

ECE 595: Machine Learning I

Lecture 09 Bayesian Decision

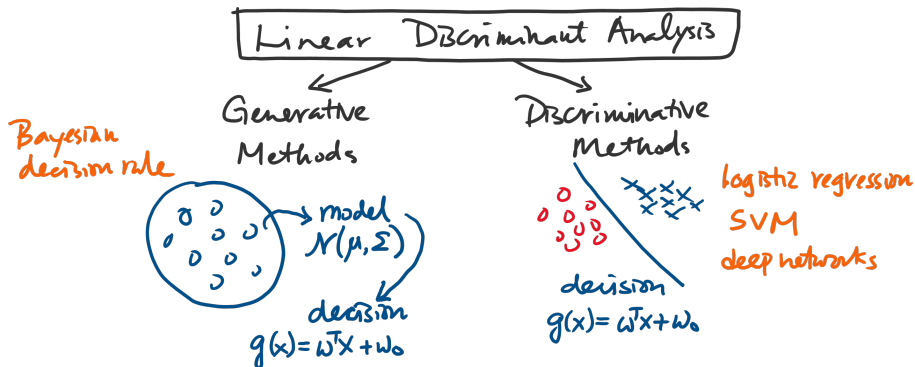
Spring 2020

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

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

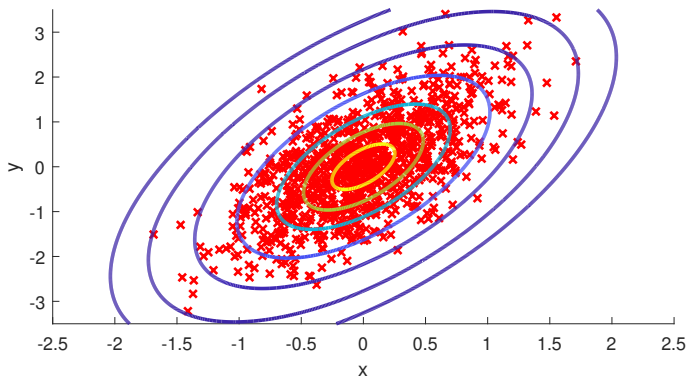
- And hence, the PDF is

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

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**.

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)

Interaction between Likelihood and Prior

- According to **Bayes Theorem**, we have that

$$p_{Y|X}(i|x) = \frac{p_{X|Y}(x|i)p_Y(i)}{p_X(x)}$$

- Posterior: **After** you have seen x
- Likelihood: **Before** you see x
- Prior: You subjective believe of class label

- You cannot just use $p_Y(i)$; Otherwise you are not using data
- You cannot just use $p_{X|Y}(x|i)$; Otherwise you cannot explain “ Y given X ”

Making the Bayesian Decision

Which class is more likely?

$$\begin{aligned}i^* &= \operatorname{argmax}_i p_{Y|X}(i|\mathbf{x}) \\&= \operatorname{argmax}_i \frac{p_{X|Y}(\mathbf{x}|i)p_Y(i)}{p_X(\mathbf{x})} \\&= \operatorname{argmax}_i \log p_{X|Y}(\mathbf{x}|i) + \log \pi_i - \log p_X(\mathbf{x}) \\&= \operatorname{argmax}_i \log p_{X|Y}(\mathbf{x}|i) + \log \pi_i - \cancel{\log p_X(\mathbf{x})} \text{ remove}\end{aligned}$$

- Solution = the most likely class according to posterior
- This involves a likelihood which depends on the model you choose
- This involves a prior term which is subjective

Let us Plug-in Multi-dimensional Gaussian

Recall d -dimensional Gaussian.

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

Plug this into the discriminant function

$$\begin{aligned} i^* &= \operatorname{argmax}_i \log p_{\mathbf{X}|Y}(\mathbf{x} | i) + \log \pi_i \\ &= \operatorname{argmax}_i -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_i)^T \boldsymbol{\Sigma}_i^{-1} (\mathbf{x} - \boldsymbol{\mu}_i) - \cancel{\frac{d}{2} \log(2\pi)} - \frac{1}{2} \log |\boldsymbol{\Sigma}_i| + \log \pi_i \\ &= \operatorname{argmax}_i \underbrace{-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_i)^T \boldsymbol{\Sigma}_i^{-1} (\mathbf{x} - \boldsymbol{\mu}_i)}_{\text{depend on } \mathbf{x}} \underbrace{-\frac{1}{2} \log |\boldsymbol{\Sigma}_i| + \log \pi_i}_{\text{does not depend on } \mathbf{x}}. \end{aligned}$$

Special Case: 1D; Two classes

The decision rule is

$$i^* = \operatorname{argmax}_i \underbrace{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_i)^T \boldsymbol{\Sigma}_i^{-1}(\mathbf{x} - \boldsymbol{\mu}_i)}_{\text{depend on } \mathbf{x}} \underbrace{-\frac{1}{2} \log |\boldsymbol{\Sigma}_i| + \log \pi_i}_{\text{does not depend on } \mathbf{x}}.$$

Substitute $\boldsymbol{\Sigma}_i = \sigma^2$, and $\boldsymbol{\mu}_i = \mu_i$. Do two classes.

$$\begin{aligned} -\frac{(x-\mu_1)^2}{2\sigma^2} - \log \sigma + \log \pi_1 &\geq_{C_2}^{\geq C_1} -\frac{(x-\mu_2)^2}{2\sigma^2} - \log \sigma + \log \pi_2 \\ -\frac{(x-\mu_1)^2}{2\sigma^2} - \cancel{\log \sigma} + \log \pi_1 &\geq_{C_2}^{\geq C_1} -\frac{(x-\mu_2)^2}{2\sigma^2} - \cancel{\log \sigma} + \log \pi_2 \\ &\vdots \\ x &\geq_{C_2}^{\geq C_1} \underbrace{\frac{\mu_1 - \mu_2}{2} - \frac{\sigma^2}{\mu_1 - \mu_2} \log \frac{\pi_1}{\pi_2}}_{\text{does not depend on } x}. \end{aligned}$$

Connecting to Linear Discriminant Function

Recall: A hypothesis function is

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

If there are only two classes, then we can define

$$g(\mathbf{x}) = g_i(\mathbf{x}) - g_j(\mathbf{x}).$$

where the i -th discriminant function is

$$g_i(\mathbf{x}) = \log p_{\mathbf{X}|Y}(\mathbf{x}|i) + \log \pi_i.$$

- Class i if $g(\mathbf{x}) > 0 \iff g_i(\mathbf{x}) > g_j(\mathbf{x})$
- Class j if $g(\mathbf{x}) < 0 \iff g_i(\mathbf{x}) < g_j(\mathbf{x})$
- Either if $g(\mathbf{x}) = 0 \iff g_i(\mathbf{x}) = g_j(\mathbf{x})$

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:

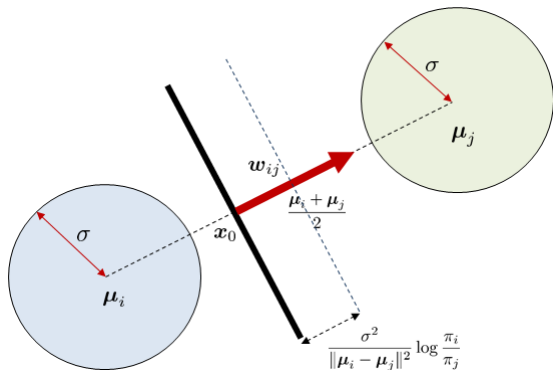
$$\begin{aligned} g(\mathbf{x}) &= \left(\frac{\boldsymbol{\mu}_i - \boldsymbol{\mu}_j}{\sigma^2} \right)^T \mathbf{x} - \left(\frac{\|\boldsymbol{\mu}_i\|^2}{2\sigma^2} - \frac{\|\boldsymbol{\mu}_j\|^2}{2\sigma^2} \right) + \log \frac{\pi_i}{\pi_j} \\ &= \left(\frac{\boldsymbol{\mu}_i - \boldsymbol{\mu}_j}{\sigma^2}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

Therefore, we have

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

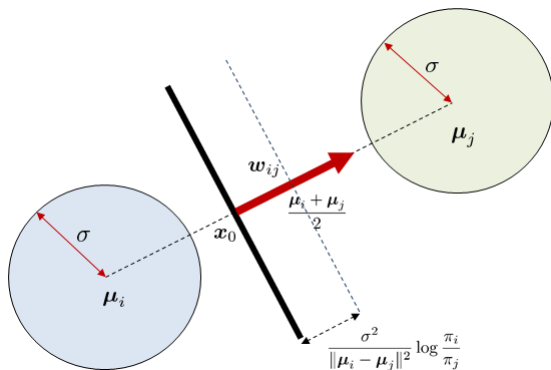
Case 1: $\Sigma_i = \sigma^2 I$: Geometry

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



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>