

ECE 595: Machine Learning I

Lecture 9.1: Bayesian Decision - Review of High-Dimensional Gaussian

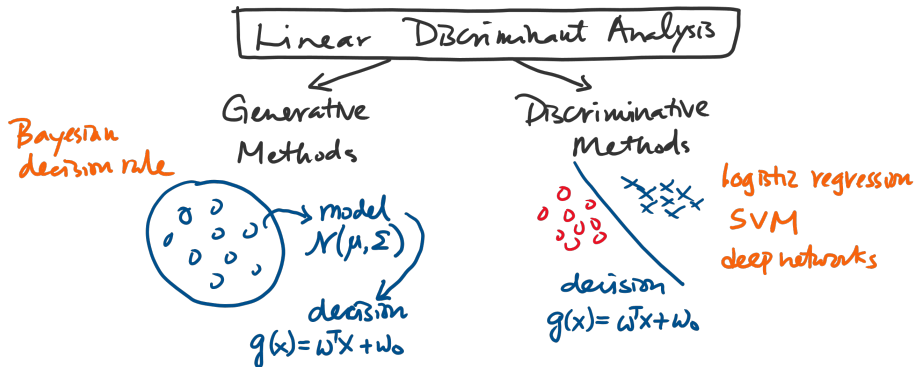
Spring 2020

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

Generative Approach

Goal: Construct a discriminant function $g(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$ from the data.

- Suppose there are two classes C_1 and C_2 .
- Each class is modeled as a Gaussian.
- We are going to utilize two concepts:
- **likelihood function**

$$p_{\mathbf{X}|Y}(\mathbf{x}|i) = \mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)$$

- **prior** distribution

$$p_Y(i) = \pi_i$$

Outline

Generative Approaches

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

Today's Lecture

- Review of High-Dimensional Gaussian
 - Likelihood and prior
 - Gaussian PDF
- Basic Principle
 - Making the Bayesian decision
 - 1D Illustration
- The Three Cases
 - $\Sigma_j = \sigma^2 I$
 - $\Sigma_j = \Sigma$ (Next Lecture)
 - General Σ_j (Next Lecture)

High-dimensional Gaussian

An d -dimensional **Gaussian** has a PDF

$$p_{\mathbf{X}}(\mathbf{x}) = \frac{1}{\sqrt{(2\pi)^d |\boldsymbol{\Sigma}|}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\},$$

where d denotes the dimensionality of the vector \mathbf{x} .

- The **mean vector** $\boldsymbol{\mu}$ is

$$\boldsymbol{\mu} = \mathbb{E}[\mathbf{X}] = \begin{bmatrix} \mathbb{E}[X_1] \\ \vdots \\ \mathbb{E}[X_d] \end{bmatrix}$$

- The **covariance matrix** $\boldsymbol{\Sigma}$ is

$$\boldsymbol{\Sigma} = \mathbb{E}[(\mathbf{X} - \boldsymbol{\mu})(\mathbf{X} - \boldsymbol{\mu})^T] = \begin{bmatrix} \text{Var}[X_1] & \text{Cov}(X_1, X_2) & \dots & \text{Cov}(X_1, X_d) \\ \text{Cov}(X_2, X_1) & \text{Var}[X_2] & \dots & \text{Cov}(X_2, X_d) \\ \vdots & \vdots & \ddots & \vdots \\ \text{Cov}(X_d, X_1) & \text{Cov}(X_d, X_2) & \dots & \text{Var}[X_d] \end{bmatrix}$$

- $\boldsymbol{\Sigma}$ is always positive semi-definite. (Why?)

Special Case: Diagonal Covariance

- Suppose that X_i and X_j are independent for all $i \neq j$.
- This implies $\text{Cov}(X_i, X_j) = 0$
- Simplify Σ

$$\Sigma = \begin{bmatrix} \sigma_1^2 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \sigma_d^2 \end{bmatrix},$$

- Then, the exponential is

$$(\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}) = \sum_{i=1}^n \frac{(x_i - \mu_i)^2}{\sigma_i^2}.$$

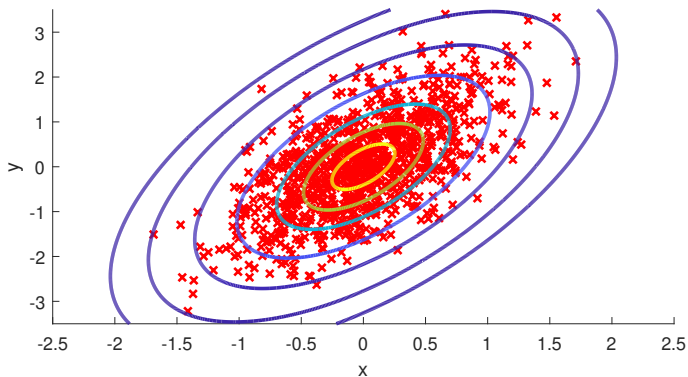
- And hence, the PDF is

$$p_{\mathbf{X}}(\mathbf{x}) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma_i^2}} \exp \left\{ -\frac{(x_i - \mu_i)^2}{2\sigma_i^2} \right\}.$$

Visualization

- Generate 1000 random samples from a 2D Gaussian

- $\boldsymbol{\mu} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$, and $\boldsymbol{\Sigma} = \begin{bmatrix} 0.25 & 0.3 \\ 0.3 & 1 \end{bmatrix}$



Conditional Gaussian

- Data $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$.
- Class $Y \in \{1, 2, \dots, K\}$.
- **Likelihood:**

$p_{\mathbf{X}|Y}(\mathbf{x}|k)$ = Probability of getting \mathbf{X} given Y

- **Prior:**

$p_Y(k)$ = Probability of getting Y

- **Posterior:**

$p_{Y|\mathbf{X}}(k|\mathbf{x})$ = Probability of getting Y given \mathbf{X}

- Related by

$$p_{Y|\mathbf{X}}(k|\mathbf{x}) = \frac{p_{\mathbf{X}|Y}(\mathbf{x}|k)p_Y(k)}{p_{\mathbf{X}}(\mathbf{x})} = \frac{p_{\mathbf{X}|Y}(\mathbf{x}|k)p_Y(k)}{\sum_k p_{\mathbf{X}|Y}(\mathbf{x}|k)p_Y(k)}$$

Example

- Two Gaussian $\mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1)$ and $\mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2)$.
- **Prior** probability of getting a class is

$$p_Y(1) = \pi_1 \quad \text{and} \quad p_Y(2) = \pi_2.$$

- The **likelihood** term is

$$\begin{aligned} p_{\mathbf{X}|Y}(\mathbf{x}|k) &= \mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \\ &= \frac{1}{\sqrt{(2\pi)^d |\boldsymbol{\Sigma}_k|}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1} (\mathbf{x} - \boldsymbol{\mu}_k) \right\} \end{aligned}$$

- The **posterior** is

$$\begin{aligned} p_{Y|\mathbf{X}}(k|\mathbf{x}) &= \frac{p_{\mathbf{X}|Y}(\mathbf{x}|k)p_Y(k)}{p_{\mathbf{X}}(\mathbf{x})} \\ &= \frac{\frac{1}{\sqrt{(2\pi)^d |\boldsymbol{\Sigma}_k|}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1} (\mathbf{x} - \boldsymbol{\mu}_k) \right\} \cdot \pi_k}{\sum_{k=1}^K \frac{1}{\sqrt{(2\pi)^d |\boldsymbol{\Sigma}_k|}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1} (\mathbf{x} - \boldsymbol{\mu}_k) \right\} \cdot \pi_k} \end{aligned}$$

Negative Log-Likelihood

Negative Log-Likelihood for Gaussian:

$$\begin{aligned} & -\log p_{\mathbf{X}|Y}(\mathbf{x}|k) \\ &= -\log \left(\frac{1}{\sqrt{(2\pi)^d |\boldsymbol{\Sigma}_k|}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1} (\mathbf{x} - \boldsymbol{\mu}_k) \right\} \right) \\ &= \underbrace{\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1} (\mathbf{x} - \boldsymbol{\mu}_k)}_{\text{contains } \mathbf{x}} \underbrace{- \frac{n}{2} \log 2\pi - \frac{1}{2} \log |\boldsymbol{\Sigma}_k|}_{\text{no } \mathbf{x}}. \end{aligned}$$

- $(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \geq 0$, always.
- $\sqrt{(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})}$ is called **Mahalanobis distance**.