

Statistical Pattern Recognition

Lecture4

Bayesian Learning

Dr Zohreh Azimifar

School of Electrical and Computer Engineering

Shiraz University

Fall2014

Table of contents

- 1 Introduction
- 2 Generative Learning vs Discriminative Learning
 - Generative Learning and Discriminative Learning
- 3 Linear Discriminant Analysis
 - Gaussian Linear Discriminant Analysis
 - Boundary Decision for GLDA
 - Analysis of GLDA
- 4 Quadratic Discriminant Analysis
 - Quadratic Discriminant Analysis
 - Analysis of QDA
- 5 GLAD and QDA, Another point of view
 - GLAD and QDA, Another point of view
- 6 Naive Bayes
 - Naive Bayes
 - Naive Bayes: An Example
- 7 Lecture Summary
 - Summary

Introduction

- Classification based on the theory **Bayesian Learning**

$$P(y = 0|\mathbf{X}) \geq_{y=1}^{y=0} P(y = 1|\mathbf{X})$$

- Classification involves determining $P(y|\mathbf{X})$, from different perspectives.

Generative Learning and Discriminative Learning

① Discriminative Learning:

- Direct learning of $P(y|\mathbf{X})$.
- Modelling of **decision boundary**, to which side a new sample is assigned.
- Logistic and softmax regression are called discriminative learners.

② Generative Learning

- Explicit modelling of each class separately.
- Compare new sample with each class probability, based on **Bayesian rule**:

$$P(y|\mathbf{X}) = \frac{\overbrace{P(\mathbf{X}|y)}^{\text{Likelihood}} \overbrace{P(y)}^{\text{Prior}}}{\underbrace{P(\mathbf{X})}_{\text{normalizing factor}}}$$

Gaussian Linear Discriminant Analysis

- Model $P(y)$ and $P(\mathbf{X}|y)$ for each class y :

$$P(y) = \phi_1^{1\{y=1\}} \phi_2^{1\{y=2\}} \dots \phi_c^{1\{y=c\}}$$

$$P(\mathbf{X}|y = i) = \frac{1}{(\sqrt{2\pi})^n |\Sigma|^{\frac{1}{2}}} \exp\left(\frac{-1}{2}(\mathbf{X} - \boldsymbol{\mu}_i)^T \Sigma^{-1}(\mathbf{X} - \boldsymbol{\mu}_i)\right)$$

- Parameter set: $\boldsymbol{\theta} = \{\phi_1, \phi_2, \dots, \phi_c, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \dots, \boldsymbol{\mu}_c, \Sigma\}$

Gaussian Linear Discriminant Analysis

- Estimate parameters: $\theta = \{\phi_1, \phi_2, \dots, \phi_c, \mu_1, \mu_2, \dots, \mu_c, \Sigma\}$

$$l(\theta) = \log \prod_{j=1}^m P(\mathbf{x}^{(j)}, y^{(j)}) = \log \prod_{j=1}^m P(\mathbf{x}^{(j)} | P(y^{(j)})) P(y^{(j)})$$

Gaussian Linear Discriminant Analysis

- Estimate parameters: $\theta = \{\phi_1, \phi_2, \dots, \phi_c, \mu_1, \mu_2, \dots, \mu_c, \Sigma\}$

$$l(\theta) = \log \prod_{j=1}^m P(\mathbf{x}^{(j)}, y^{(j)}) = \log \prod_{j=1}^m P(\mathbf{x}^{(j)} | P(y^{(j)})) P(y^{(j)})$$

- Take partial derivative in terms of each individual parameter:

$$\phi_i^{MLE} = \frac{\sum_{j=1}^m 1\{y^{(j)} = i\}}{m}$$

Gaussian Linear Discriminant Analysis

- Estimate parameters: $\theta = \{\phi_1, \phi_2, \dots, \phi_c, \mu_1, \mu_2, \dots, \mu_c, \Sigma\}$

$$l(\theta) = \log \prod_{j=1}^m P(\mathbf{x}^{(j)}, y^{(j)}) = \log \prod_{j=1}^m P(\mathbf{x}^{(j)} | P(y^{(j)})) P(y^{(j)})$$

- Take partial derivative in terms of each individual parameter:

$$\phi_i^{MLE} = \frac{\sum_{j=1}^m 1\{y^{(j)} = i\}}{m}$$

$$\mu_i^{MLE} = \frac{\sum_{j=1}^m 1\{y^{(j)} = i\} \mathbf{x}^{(j)}}{\sum_{j=1}^m 1\{y^{(j)} = i\}}$$

Gaussian Linear Discriminant Analysis

- Estimate parameters: $\theta = \{\phi_1, \phi_2, \dots, \phi_c, \mu_1, \mu_2, \dots, \mu_c, \Sigma\}$

$$l(\theta) = \log \prod_{j=1}^m P(\mathbf{x}^{(j)}, y^{(j)}) = \log \prod_{j=1}^m P(\mathbf{x}^{(j)} | P(y^{(j)})) P(y^{(j)})$$

- Take partial derivative in terms of each individual parameter:

$$\phi_i^{MLE} = \frac{\sum_{j=1}^m 1\{y^{(j)} = i\}}{m}$$

$$\mu_i^{MLE} = \frac{\sum_{j=1}^m 1\{y^{(j)} = i\} \mathbf{x}^{(j)}}{\sum_{j=1}^m 1\{y^{(j)} = i\}}$$

$$\Sigma^{MLE} = \frac{1}{m} \sum_{j=1}^m (\mathbf{x}^{(j)} - \mu_{y^{(j)}})(\mathbf{x}^{(j)} - \mu_{y^{(j)}})^T$$

Gaussian Linear Discriminant Analysis

- Determine class label of a new sample \mathbf{X}^{new} :

$$\begin{aligned}y^{new} &= \operatorname{argmax}_y P(y|\mathbf{X}) \\&= \operatorname{argmax}_y \frac{P(\mathbf{X}|y)P(y)}{P(\mathbf{X})} \\&= \operatorname{argmax}_y P(\mathbf{X}|y)P(y)\end{aligned}$$

- Note that $P(\mathbf{X}|y)$ is a class dependent density.

Boundary Decision for GLDA

- Decision boundary is a line, a plane, or a hyper-plane. Why?

$$P(y = i|\mathbf{X}) = P(y = j|\mathbf{X})$$
$$\frac{P(\mathbf{X}|y = i)P(y = i)}{P(\mathbf{X})} = \frac{P(\mathbf{X}|y = j)P(y = j)}{P(\mathbf{X})}$$

Boundary Decision for GLDA

- Decision boundary is a line, a plane, or a hyper-plane. Why?

$$P(y = i|\mathbf{X}) = P(y = j|\mathbf{X})$$
$$\frac{P(\mathbf{X}|y = i)P(y = i)}{P(\mathbf{X})} = \frac{P(\mathbf{X}|y = j)P(y = j)}{P(\mathbf{X})}$$

$$P(\mathbf{X}|y = i)P(y = i) = P(\mathbf{X}|y = j)P(y = j)$$

Boundary Decision for GLDA

- Decision boundary is a line, a plane, or a hyper-plane. Why?

$$P(y = i|\mathbf{X}) = P(y = j|\mathbf{X})$$

$$\frac{P(\mathbf{X}|y = i)P(y = i)}{P(\mathbf{X})} = \frac{P(\mathbf{X}|y = j)P(y = j)}{P(\mathbf{X})}$$

$$P(\mathbf{X}|y = i)P(y = i) = P(\mathbf{X}|y = j)P(y = j)$$

$$\frac{1}{(2\pi)^{\frac{n}{2}}|\Sigma|^{\frac{1}{2}}} \exp\left(-\frac{1}{2}(\mathbf{X} - \mu_i)^T \Sigma^{-1}(\mathbf{X} - \mu_i)\right)P(y = i)$$

$$= \frac{1}{(2\pi)^{\frac{n}{2}}|\Sigma|^{\frac{1}{2}}} \exp\left(-\frac{1}{2}(\mathbf{X} - \mu_j)^T \Sigma^{-1}(\mathbf{X} - \mu_j)\right)P(y = j)$$

Boundary Decision for GLDA

- Decision boundary cont'ed:

$$\exp\left(\frac{-1}{2}(\mathbf{X} - \boldsymbol{\mu}_i)^T \boldsymbol{\Sigma}^{-1}(\mathbf{X} - \boldsymbol{\mu}_i)\right)P(y = i) = \exp\left(\frac{-1}{2}(\mathbf{X} - \boldsymbol{\mu}_j)^T \boldsymbol{\Sigma}^{-1}(\mathbf{X} - \boldsymbol{\mu}_j)\right)P(y = j)$$

Boundary Decision for GLDA

- Decision boundary cont'ed:

$$\exp\left(\frac{-1}{2}(\mathbf{X} - \boldsymbol{\mu}_i)^T \boldsymbol{\Sigma}^{-1}(\mathbf{X} - \boldsymbol{\mu}_i)\right)P(y = i) = \exp\left(\frac{-1}{2}(\mathbf{X} - \boldsymbol{\mu}_j)^T \boldsymbol{\Sigma}^{-1}(\mathbf{X} - \boldsymbol{\mu}_j)\right)P(y = j)$$

$$\frac{-1}{2}(\mathbf{X} - \boldsymbol{\mu}_i)^T \boldsymbol{\Sigma}^{-1}(\mathbf{X} - \boldsymbol{\mu}_i) + \log P(y = i) = \frac{-1}{2}(\mathbf{X} - \boldsymbol{\mu}_j)^T \boldsymbol{\Sigma}^{-1}(\mathbf{X} - \boldsymbol{\mu}_j) + \log P(y = j)$$

Boundary Decision for GLDA

- Decision boundary cont'ed:

$$\exp\left(\frac{-1}{2}(\mathbf{X} - \boldsymbol{\mu}_i)^T \boldsymbol{\Sigma}^{-1}(\mathbf{X} - \boldsymbol{\mu}_i)\right)P(y = i) = \exp\left(\frac{-1}{2}(\mathbf{X} - \boldsymbol{\mu}_j)^T \boldsymbol{\Sigma}^{-1}(\mathbf{X} - \boldsymbol{\mu}_j)\right)P(y = j)$$

$$\frac{-1}{2}(\mathbf{X} - \boldsymbol{\mu}_i)^T \boldsymbol{\Sigma}^{-1}(\mathbf{X} - \boldsymbol{\mu}_i) + \log P(y = i) = \frac{-1}{2}(\mathbf{X} - \boldsymbol{\mu}_j)^T \boldsymbol{\Sigma}^{-1}(\mathbf{X} - \boldsymbol{\mu}_j) + \log P(y = j)$$

$$\begin{aligned} \Rightarrow \log \frac{P(y = i)}{P(y = j)} - \frac{1}{2}[\mathbf{X}^T \boldsymbol{\Sigma}^{-1} \mathbf{X} - 2\mathbf{X}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_i + \boldsymbol{\mu}_i^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_i] \\ + \frac{1}{2}[\mathbf{X}^T \boldsymbol{\Sigma}^{-1} \mathbf{X} - 2\mathbf{X}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_j + \boldsymbol{\mu}_j^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_j] = 0 \end{aligned}$$

Boundary Decision for GLDA

- Decision boundary cont'ed:

$$\exp\left(\frac{-1}{2}(\mathbf{X} - \mu_i)^T \Sigma^{-1}(\mathbf{X} - \mu_i)\right)P(y = i) = \exp\left(\frac{-1}{2}(\mathbf{X} - \mu_j)^T \Sigma^{-1}(\mathbf{X} - \mu_j)\right)P(y = j)$$

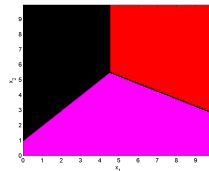
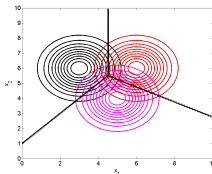
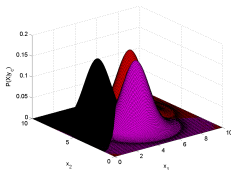
$$\frac{-1}{2}(\mathbf{X} - \mu_i)^T \Sigma^{-1}(\mathbf{X} - \mu_i) + \log P(y = i) = \frac{-1}{2}(\mathbf{X} - \mu_j)^T \Sigma^{-1}(\mathbf{X} - \mu_j) + \log P(y = j)$$

$$\Rightarrow \log \frac{P(y = i)}{P(y = j)} - \frac{1}{2}[\mathbf{X}^T \Sigma^{-1} \mathbf{X} - 2\mathbf{X}^T \Sigma^{-1} \mu_i + \mu_i^T \Sigma^{-1} \mu_i] \\ + \frac{1}{2}[\mathbf{X}^T \Sigma^{-1} \mathbf{X} - 2\mathbf{X}^T \Sigma^{-1} \mu_j + \mu_j^T \Sigma^{-1} \mu_j] = 0$$

$$\Rightarrow \underbrace{\mathbf{X}^T \Sigma^{-1}(\mu_i - \mu_j)}_{a\mathbf{X}} + \underbrace{\frac{1}{2}\mu_i^T \Sigma^{-1} \mu_i - \frac{1}{2}\mu_j^T \Sigma^{-1} \mu_j + \log \frac{P(y = i)}{P(y = j)}}_b = 0$$

Analysis of GLDA when $\Sigma = \sigma^2 I$

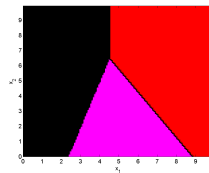
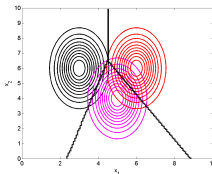
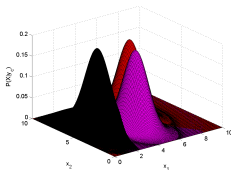
- Classes are of identical distribution, but different means.
- Cross-section of classes distribution is spherical.
- Decision boundary is linear.
- Called classifier with nearest Euclidean distance to the class mean, when?



$$\Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

Analysis of GLDA when Σ is not identity

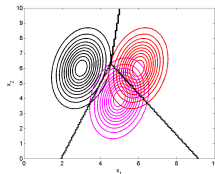
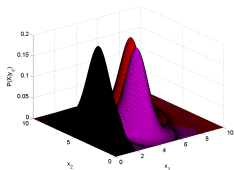
- Classes are of identical distribution, but different means.
- Cross-section of classes distribution is ellipsoidal.
- Decision boundary is linear.



$$\Sigma = \begin{bmatrix} 0.5 & 0 \\ 0 & 1.5 \end{bmatrix}$$

Analysis of GLDA with arbitrary Σ

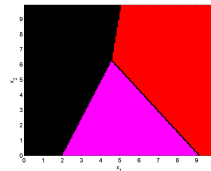
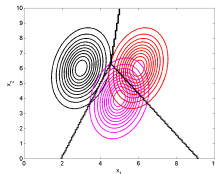
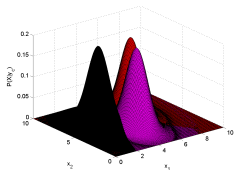
- Classes are of identical distribution, but different means.
- Cross-section of classes distribution is ellipsoidal. Linear decision boundary.
- Classes are aligned with direction of covariance eigenvectors.



$$\Sigma = \begin{bmatrix} 0.5 & 0.2 \\ 0.2 & 1.5 \end{bmatrix}$$

Analysis of GLDA with arbitrary Σ

- Classes are of identical distribution, but different means.
- Cross-section of classes distribution is ellipsoidal. Linear decision boundary.
- Classes are aligned with direction of covariance eigenvectors.



$$\Sigma = \begin{bmatrix} 0.5 & 0.2 \\ 0.2 & 1.5 \end{bmatrix}$$

- Called classifier with nearest Mahalanobis distance to the class mean, when? $\text{Dist}(\mathbf{X}, \mu_i) = (\mathbf{X} - \mu_i)^T \Sigma^{-1} (\mathbf{X} - \mu_i)$

Quadratic Discriminant Analysis

- Another generative learning model; a Bayesian classifier
- Here, classes are of multinomial distribution, and likelihoods are multivariate Gaussian with separate covariance Σ_i .
- Decision boundary becomes non-linear, Why?

$$P(\mathbf{X}|y = i) = \frac{1}{(\sqrt{2\pi})^n |\Sigma_i|^{\frac{1}{2}}} \exp\left(-\frac{1}{2}(\mathbf{X} - \mu_i)^T \Sigma_i^{-1}(\mathbf{X} - \mu_i)\right)$$

Quadratic Discriminant Analysis

Decision boundary is parabolic:

$$\begin{aligned}P(y = i|\mathbf{X}) &= P(y = j|\mathbf{X}) \\ \frac{P(\mathbf{X}|y = i)P(y = i)}{P(\mathbf{X})} &= \frac{P(\mathbf{X}|y = j)P(y = j)}{P(\mathbf{X})} \\ P(\mathbf{X}|y = i)P(y = i) &= P(\mathbf{X}|y = j)P(y = j)\end{aligned}$$

Quadratic Discriminant Analysis

Decision boundary is parabolic:

$$\begin{aligned}P(y = i|\mathbf{X}) &= P(y = j|\mathbf{X}) \\ \frac{P(\mathbf{X}|y = i)P(y = i)}{P(\mathbf{X})} &= \frac{P(\mathbf{X}|y = j)P(y = j)}{P(\mathbf{X})} \\ P(\mathbf{X}|y = i)P(y = i) &= P(\mathbf{X}|y = j)P(y = j)\end{aligned}$$

$$\begin{aligned}& \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma_i|^{\frac{1}{2}}} \exp\left(\frac{-1}{2}(\mathbf{X} - \mu_i)^T \Sigma_i^{-1}(\mathbf{X} - \mu_i)\right) P(y = i) \\ &= \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma_j|^{\frac{1}{2}}} \exp\left(\frac{-1}{2}(\mathbf{X} - \mu_j)^T \Sigma_j^{-1}(\mathbf{X} - \mu_j)\right) P(y = j)\end{aligned}$$

Quadratic Discriminant Analysis

Decision boundary cont'ed:

$$\begin{aligned} & \frac{1}{|\Sigma_i|^{\frac{1}{2}}} \exp\left(\frac{-1}{2}(\mathbf{X} - \mu_i)^T \Sigma_i^{-1}(\mathbf{X} - \mu_i)\right) P(y = i) \\ &= \frac{1}{|\Sigma_j|^{\frac{1}{2}}} \exp\left(\frac{-1}{2}(\mathbf{X} - \mu_j)^T \Sigma_j^{-1}(\mathbf{X} - \mu_j)\right) P(y = j) \end{aligned}$$

Quadratic Discriminant Analysis

Decision boundary cont'ded:

$$\begin{aligned} & \frac{1}{|\Sigma_i|^{\frac{1}{2}}} \exp\left(\frac{-1}{2}(\mathbf{X} - \mu_i)^T \Sigma_i^{-1}(\mathbf{X} - \mu_i)\right) P(y = i) \\ &= \frac{1}{|\Sigma_j|^{\frac{1}{2}}} \exp\left(\frac{-1}{2}(\mathbf{X} - \mu_j)^T \Sigma_j^{-1}(\mathbf{X} - \mu_j)\right) P(y = j) \\ & \frac{-1}{2} \log |\Sigma_i| - \frac{1}{2}(\mathbf{X} - \mu_i)^T \Sigma_i^{-1}(\mathbf{X} - \mu_i) + \log P(y = i) \\ &= \frac{-1}{2} \log |\Sigma_j| - \frac{1}{2}(\mathbf{X} - \mu_j)^T \Sigma_j^{-1}(\mathbf{X} - \mu_j) + \log P(y = j) \end{aligned}$$

Quadratic Discriminant Analysis

Decision boundary cont'ed:

$$\begin{aligned}
 & \frac{1}{|\Sigma_i|^{\frac{1}{2}}} \exp\left(\frac{-1}{2}(\mathbf{X} - \mu_i)^T \Sigma_i^{-1}(\mathbf{X} - \mu_i)\right) P(y = i) \\
 &= \frac{1}{|\Sigma_j|^{\frac{1}{2}}} \exp\left(\frac{-1}{2}(\mathbf{X} - \mu_j)^T \Sigma_j^{-1}(\mathbf{X} - \mu_j)\right) P(y = j) \\
 \\
 & \frac{-1}{2} \log |\Sigma_i| - \frac{1}{2}(\mathbf{X} - \mu_i)^T \Sigma_i^{-1}(\mathbf{X} - \mu_i) + \log P(y = i) \\
 &= \frac{-1}{2} \log |\Sigma_j| - \frac{1}{2}(\mathbf{X} - \mu_j)^T \Sigma_j^{-1}(\mathbf{X} - \mu_j) + \log P(y = j) \\
 \\
 & \Rightarrow \log \frac{P(y = i)}{P(y = j)} - \frac{1}{2} \log \frac{|\Sigma_i|}{|\Sigma_j|} - \frac{1}{2} [\mathbf{X}^T \Sigma_i^{-1} \mathbf{X} + \mu_i^T \Sigma_i^{-1} \mu_i - \\
 & \quad 2\mathbf{X}^T \Sigma_i^{-1} \mu_i - \mathbf{X}^T \Sigma_j^{-1} \mathbf{X} - \mu_j^T \Sigma_j^{-1} \mu_j + 2\mathbf{X}^T \Sigma_j^{-1} \mu_j] = 0
 \end{aligned}$$

Quadratic Discriminant Analysis

Decision boundary cont'ed:

$$\log \frac{P(y=i)}{P(y=j)} - \frac{1}{2} \log \frac{|\Sigma_i|}{|\Sigma_j|} - \frac{1}{2} [\mathbf{X}^T (\Sigma_i^{-1} - \Sigma_j^{-1}) \mathbf{X} + \mu_i^T \Sigma_i^{-1} \mu_i - \mu_j^T \Sigma_j^{-1} \mu_j - 2 \mathbf{X}^T (\Sigma_i^{-1} \mu_i - \Sigma_j^{-1} \mu_j)] = 0$$

Quadratic Discriminant Analysis

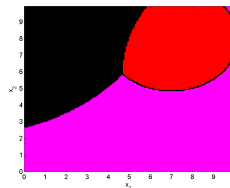
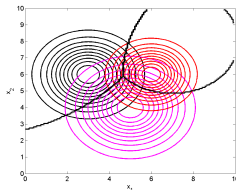
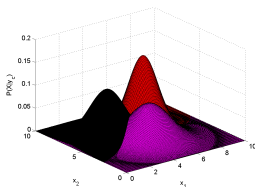
Decision boundary cont'ed:

$$\log \frac{P(y=i)}{P(y=j)} - \frac{1}{2} \log \frac{|\Sigma_i|}{|\Sigma_j|} - \frac{1}{2} [\mathbf{X}^T (\Sigma_i^{-1} - \Sigma_j^{-1}) \mathbf{X} + \mu_i^T \Sigma_i^{-1} \mu_i - \mu_j^T \Sigma_j^{-1} \mu_j - 2 \mathbf{X}^T (\Sigma_i^{-1} \mu_i - \Sigma_j^{-1} \mu_j)] = 0$$

$$\Rightarrow \mathbf{X}^T a \mathbf{X} + b^T \mathbf{X} + c = 0$$

Analysis of QDA when $\Sigma = \sigma_i^2 I$

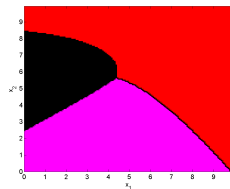
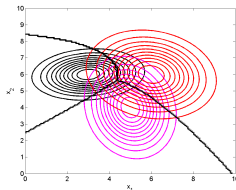
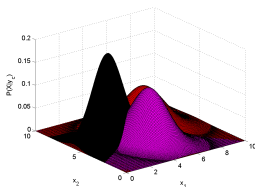
- Classes are of different distributions and different means.
- Cross-sections of classes distribution are spherical but of different sizes.



$$\Sigma_1 = \begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \end{bmatrix}, \quad \Sigma_2 = \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix}, \quad \Sigma_3 = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

Analysis of QDA with arbitrary $\Sigma_i \neq \Sigma_j$

- Classes are of different distributions and different means.
- Cross-sections of classes distribution are ellipsoidal and of different sizes.



$$\Sigma_1 = \begin{bmatrix} 1.5 & 0.1 \\ 0.1 & 0.5 \end{bmatrix}, \Sigma_2 = \begin{bmatrix} 1 & -0.2 \\ -0.2 & 2 \end{bmatrix}, \Sigma_3 = \begin{bmatrix} 2 & -0.25 \\ -0.25 & 1.5 \end{bmatrix}$$

GLAD and QDA, Another point of view

-
-
-

Naive Bayes

-
-
-

Naive Bayes: An Example

Summary

