

ECE 595: Machine Learning I

Lecture 9.3: Bayesian Decision - The Three Cases I

Spring 2020

Stanley Chan

School of Electrical and Computer Engineering
Purdue University



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)

Three Cases of Gaussians

Discriminant function of Gaussian:

$$\begin{aligned}g_i(\mathbf{x}) &= \log p_{\mathbf{X}|Y}(\mathbf{x}|i) + \log \pi_i \\ &= -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_i)^T \boldsymbol{\Sigma}_i^{-1}(\mathbf{x} - \boldsymbol{\mu}_i) - \frac{1}{2} \log |\boldsymbol{\Sigma}_i| + \log \pi_i.\end{aligned}$$

- $\boldsymbol{\Sigma}_i = \sigma^2 \mathbf{I}$
 - All Gaussians have the same covariance matrix
 - The covariance matrix is diagonal and same variance
- $\boldsymbol{\Sigma}_i = \boldsymbol{\Sigma}$
 - All Gaussians have the same covariance matrix
 - The covariance matrix can be anything
- arbitrary $\boldsymbol{\Sigma}_i$
 - Any positive semi-definite covariance matrix

Case 1: $\Sigma_i = \sigma^2 I$

Put $\Sigma_i = \Sigma$:

$$g_i(\mathbf{x}) = -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_i)^T \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu}_i) - \frac{1}{2} \log |\Sigma| + \log \pi_i.$$

Let us do some simplification:

$$\begin{aligned} g_i(\mathbf{x}) &= -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_i)^T \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu}_i) - \cancel{\frac{1}{2} \log |\Sigma|} + \log \pi_i \\ &= -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_i)^T \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu}_i) + \log \pi_i \\ &= -\frac{1}{2\sigma^2} \|\mathbf{x} - \boldsymbol{\mu}_i\|^2 + \log \pi_i \\ &= -\frac{1}{2\sigma^2} \left(\|\mathbf{x}\|^2 - 2\mathbf{x}^T \boldsymbol{\mu}_i + \|\boldsymbol{\mu}_i\|^2 \right) + \log \pi_i \\ &= -\frac{1}{2\sigma^2} \left(\cancel{\|\mathbf{x}\|^2} - 2\mathbf{x}^T \boldsymbol{\mu}_i + \|\boldsymbol{\mu}_i\|^2 \right) + \log \pi_i \\ &= \left(\frac{\boldsymbol{\mu}_i}{\sigma^2} \right)^T \mathbf{x} - \left(\frac{\|\boldsymbol{\mu}_i\|^2}{2\sigma^2} - \log \pi_i \right). \end{aligned}$$

Case 1: $\Sigma_i = \sigma^2 I$

$$\begin{aligned} g_i(\mathbf{x}) &= \underbrace{\left(\frac{\boldsymbol{\mu}_i}{\sigma^2}\right)^T}_{\mathbf{w}_i} \mathbf{x} - \underbrace{\left(\frac{\|\boldsymbol{\mu}_i\|^2}{2\sigma^2} - \log \pi_i\right)}_{w_{i0}} \\ &= \mathbf{w}_i^T \mathbf{x} + w_{i0} \end{aligned}$$

So if the i -th and the j -th discriminant functions are

$$g_i(\mathbf{x}) = \mathbf{w}_i^T \mathbf{x} + w_{i0}$$

$$g_j(\mathbf{x}) = \mathbf{w}_j^T \mathbf{x} + w_{j0},$$

then,

$$\begin{aligned} g(\mathbf{x}) &= g_i(\mathbf{x}) - g_j(\mathbf{x}) \\ &= \underbrace{\left(\frac{\boldsymbol{\mu}_i - \boldsymbol{\mu}_j}{\sigma^2}\right)^T}_{\mathbf{w}_i - \mathbf{w}_j} \mathbf{x} + \underbrace{\left(-\frac{\|\boldsymbol{\mu}_i\|^2 - \|\boldsymbol{\mu}_j\|^2}{2\sigma^2} + \log \frac{\pi_i}{\pi_j}\right)}_{w_{i0} - w_{j0}}. \end{aligned}$$

Case 1: $\Sigma_i = \sigma^2 I$

Theorem

If $\Sigma_i = \sigma^2 I$, then the separating hyperplane is given by

$$g(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0 = 0,$$

where

$$\mathbf{w} = \frac{\boldsymbol{\mu}_i - \boldsymbol{\mu}_j}{\sigma^2}, \quad \text{and} \quad w_0 = -\frac{\|\boldsymbol{\mu}_i\|^2 - \|\boldsymbol{\mu}_j\|^2}{2\sigma^2} + \log \frac{\pi_i}{\pi_j}.$$

- You tell me the two Gaussians: $\boldsymbol{\mu}_i, \boldsymbol{\mu}_j, \pi_i, \pi_j, \sigma$
- I return you a separating hyperplane

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

- This is the best possible hyperplane according to posterior distribution

Case 1: $\Sigma_i = \sigma^2 I$: Geometry

Can we write $g(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$ in terms of

$$g(\mathbf{x}) = \mathbf{w}^T (\mathbf{x} - \mathbf{x}_0).$$

Not too difficult:

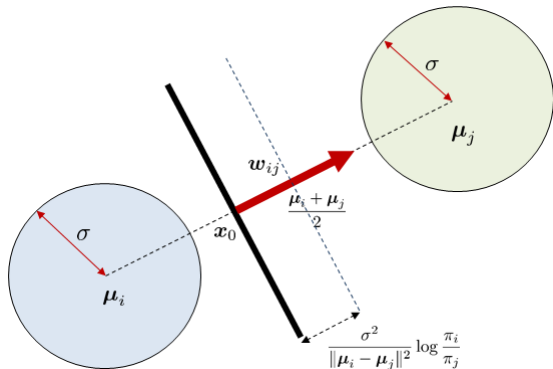
$$\begin{aligned} g(\mathbf{x}) &= \left(\frac{\boldsymbol{\mu}_i - \boldsymbol{\mu}_j}{\sigma^2} \right)^T \mathbf{x} - \left(\frac{\|\boldsymbol{\mu}_i\|^2}{2\sigma^2} - \frac{\|\boldsymbol{\mu}_j\|^2}{2\sigma^2} \right) + \log \frac{\pi_i}{\pi_j} \\ &= \left(\frac{\boldsymbol{\mu}_i - \boldsymbol{\mu}_j}{\sigma^2} \right)^T \left[\mathbf{x} - \underbrace{\frac{\boldsymbol{\mu}_i + \boldsymbol{\mu}_j}{2} + \sigma^2 \left(\log \frac{\pi_i}{\pi_j} \right) \frac{\boldsymbol{\mu}_i - \boldsymbol{\mu}_j}{\|\boldsymbol{\mu}_i - \boldsymbol{\mu}_j\|^2}}_{\mathbf{x}_0} \right] \end{aligned}$$

Therefore, we have

$$\mathbf{w} = \frac{\boldsymbol{\mu}_i - \boldsymbol{\mu}_j}{\sigma^2}, \quad \text{and} \quad \mathbf{x}_0 = \frac{\boldsymbol{\mu}_i + \boldsymbol{\mu}_j}{2} - \frac{\sigma^2}{\|\boldsymbol{\mu}_i - \boldsymbol{\mu}_j\|^2} \left(\log \frac{\pi_i}{\pi_j} \right) (\boldsymbol{\mu}_i - \boldsymbol{\mu}_j),$$

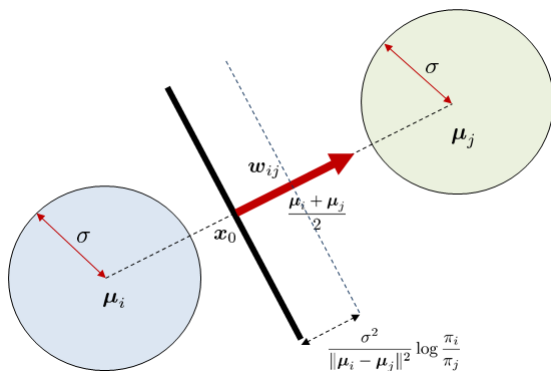
Case 1: $\Sigma_i = \sigma^2 I$: Geometry

$$\mathbf{w} = \frac{\boldsymbol{\mu}_i - \boldsymbol{\mu}_j}{\sigma^2}, \quad \text{and} \quad \mathbf{x}_0 = \frac{\boldsymbol{\mu}_i + \boldsymbol{\mu}_j}{2} - \frac{\sigma^2}{\|\boldsymbol{\mu}_i - \boldsymbol{\mu}_j\|^2} \left(\log \frac{\pi_i}{\pi_j} \right) (\boldsymbol{\mu}_i - \boldsymbol{\mu}_j),$$



Interpreting Results

Here are the geometric interpretations:



- Normal vector is $\mathbf{w} = \frac{\mu_i - \mu_j}{\sigma^2}$. It points from one center to another.
- Midpoint is $\mathbf{x}_0 = \frac{\mu_i + \mu_j}{2}$
- The prior creates an offset. Offset direction is also $\mu_i - \mu_j$. If $\pi_i = \pi_j = 1/2$, then $\log(\pi_i/\pi_j) = 0$.

Reading List

High Dimensional Gaussian

- Bishop, Pattern Recognition and Machine Learning, Chapter 2.3
- Stanford CS 229 Tutorial on Gaussian
<http://cs229.stanford.edu/section/gaussians.pdf>

Bayesian Decision Rule

- Bishop, Pattern Recognition and Machine Learning, Chapter 4.1
- Duda, Hart and Stork's Pattern Classification, Chapter 2.1, 2.2, 2.6
- Stanford CS 229 Generative Algorithms
<http://cs229.stanford.edu/notes/cs229-notes2.pdf>
- UCSD ECE 271A, Lecture 4 and 5
<http://www.svcl.ucsd.edu/courses/ece271A/ece271A.htm>