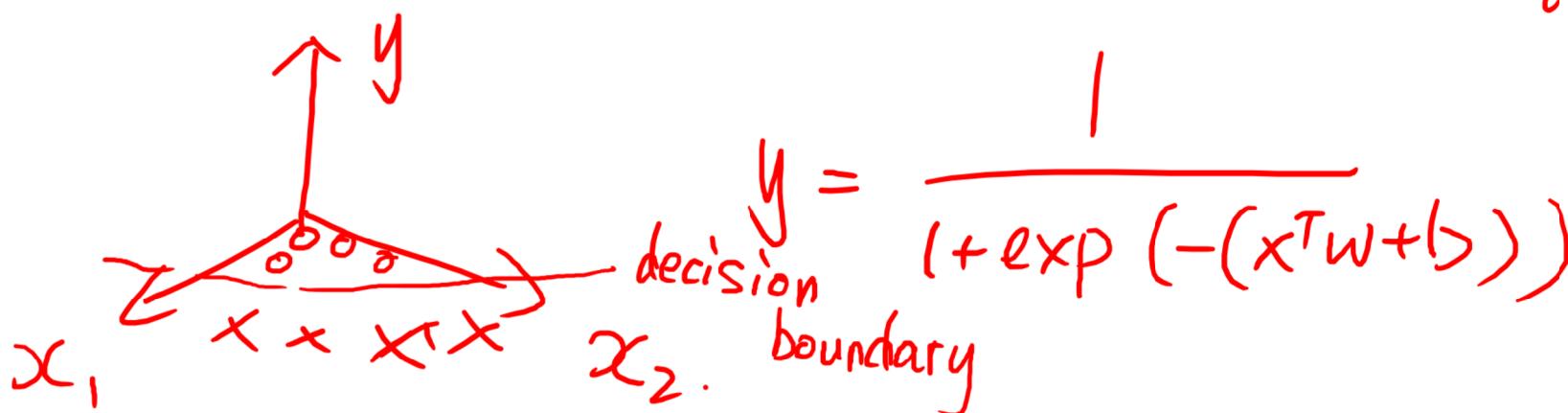
# Example of Using Decision Boundary



$$f_{\mathbf{w}, b}(x) = \hat{y} = \frac{1}{1 + \exp(-(\mathbf{x}^T \mathbf{w} + b))}$$
$$= \frac{1}{1 + \exp(b + w_1 x_1 + w_2 x_2)}$$

- Predict " $\hat{y} = 1$ " if  $-3 + x_1 + x_2 \geq 0$
- Decision boundary  $x_1 + x_2 = 3$

$$\sum_i^T \hat{w} = 0$$





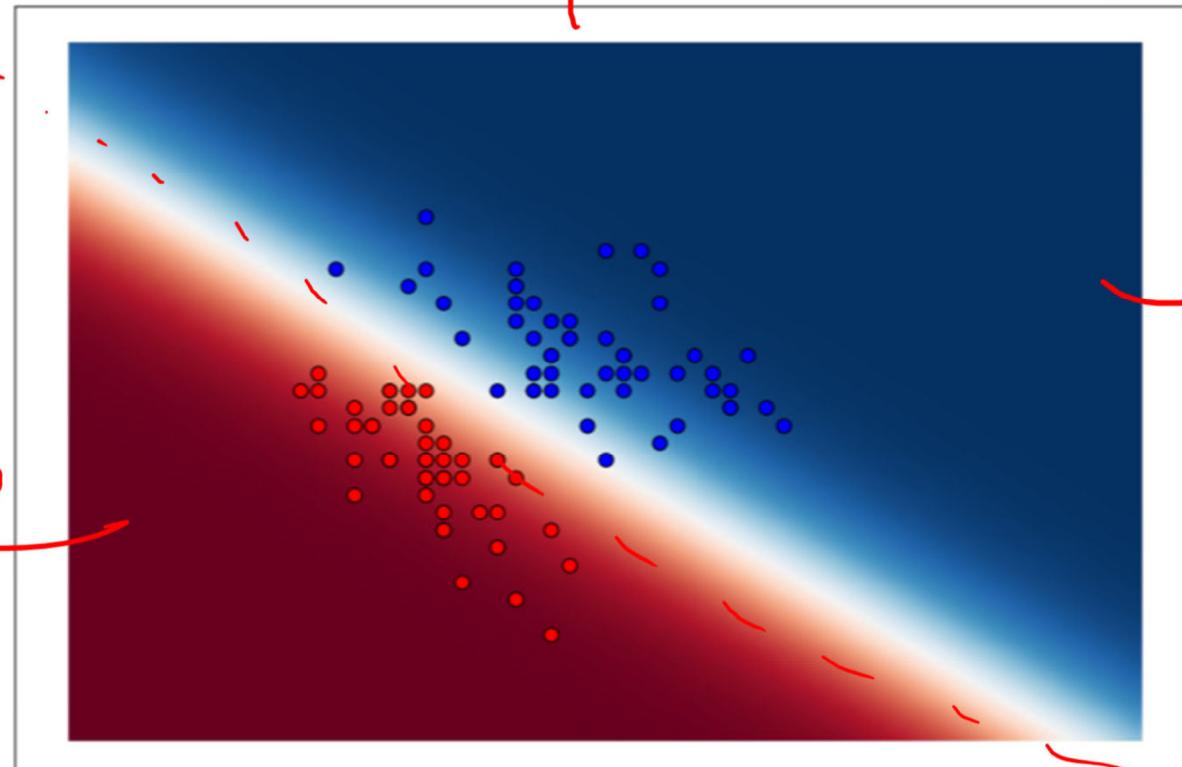
# Visualization

$$p = 0.5$$

$$p(y=1|x) \geq 0$$

$\downarrow$

$$y \rightarrow 0$$



$$p(y=1|x)$$

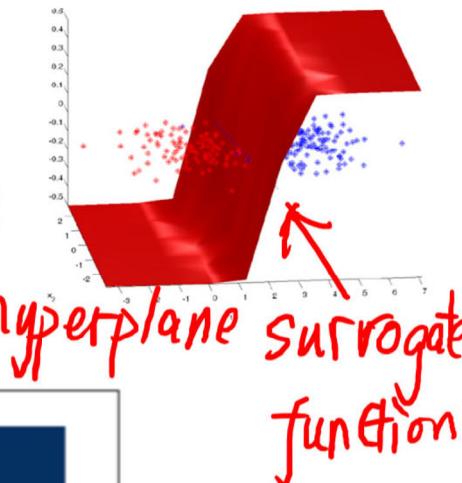
$$CSE 4820/5819$$

Introduction to Machine Learning  
University of Connecticut

$$p = \frac{1}{1 + e^{-(\theta_0 - 3 + x_1 + x_2)}}$$

if  $-3 + x_1 + x_2 = 0$

$$p = 0.5$$



$$p(y=1|x)$$

$\downarrow$

$$y = 1$$

decision  
boundary  
 $-3 + x_1 + x_2 = 0$



# More on Logistic Regression

- The logistic distribution constrains the estimated probabilities to lie between 0 and 1.
- The estimated probability is:

$$\hat{y} = \frac{1}{1+\exp(-(\mathbf{x}^T \mathbf{w} + b))}$$

- If you let  $\mathbf{x}^T \mathbf{w} + b = 0$ , then  $\hat{y} = 0.5$ 
  - As  $\mathbf{x}^T \mathbf{w} + b$  gets really big,  $\hat{y}$  approaches 1
  - As  $\mathbf{x}^T \mathbf{w} + b$  gets really small,  $\hat{y}$  approaches 0



# Understanding with Maximum Likelihood

- When  $y_i = 1$ , find  $\mathbf{w}$  that maximizes

$$\underbrace{P(y_i = 1 | \mathbf{x}_i)}$$

- When  $y_i = 0$ , find  $\mathbf{w}$  that maximizes

$$P(y_i = 0 | \mathbf{x}_i) = 1 - P(y_i = 1 | \mathbf{x}_i)$$

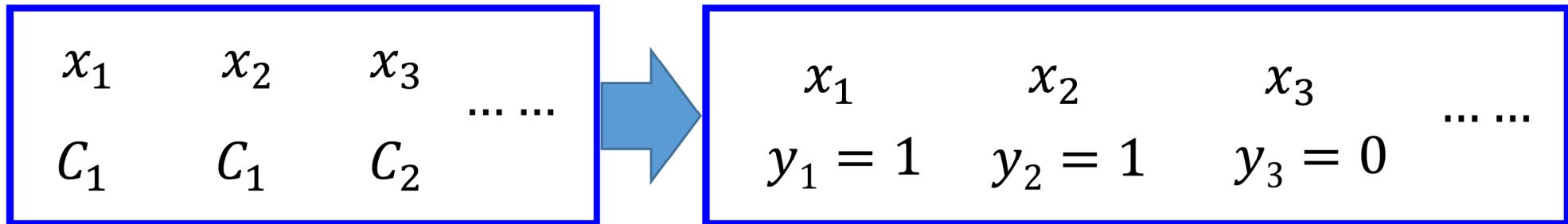
- Overall

$$\max_{\mathbf{w}, b} \left( \prod_{i:y_i=1} P(y_i = 1 | \mathbf{x}_i) \prod_{i:y_i=0} (1 - P(y_i = 1 | \mathbf{x}_i)) \right)$$



# Example

Training Data



Assume the data is generated based on

$$f_{\mathbf{w}, b}(x) = P_{\mathbf{w}, b}(C_1|x) = P(y_i = 1|\mathbf{x}_i)$$

Given a set of  $\mathbf{w}$  and  $b$ , what is the joint probability of generating the data?

$$L(\mathbf{w}, b) = f_{\mathbf{w}, b}(x_1)f_{\mathbf{w}, b}(x_2)\left(1 - f_{\mathbf{w}, b}(x_3)\right)\cdots f_{\mathbf{w}, b}(x_N)$$

The most likely  $\mathbf{w}^*$  and  $b^*$  is the one with the largest  $L(\mathbf{w}, b)$ .

$$\mathbf{w}^*, b^* = \arg \max_{\mathbf{w}, b} L(\mathbf{w}, b)$$

CSE 4820/5819



# y, Maximum Likelihood Example



$$\hat{y}_3 = 0$$

$$\hat{y}_4 = 1$$



2024/9/19

$$y_2 = 1$$

$$\hat{y}_1 \rightarrow 0$$

$$\hat{y}_2 \rightarrow 1$$

$$\hat{y}_3 \rightarrow 0$$

$$\hat{y}_4 \rightarrow 1$$

$$\hat{J}(w, b)$$

$$= \hat{\ell}((1 - f_{w,b}(x_1)))$$

$$\cdot f_{w,b}(x_2)$$

$$\cdot \hat{\ell}((1 - f_{w,b}(x_3)))$$

$$\cdot \hat{\ell}((1 - f_{w,b}(x_4)))$$

$$1 - f_{w,b}(x_1) = \left( f_{w,b}(x_1) \right)^{y_1} (1 - f_{w,b}(x_1))^{1 - y_1}$$

$$y_1 = 0.$$

$$f_{w,b}(x_2) = \left( f_{w,b}(x_2) \right)^{y_2} (1 - f_{w,b}(x_2))^{1 - y_2}$$

$$y_2 = 1$$



# Bernoulli Distribution

- We can write it in a way of **Bernoulli distribution**

$$\max_{\mathbf{w}, b} \left( \prod_i (f_{\mathbf{w}, b}(x))^{y_i} (1 - f_{\mathbf{w}, b}(x))^{(1-y_i)} \right)$$

objective

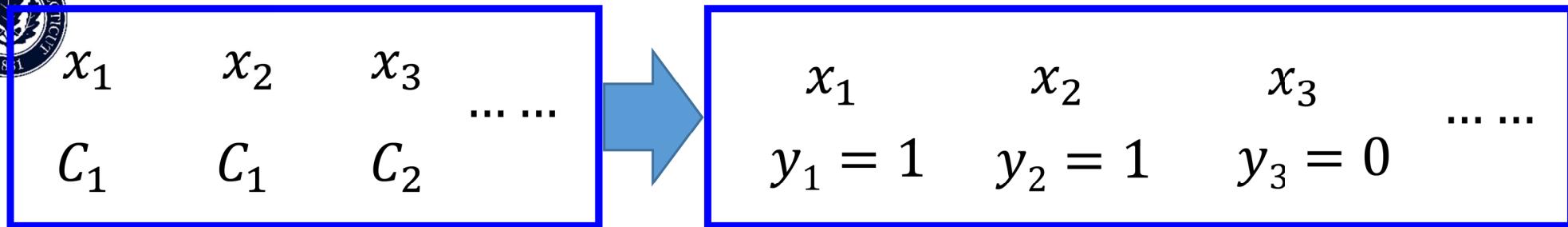
- $f_{\mathbf{w}, b}(x)$  represent the probability that the  $y_i$  is 1 (belong to  $C_1$ )



# Example



2024/9/19



$y_n$ : 1 for class 1, 0 for class 2

$$L(w, b) = f_{w,b}(x_1)f_{w,b}(x_2)\left(1 - f_{w,b}(x_3)\right)\cdots$$

$$\mathbf{w}^*, b^* = \arg \max_{\mathbf{w}, b} L(\mathbf{w}, b)$$

$$\mathbf{w}^*, b^* = \arg \min_{\mathbf{w}, b} -\log(L(\mathbf{w}, b))$$

$$-\log L(\mathbf{w}, b)$$

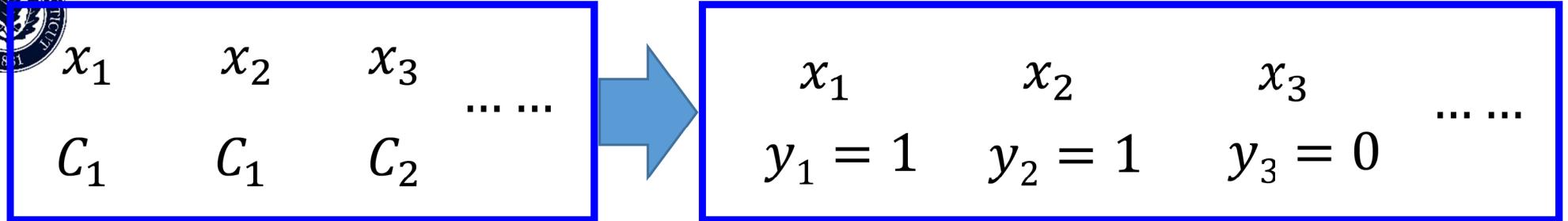
$$f_{w,b}(x_i) = \frac{(f(x_i))^y}{(1-f_{w,b}(x_i))^{1-y}}$$

$$= -\log f_{w,b}(x_1) \rightarrow -[1 \log f(x_1) + (1 - 0y_1) \log(1 - f(x_1))]$$

$$-\log f_{w,b}(x_2) \rightarrow -[1 \log f(x_2) + (1 - 0y_2) \log(1 - f(x_2))]$$

$$-\log(1 - f_{w,b}(x_3)) \rightarrow -[0 \log f(x_3) + (1 - 1y_3) \log(1 - f(x_3))]$$

⋮



$y_n$ : 1 for class 1, 0 for class 2



# Maximum Likelihood Estimation

- Parameter estimation of  $w$  and  $b$

$$f_{w,b}(x) \underbrace{P(y_i = 1 | \mathbf{x}_i)}_{\text{Red underline}} = \frac{\exp(\mathbf{x}_i^T \mathbf{w} + b)}{1 + \exp(\mathbf{x}_i^T \mathbf{w} + b)}, \underbrace{P(y_i = 0 | \mathbf{x}_i)}_{\text{Red underline}} = \frac{1}{1 + \exp(\mathbf{x}_i^T \mathbf{w} + b)}$$

- Maximum likelihood: find the parameters that can maximize the log likelihood

$$L(\mathbf{w}, b) = \sum_{i=1}^m \log p(y_i | \mathbf{x}_i; \mathbf{w}, b)$$



# Preparing Maximum Likelihood

$$\max_{\mathbf{w}, b} \left( \prod_{i:y_i=1} P(y_i = 1 | \mathbf{x}_i) \prod_{i:y_i=0} (1 - P(y_i = 1 | \mathbf{x}_i)) \right)$$



Let  $p_i = P(y_i = 1 | \mathbf{x}_i)$

$$\max_{\mathbf{w}, b} \left( \prod_i p_i^{y_i} (1 - p_i)^{(1-y_i)} \right)$$



# Preparing Maximum Log-Likelihood

$$\max_{\mathbf{w}, b} \left( \prod_i p_i^{y_i} (1 - p_i)^{(1-y_i)} \right)$$

↓ Take log

$$\max_{\mathbf{w}, b} \left( \sum_i y_i \log p_i + (1 - y_i) \log(1 - p_i) \right)$$

$$p_i = \frac{\exp(\mathbf{x}_i^T \hat{\mathbf{w}})}{1 + \exp(\mathbf{x}_i^T \hat{\mathbf{w}})}, \text{ where } \hat{\mathbf{w}} = (\mathbf{w}; b)$$

The loss function  $L(w)$

Solved with gradient descent

$$L(\mathbf{w}, b) = \sum_{i=1}^m \log p(y_i | x_i; \mathbf{w}, b)$$



# Journey for Solving the Logistic Regression

- Loss function

$$L(\hat{\mathbf{w}}) = \max_{\hat{\mathbf{w}}} \left( \sum_i y_i \log p_i + (1 - y_i) \log(1 - p_i) \right)$$

$$\downarrow \quad p_i = \frac{\exp(\mathbf{x}_i^T \hat{\mathbf{w}})}{1 + \exp(\mathbf{x}_i^T \hat{\mathbf{w}})}$$

$$\max_{\hat{\mathbf{w}}} \left( \sum_i y_i \mathbf{x}_i^T \hat{\mathbf{w}} - \log(1 + \exp(\mathbf{x}_i^T \hat{\mathbf{w}})) \right)$$

$$\min_{\hat{\mathbf{w}}} \left( \sum_i \log(1 + \exp(\mathbf{x}_i^T \hat{\mathbf{w}})) - y_i \mathbf{x}_i^T \hat{\mathbf{w}} \right)$$

Solved with gradient descent

The loss function  $L(\hat{\mathbf{w}})$



# Deduction (Again!)



# Understanding the Loss Function: Cross Entropy

$$L(\mathbf{w}, b) = f_{\mathbf{w}, b}(x_1)f_{\mathbf{w}, b}(x_2)\left(1 - f_{\mathbf{w}, b}(x_3)\right)\cdots f_{\mathbf{w}, b}(x_N)$$

$$-\log(L(\mathbf{w}, b)) = -[\log(f_{\mathbf{w}, b}(x_1)) + \log(f_{\mathbf{w}, b}(x_2)) + \log((1 - f_{\mathbf{w}, b}(x_3)))\cdots]$$

$y_i$ : 1 for class 1, 0 for class 2

$$= \sum_i -[y_i \log f_{\mathbf{w}, b}(x_i) + (1 - y_i) \log(1 - f_{\mathbf{w}, b}(x_i))]$$

Cross entropy between two Bernoulli distributions

Distribution  $p$ :

$$p(x = 1) = y_i$$

$$p(x = 0) = 1 - y_i$$

cross  
entropy

Distribution  $q$ :

$$q(x = 1) = f(x_i)$$

$$q(x = 0) = 1 - f(x_i)$$

$$H(p, q) = -\sum_{x \in \{0, 1\}} p(x) \log(q(x))$$



# Entropy and Cross Entropy

- In information theory, we like to describe the “*surprise*” of an event. An event is more surprising the less likely it is, meaning it contains more information.
  - Low Probability Event (surprising): More information.
  - Higher Probability Event (unsurprising): Less information.
- Information  $h(x)$  can be calculated for an event  $x$ , given the probability of the event  $P(x)$  as follows:
  - $h(x) = -\log(P(x))$

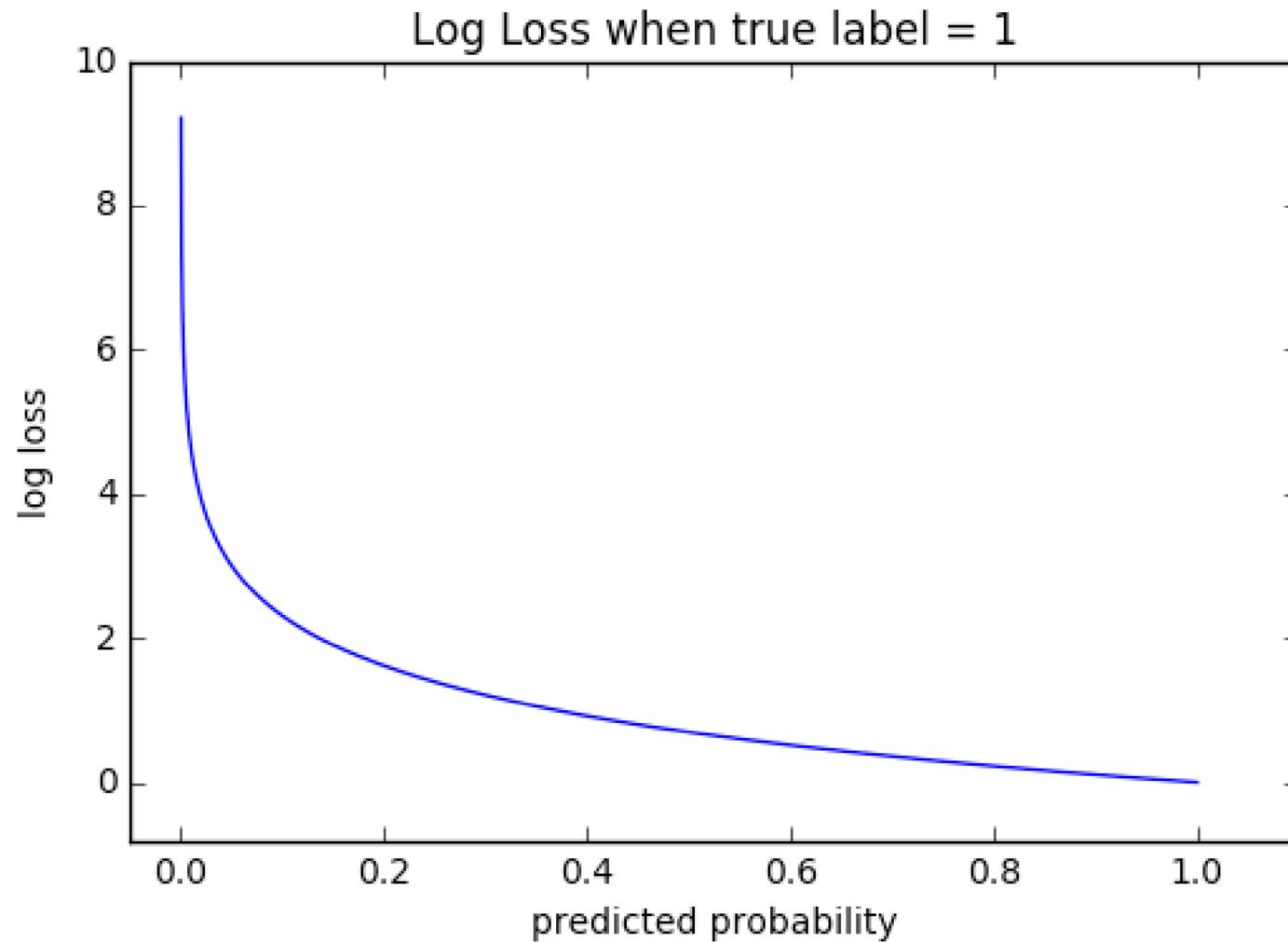


# Entropy and Cross Entropy

- **Entropy** is the number of bits required to transmit a randomly selected event from a probability distribution.
  - A skewed distribution has a low entropy, whereas a distribution where events have equal probability has a larger entropy.
- **Cross-entropy**: cross-entropy loss increases as the predicted probability diverges from the actual label.
  - Calculates the number of bits required to represent or transmit from one distribution (prediction) compared to another distribution (ground-truth).

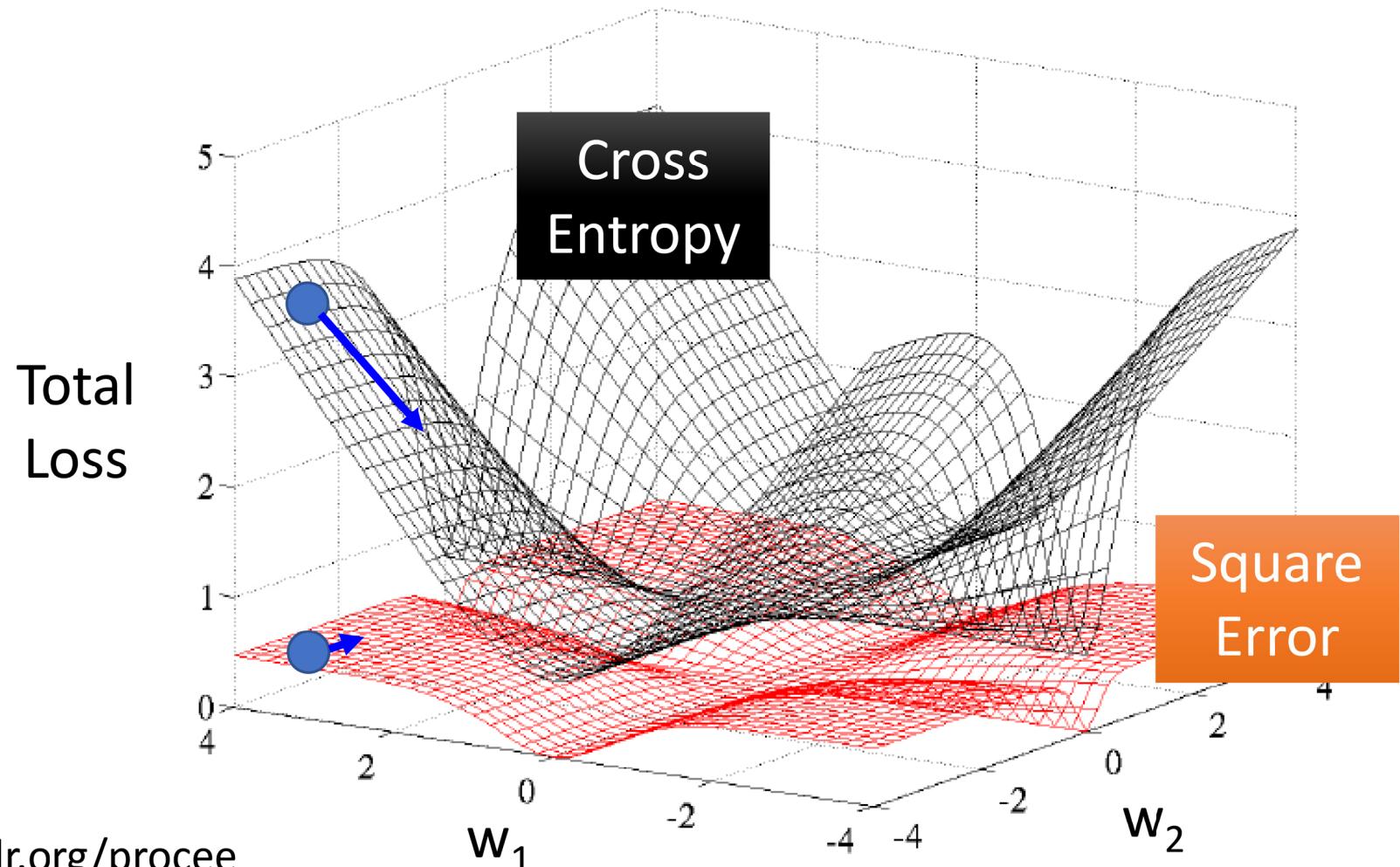


# Cross-Entropy





# Cross Entropy v.s. Square Error



<http://jmlr.org/proceedings/papers/v9/glorot10a/glorot10a.pdf>  
2024/9/19



# How to Find Gradient?

$$\min_{\hat{\mathbf{w}}} \left( \sum_i \log(1 + \exp(\mathbf{x}_i^T \hat{\mathbf{w}})) - y_i \mathbf{x}_i^T \hat{\mathbf{w}} \right)$$




# How to Find Gradient?

$$\min_{\hat{\mathbf{w}}} \left( \sum_i \log(1 + \exp(\mathbf{x}_i^T \hat{\mathbf{w}})) - y_i \mathbf{x}_i^T \hat{\mathbf{w}} \right)$$




# Gradient Descent

Iteratively find a  $\hat{w}$  that the gradient of the log probability function close to 0

1.  $\hat{w}^0 = 0$
2.  $\hat{w}^{(k+1)} = \hat{w}^{(k)} - \alpha \nabla_{\hat{w}^{(k)}} L(\hat{w}^{(k)})$

Until reach termination condition



# Relation with Neural Network (Spoiler Alert...)



# Just Another Example

- <https://www.youtube.com/watch?v=zHYF-o7tQ4o>



CSE 4820/5819

Introduction to Machine Learning  
University of Connecticut



# What We can Consider

- Shang Dynasty: Bone Crack (maybe related to temperature and humidity...) and Fortune of the Dynasty
- Temperature, humidity, and success: [80F, 10%] and 0; [60F, 70%] and 1;



# Example for Deduction: Logistic Regression Training

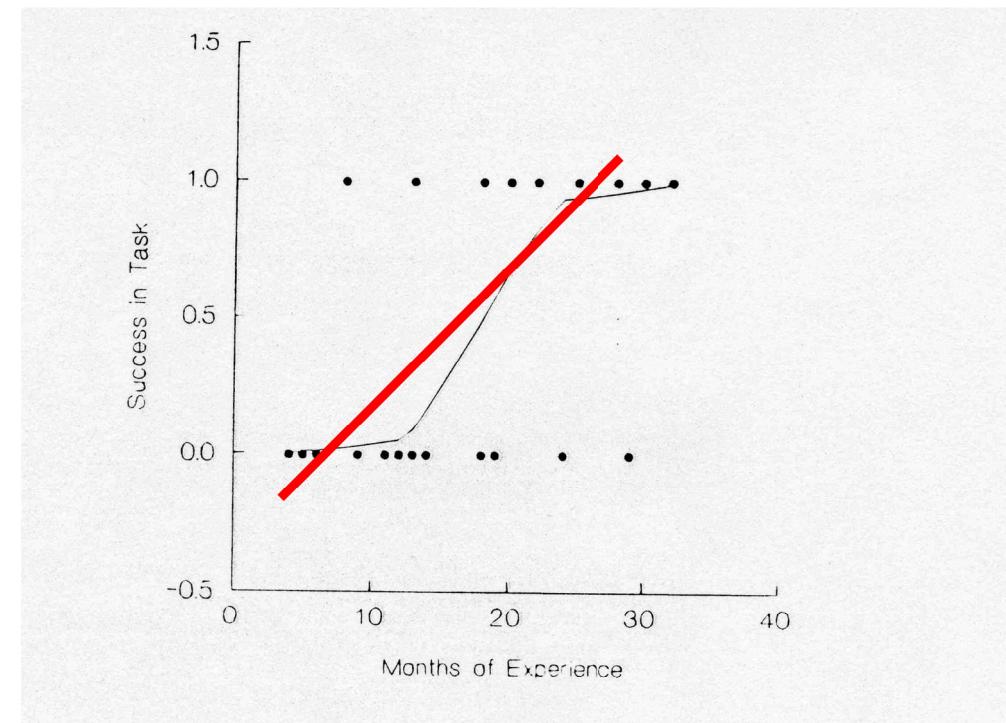


# Example for Deduction



# Another Example

- A systems analyst studied the effect of computer programming experience on ability to complete a task within a specified time
- Twenty-five persons selected for the study, with varying amounts of computer experience (in months)
- Results are coded in binary fashion:  $Y = 1$  if task completed successfully;  $Y = 0$ , otherwise





# Another Example

- Results from a standard package give:
  - $b = -3.0597$  and  $w = 0.1615$

- Estimated logistic regression function:

$$\hat{y} = \frac{1}{1 + e^{3.0597 - 0.1615X}}$$

- For example, the fitted value for  $X = 14$  is:

$$\hat{y} = \frac{1}{1 + e^{3.0597 - 0.1615(14)}} = 0.31$$

(Estimated probability that a person with 14 months experience will successfully complete the task)

CSE 4820/5819



# Similarities

- Linear Regression and Logistic Regression both are *supervised learning* algorithms.
- Linear Regression and Logistic Regression, both the models use *linear* models within the parameterization process



# Differences

- Linear Regression is used to handle *regression* problems whereas Logistic regression is used to handle the *classification* problems.
- Linear regression provides a *continuous* output, but Logistic regression provides *discrete* output.
- The purpose of Linear Regression is to find the best-fitted line while Logistic regression is one step ahead and fitting the line values to the sigmoid curve.
- The method for calculating loss function in linear regression is the *mean squared error* whereas for logistic regression it is *maximum likelihood estimation*.



# Conclusion

- Background
- Key Ideas
- Gradient Descent for Logistic Regression