

# ECE595 / STAT598: Machine Learning I

## Lecture 13.1: Connecting Bayesian with Linear Regression - Linear Regression Review

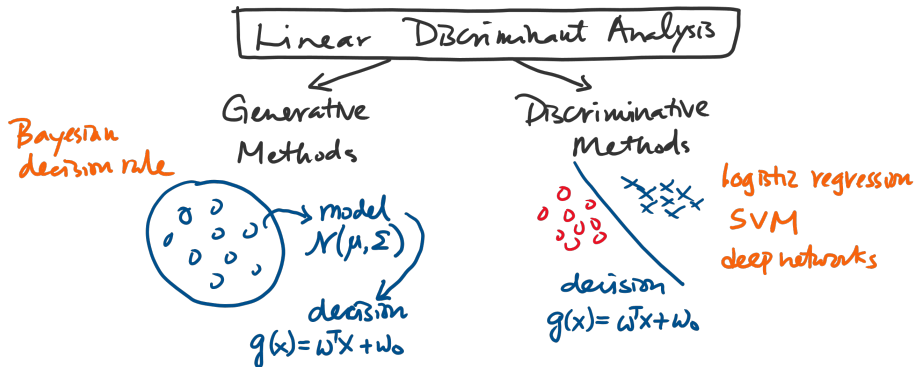
Spring 2020

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

## Generative Approaches

- Lecture 9 Bayesian Decision Rules
- Lecture 10 Evaluating Performance
- Lecture 11 Parameter Estimation
- Lecture 12 Bayesian Prior
- **Lecture 13 Connecting Bayesian and Linear Regression**

## Today's Lecture

- **Linear Regression Review**
  - **Linear regression in the context of classification**
  - **Linking linear regression with MLE and MAP**
- Connection between Linear Regression and Bayesian
  - Expected Loss
  - Main Result
  - Implications

# Linear Regression Reviewed

- Linear regression is actually a **discriminative method**.
- Do not require a distributional model.
- Construct the hypothesis function directly:

$$h(\mathbf{x}) = \begin{cases} +1, & \text{if } g(\mathbf{x}) > 0, \\ -1, & \text{if } g(\mathbf{x}) < 0. \end{cases}$$

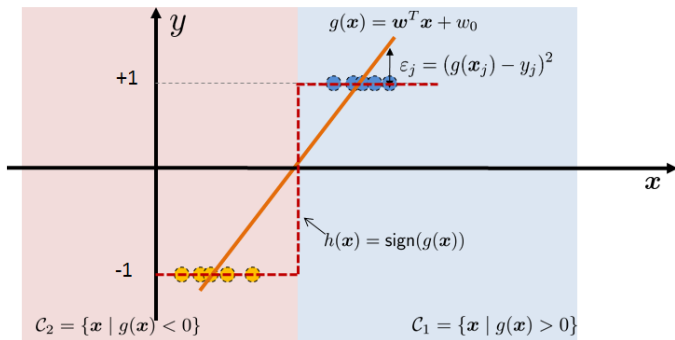
- Consider a binary classification problem with discriminant function:

$$g(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$$

- The goal is to determine the parameters  $\theta = \{\mathbf{w}, w_0\}$
- Training data:  $(\mathbf{x}_n, y_n)_{n=1}^N$ 
  - $\mathbf{x}_n \in \mathbb{R}^d$  is the input vector
  - $y_n \in \{-1, +1\}$  is the corresponding label

# 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



# Loss Function

- 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,

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

# Solution of Linear Regression

## Theorem (Linear Regression Solution)

The loss function of a linear regression model is given by

$$J(\theta) = \|\mathbf{A}\theta - \mathbf{y}\|^2,$$

of which the minimizer is

$$\theta^* = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{y}.$$

- Take derivative and setting to zero:

$$\begin{aligned} \nabla_{\theta} J(\theta) &= \nabla_{\theta} \{ \|\mathbf{A}\theta - \mathbf{y}\|^2 \} \\ &= 2\mathbf{A}^T (\mathbf{A}\theta - \mathbf{y}) = \mathbf{0}. \end{aligned}$$

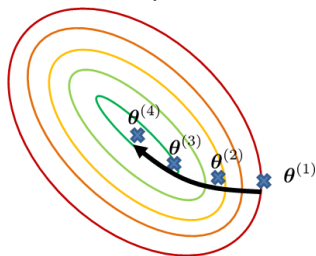
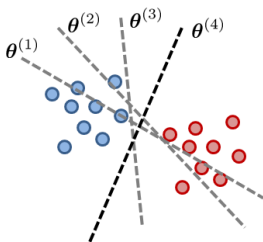
- So solution is  $\theta^* = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{y}$ , assuming  $\mathbf{A}^T \mathbf{A}$  is invertible.

## When $\mathbf{A}^T \mathbf{A}$ is large

- Computing  $(\mathbf{A}^T \mathbf{A})^{-1}$  directly is infeasible for large-scale datasets with a large number of variables
- Consider using iterative algorithms such as gradient descent
- The gradient descent is given by the iteration:

$$\begin{aligned}\boldsymbol{\theta}^{(k+1)} &= \boldsymbol{\theta}^{(k)} - \eta \nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}^{(k)}) \\ &= \boldsymbol{\theta}^{(k)} - \eta (2\mathbf{A}^T \mathbf{A} \boldsymbol{\theta}^{(k)} - 2\mathbf{A}^T \mathbf{y})\end{aligned}$$

- A pictorial illustration of the gradient descent step:





# Treating Linear Regression as Maximum-Likelihood

- Minimizing  $J(\theta)$  is the same as solving a **maximum-likelihood**:

$$\begin{aligned}\theta^* &= \operatorname{argmin}_{\theta} \|\mathbf{A}\theta - \mathbf{y}\|^2 \\ &= \operatorname{argmin}_{\theta} \sum_{n=1}^N (\mathbf{a}_n^T \theta - y_n)^2 \\ &= \operatorname{argmax}_{\theta} \exp \left\{ - \sum_{n=1}^N (\mathbf{a}_n^T \theta - y_n)^2 \right\} \\ &= \operatorname{argmax}_{\theta} \prod_{n=1}^N \left\{ \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left\{ - \frac{(\mathbf{a}_n^T \theta - y_n)^2}{2\sigma^2} \right\} \right\}\end{aligned}$$

- Assume noise is i.i.d. Gaussian with variance  $\sigma^2$ .

# Treating Linear Regression as Maximum-a-Posteriori

- We can modify the MLE by adding a prior

$$p_{\Theta}(\boldsymbol{\theta}) = \exp \left\{ -\frac{\rho(\boldsymbol{\theta})}{\beta} \right\}.$$

- Then, we have a MAP problem:

$$\begin{aligned}\boldsymbol{\theta}^* &= \operatorname{argmax}_{\boldsymbol{\theta}} \prod_{n=1}^N \left\{ \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left\{ -\frac{(\mathbf{a}_n^T \boldsymbol{\theta} - y_n)^2}{2\sigma^2} \right\} \right\} \exp \left\{ -\frac{\rho(\boldsymbol{\theta})}{\beta} \right\} \\ &= \operatorname{argmin}_{\boldsymbol{\theta}} \frac{1}{2\sigma^2} \sum_{n=1}^N (\mathbf{a}_n^T \boldsymbol{\theta} - y_n)^2 + \frac{1}{\beta} \rho(\boldsymbol{\theta}) \\ &= \operatorname{argmin}_{\boldsymbol{\theta}} \|\mathbf{A}\boldsymbol{\theta} - \mathbf{y}\|^2 + \lambda \rho(\boldsymbol{\theta}), \quad \text{where } \lambda = 2\sigma^2/\beta.\end{aligned}$$

- $\rho(\cdot)$  is called **regularization function**.
- Useful when  $\mathbf{A}^T \mathbf{A}$  is not invertible.

# Ridge Regression

- One option: Choose a Gaussian prior

$$\exp \left\{ -\frac{\rho(\boldsymbol{\theta})}{\beta} \right\} = \exp \left\{ -\frac{\|\boldsymbol{\theta}\|^2}{2\sigma_0^2} \right\}$$

- Then, the MAP becomes

$$\begin{aligned}\boldsymbol{\theta}^* &= \operatorname{argmax}_{\boldsymbol{\theta}} \prod_{n=1}^N \left\{ \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left\{ -\frac{(\mathbf{a}_n^T \boldsymbol{\theta} - y_n)^2}{2\sigma^2} \right\} \right\} \exp \left\{ -\frac{\|\boldsymbol{\theta}\|^2}{2\sigma_0^2} \right\} \\ &= \operatorname{argmin}_{\boldsymbol{\theta}} \sum_{n=1}^N (\mathbf{a}_n^T \boldsymbol{\theta} - y_n)^2 + \underbrace{\frac{\sigma^2}{\sigma_0^2}}_{=\lambda} \|\boldsymbol{\theta}\|^2 \\ &= \operatorname{argmin}_{\boldsymbol{\theta}} \|\mathbf{A}\boldsymbol{\theta} - \mathbf{y}\|^2 + \lambda \|\boldsymbol{\theta}\|^2\end{aligned}$$

- This is called **Tikhonov regularization** or **Ridge regression**.