

# ECE595 / STAT598: Machine Learning I

## Lecture .1 Logistic Regression 1 - From Linear to Logistic

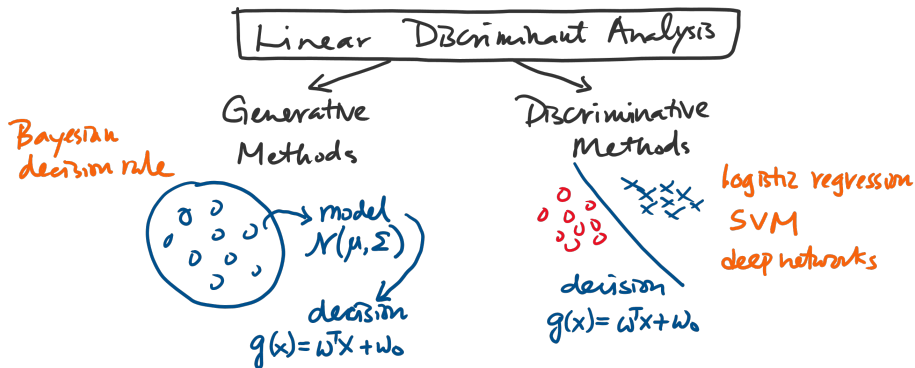
Spring 2020

Stanley Chan

School of Electrical and Computer Engineering  
Purdue University



# Overview



- In linear discriminant analysis (LDA), there are generally two types of approaches
- **Generative approach:** Estimate model, then define the classifier
- **Discriminative approach:** Directly define the classifier

# Outline

## Discriminative Approaches

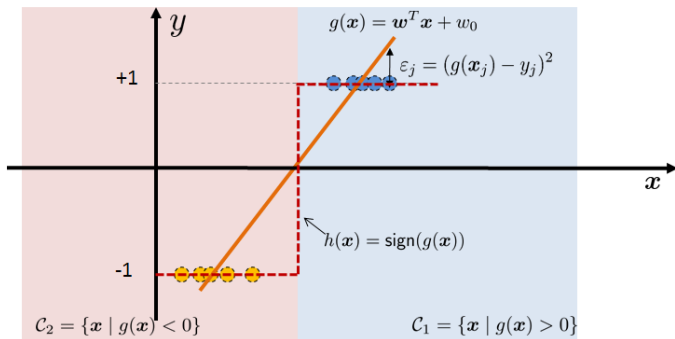
- Lecture 14 Logistic Regression 1
- Lecture 15 Logistic Regression 2

## This lecture: Logistic Regression 1

- From Linear to Logistic
  - Motivation
  - Loss Function
  - Why not L2 Loss?
- Interpreting Logistic
  - Maximum Likelihood
  - Log-odd
- Convexity
  - Is logistic loss convex?
  - Computation

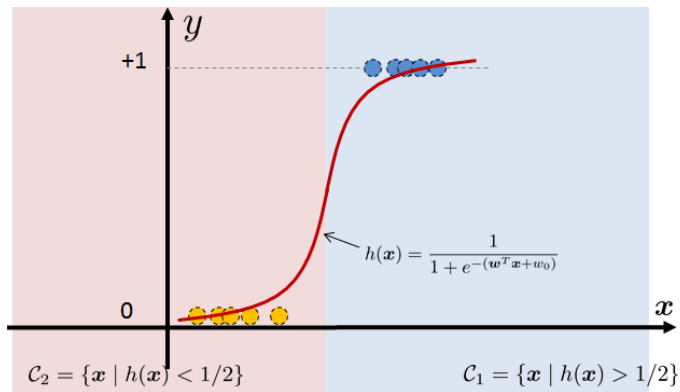
# Geometry of Linear Regression

- The discriminant function  $g(\mathbf{x})$  is linear
- The hypothesis function  $h(\mathbf{x}) = \text{sign}(g(\mathbf{x}))$  is a unit step



# From Linear to Logistic Regression

- Can we replace  $g(\mathbf{x})$  by  $\text{sign}(g(\mathbf{x}))$ ?
- How about a soft-version of  $\text{sign}(g(\mathbf{x}))$ ?
- This gives a logistic regression.



# Sigmoid Function

- The function

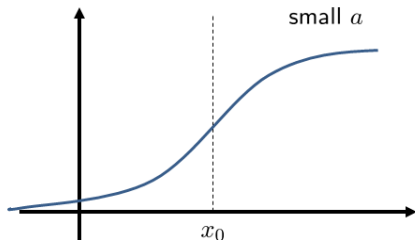
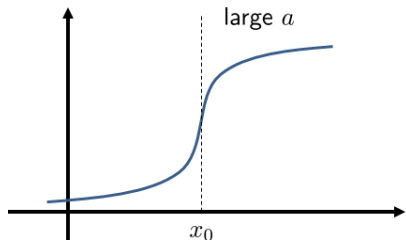
$$h(\mathbf{x}) = \frac{1}{1 + e^{-g(\mathbf{x})}} = \frac{1}{1 + e^{-(\mathbf{w}^T \mathbf{x} + w_0)}}$$

is called a **sigmoid function**.

- Its 1D form is

$$h(x) = \frac{1}{1 + e^{-a(x-x_0)}}, \quad \text{for some } a \text{ and } x_0,$$

- $a$  controls the transient speed
- $x_0$  controls the cutoff location



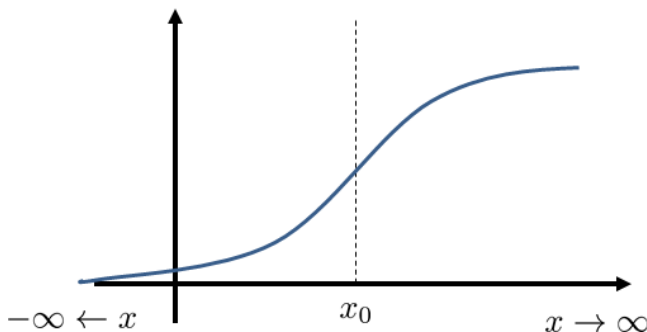
# Sigmoid Function

- Note that

$$h(x) \rightarrow 1, \quad \text{as } x \rightarrow \infty,$$

$$h(x) \rightarrow 0, \quad \text{as } x \rightarrow -\infty,$$

- So  $h(x)$  can be regarded as a “probability”.



# Sigmoid Function

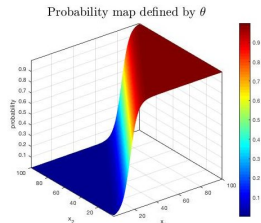
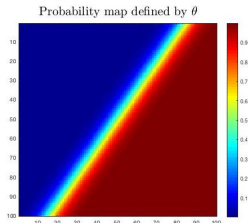
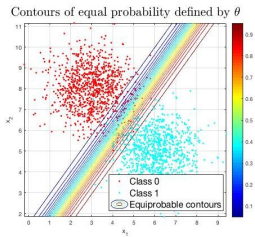
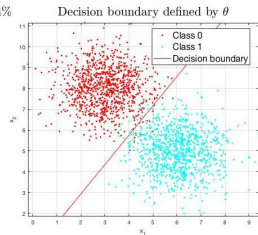
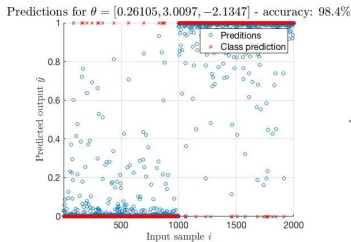
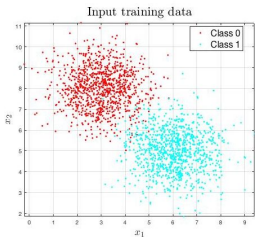
- Derivative is

$$\begin{aligned}\frac{d}{dx} \left( \frac{1}{1 + e^{-a(x-x_0)}} \right) &= - \left( 1 + e^{-a(x-x_0)} \right)^{-2} \left( e^{-a(x-x_0)} \right) (-a) \\ &= a \left( \frac{e^{-a(x-x_0)}}{1 + e^{-a(x-x_0)}} \right) \left( \frac{1}{1 + e^{-a(x-x_0)}} \right) \\ &= a \left( 1 - \frac{1}{1 + e^{-a(x-x_0)}} \right) \left( \frac{1}{1 + e^{-a(x-x_0)}} \right) \\ &= a[1 - h(x)][h(x)].\end{aligned}$$

- Since  $0 < h(x) < 1$ , we have  $0 < 1 - h(x) < 1$ .
- Therefore, the derivative is always positive.
- So  $h$  is an increasing function.
- Hence  $h$  can be considered as a “CDF”.

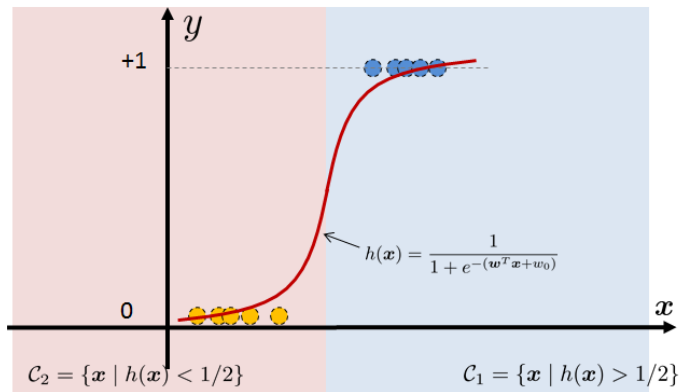


# Sigmoid Function



# From Linear to Logistic Regression

- Can we replace  $g(\mathbf{x})$  by  $\text{sign}(g(\mathbf{x}))$ ?
- How about a soft-version of  $\text{sign}(g(\mathbf{x}))$ ?
- This gives a logistic regression.



# Loss Function for Linear Regression

- All discriminant algorithms have a **Training Loss Function**

$$J(\theta) = \frac{1}{N} \sum_{n=1}^N \mathcal{L}(g(\mathbf{x}_n), y_n).$$

- In linear regression,

$$\begin{aligned} J(\theta) &= \frac{1}{N} \sum_{n=1}^N (g(\mathbf{x}_n) - y_n)^2 \\ &= \frac{1}{N} \sum_{n=1}^N (\mathbf{w}^T \mathbf{x}_n + w_0 - y_n)^2 \\ &= \frac{1}{N} \left\| \begin{bmatrix} \mathbf{x}_1^T & 1 \\ \vdots & \vdots \\ \mathbf{x}_N^T & 1 \end{bmatrix} \begin{bmatrix} \mathbf{w} \\ w_0 \end{bmatrix} - \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix} \right\|^2 = \frac{1}{N} \|\mathbf{A}\theta - \mathbf{y}\|^2. \end{aligned}$$

# Training Loss for Logistic Regression

$$\begin{aligned} J(\theta) &= \sum_{n=1}^N \mathcal{L}(h_{\theta}(\mathbf{x}_n), y_n) \\ &= \sum_{n=1}^N -\left\{ y_n \log h_{\theta}(\mathbf{x}_n) + (1 - y_n) \log(1 - h_{\theta}(\mathbf{x}_n)) \right\} \end{aligned}$$

- This loss is also called the **cross-entropy loss**.
- Why do we want to choose this cost function?
- Consider two cases

$$y_n \log h_{\theta}(\mathbf{x}_n) = \begin{cases} 0, & \text{if } y_n = 1, \text{ and } h_{\theta}(\mathbf{x}_n) = 1, \\ -\infty, & \text{if } y_n = 1, \text{ and } h_{\theta}(\mathbf{x}_n) = 0, \end{cases}$$

$$(1 - y_n)(1 - \log h_{\theta}(\mathbf{x}_n)) = \begin{cases} 0, & \text{if } y_n = 0, \text{ and } h_{\theta}(\mathbf{x}_n) = 0, \\ -\infty, & \text{if } y_n = 0, \text{ and } h_{\theta}(\mathbf{x}_n) = 1. \end{cases}$$

- No solution if mismatch

## Why Not L2 Loss?

- Why not use L2 loss?

$$J(\theta) = \sum_{n=1}^N (h_{\theta}(\mathbf{x}_n) - y_n)^2$$

- Let's look at the 1D case:

$$J(\theta) = \left( \frac{1}{1 + e^{-\theta x}} - y \right)^2.$$

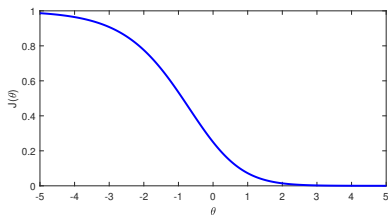
- This is NOT convex!
- How about the logistic loss?
- 

$$J(\theta) = y \log \left( \frac{1}{1 + e^{-\theta x}} \right) + (1 - y) \log \left( 1 - \frac{1}{1 + e^{-\theta x}} \right)$$

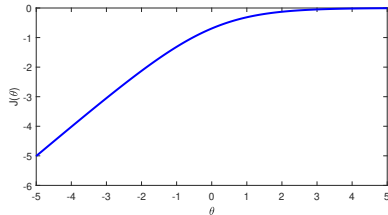
- This is convex!

## Why Not L2 Loss?

- Experiment: Set  $x = 1$  and  $y = 1$ .
- Plot  $J(\theta)$  as a function of  $\theta$ .



L2



Logistic

- So the L2 loss is not convex, but the logistic loss is concave (negative is convex)
- If you do gradient descent on L2, you will be trapped at local minima