

Lecture 5

Gaussian Models - Part 1

Luigi Freda

ALCOR Lab
DIAG
University of Rome "La Sapienza"

November 29, 2016

1 Basics

- Multivariate Gaussian

2 MLE for an MVN

- Theorem

3 Gaussian Discriminant Analysis

- Generative Classifiers
- Gaussian Discriminant Analysis (GDA)
- Quadratic Discriminant Analysis (QDA)
- Linear Discriminant Analysis (LDA)
- MLE for Gaussian Discriminant Analysis
- Diagonal LDA
- Bayesian Procedure

1 Basics

- Multivariate Gaussian

2 MLE for an MVN

- Theorem

3 Gaussian Discriminant Analysis

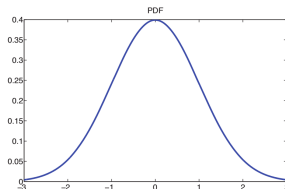
- Generative Classifiers
- Gaussian Discriminant Analysis (GDA)
- Quadratic Discriminant Analysis (QDA)
- Linear Discriminant Analysis (LDA)
- MLE for Gaussian Discriminant Analysis
- Diagonal LDA
- Bayesian Procedure

Univariate Gaussian (Normal) Distribution

- X is a continuous RV with values $x \in \mathbb{R}$
- $X \sim \mathcal{N}(\mu, \sigma^2)$, i.e. X has a **Gaussian distribution** or **normal distribution**

$$\mathcal{N}(x|\mu, \sigma^2) \triangleq \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(x-\mu)^2} \quad (= P_X(X=x))$$

- mean $\mathbb{E}[X] = \mu$
- mode μ
- variance $\text{var}[X] = \sigma^2$
- precision $\lambda = \frac{1}{\sigma^2}$
- $(\mu - 2\sigma, \mu + 2\sigma)$ is the approx 95% interval
- $(\mu - 3\sigma, \mu + 3\sigma)$ is the approx. 99.7% interval

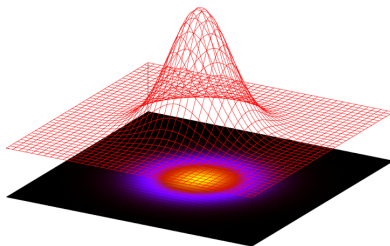


Multivariate Gaussian (Normal) Distribution

- \mathbf{X} is a continuous RV with values $\mathbf{x} \in \mathbb{R}^D$
- $\mathbf{X} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, i.e. \mathbf{X} has a **Multivariate Normal** distribution (MVN) or **multivariate Gaussian**

$$\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) \triangleq \frac{1}{(2\pi)^{D/2}|\boldsymbol{\Sigma}|^{1/2}} \exp\left[-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right]$$

- mean: $\mathbb{E}[\mathbf{x}] = \boldsymbol{\mu}$
- mode: $\boldsymbol{\mu}$
- covariance matrix: $\text{cov}[\mathbf{x}] = \boldsymbol{\Sigma} \in \mathbb{R}^{D \times D}$ where $\boldsymbol{\Sigma} = \boldsymbol{\Sigma}^T$ and $\boldsymbol{\Sigma} \geq 0$
- precision matrix: $\boldsymbol{\Lambda} \triangleq \boldsymbol{\Sigma}^{-1}$
- spherical isotropic covariance with $\boldsymbol{\Sigma} = \sigma^2 \mathbf{I}_D$



1 Basics

- Multivariate Gaussian

2 MLE for an MVN

- Theorem

3 Gaussian Discriminant Analysis

- Generative Classifiers
- Gaussian Discriminant Analysis (GDA)
- Quadratic Discriminant Analysis (QDA)
- Linear Discriminant Analysis (LDA)
- MLE for Gaussian Discriminant Analysis
- Diagonal LDA
- Bayesian Procedure

MLE for an MVN

Theorem

Theorem 1

If we have N iid samples $\mathbf{x}_i \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, then the MLE for the parameters is given by

$$\textcircled{1} \quad \boldsymbol{\mu}_{MLE} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i \triangleq \bar{\mathbf{x}}$$

$$\textcircled{2} \quad \boldsymbol{\Sigma}_{MLE} = \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^T = \frac{1}{N} \left(\sum_{i=1}^N \mathbf{x}_i \mathbf{x}_i^T \right) - \bar{\mathbf{x}} \bar{\mathbf{x}}^T$$

- this theorem states the MLE parameter estimates for an MVN are just the **empirical mean** and the **empirical covariance**
- in the **univariate case**, one has

$$\hat{\mu} = \frac{1}{N} \sum_{i=1}^N x_i \triangleq \bar{x}$$

$$\hat{\sigma}^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})(x_i - \bar{x})^T = \frac{1}{N} \left(\sum_{i=1}^N x_i x_i^T \right) - \bar{x}^2$$

MLE for an MVN

Theorem

proof sketch

- in order to find the MLE one should maximize the log-likelihood of the dataset
- given that $\mathbf{x}_i \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$

$$p(\mathcal{D}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \prod_i \mathcal{N}(\mathbf{x}_i|\boldsymbol{\mu}, \boldsymbol{\Sigma})$$

- the log-likelihood (dropping additive constants) is

$$l(\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \log p(\mathcal{D}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{N}{2} \log |\boldsymbol{\Lambda}| - \frac{1}{2} \sum_i (\mathbf{x}_i - \boldsymbol{\mu}) \boldsymbol{\Lambda} (\mathbf{x}_i - \boldsymbol{\mu})^T + \text{const}$$

- the MLE estimates can be obtained by maximizing $l(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ w.r.t. $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}$

homework: continue the proof for the univariate case

1 Basics

- Multivariate Gaussian

2 MLE for an MVN

- Theorem

3 Gaussian Discriminant Analysis

- **Generative Classifiers**
 - Gaussian Discriminant Analysis (GDA)
 - Quadratic Discriminant Analysis (QDA)
 - Linear Discriminant Analysis (LDA)
 - MLE for Gaussian Discriminant Analysis
 - Diagonal LDA
 - Bayesian Procedure

Generative Classifiers

probabilistic classifier

- we are given a dataset $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$
- the goal is to compute the **class posterior** $p(y = c|\mathbf{x})$ which models the mapping $y = f(\mathbf{x})$

generative classifiers

- $p(y = c|\mathbf{x})$ is computed starting from the **class-conditional density** $p(\mathbf{x}|y = c, \theta)$ and the **class prior** $p(y = c|\theta)$ given that

$$p(y = c|\mathbf{x}, \theta) \propto p(\mathbf{x}|y = c, \theta)p(y = c|\theta) \quad (= p(y = c, \mathbf{x}|\theta))$$

- this is called a **generative classifier** since it specifies how to generate the feature vector \mathbf{x} for each class $y = c$ (by using $p(\mathbf{x}|y = c, \theta)$)
- the model is usually fit by maximizing the joint log-likelihood, i.e. one computes $\theta^* = \arg \max_{\theta} \sum_i \log p(y_i, \mathbf{x}_i|\theta)$

discriminative classifiers

- the model $p(y = c|\mathbf{x})$ is directly fit to the data
- the model is usually fit by maximizing the conditional log-likelihood, i.e. one computes $\theta^* = \arg \max_{\theta} \sum_i \log p(y_i|\mathbf{x}_i, \theta)$

1 Basics

- Multivariate Gaussian

2 MLE for an MVN

- Theorem

3 Gaussian Discriminant Analysis

- Generative Classifiers
- **Gaussian Discriminant Analysis (GDA)**
- Quadratic Discriminant Analysis (QDA)
- Linear Discriminant Analysis (LDA)
- MLE for Gaussian Discriminant Analysis
- Diagonal LDA
- Bayesian Procedure

Gaussian Discriminant Analysis

GDA

- we can use the MVN for defining the class conditional densities in a generative classifier

$$p(\mathbf{x}|y = c, \boldsymbol{\theta}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_c, \boldsymbol{\Sigma}_c) \quad \text{for } c \in \{1, \dots, C\}$$

- this means the samples of each class c are characterized by a normal distribution
- this model is called **Gaussian Discriminative Analysis** (GDA) but it is a **generative classifier** (not discriminative)
- in the case $\boldsymbol{\Sigma}_c$ is diagonal for each c , this model is equivalent to a Naive Bayes Classifier (NBC) since

$$p(\mathbf{x}|y = c, \boldsymbol{\theta}) = \prod_{j=1}^D \mathcal{N}(x_j|\mu_{jc}, \sigma_{jc}^2) \quad \text{for } c \in \{1, \dots, C\}$$

- once the model is fit to the data, we can classify a feature vector by using the decision rule

$$\hat{y}(\mathbf{x}) = \operatorname{argmax}_c \log p(y = c|\mathbf{x}, \boldsymbol{\theta}) = \operatorname{argmax}_c \left[\log p(y = c|\boldsymbol{\pi}) + \log p(\mathbf{x}|y = c, \boldsymbol{\theta}_c) \right]$$

Gaussian Discriminant Analysis

GDA

- decision rule

$$\hat{y}(\mathbf{x}) = \underset{c}{\operatorname{argmax}} \left[\log p(y = c | \boldsymbol{\pi}) + \log p(\mathbf{x} | y = c, \boldsymbol{\theta}_c) \right]$$

- given that $y \sim \text{Cat}(\boldsymbol{\pi})$ and $\mathbf{x} | (y = c) \sim \mathcal{N}(\boldsymbol{\mu}_c, \boldsymbol{\Sigma}_c)$ the decision rule becomes (dropping additive constants)

$$\hat{y}(\mathbf{x}) = \underset{c}{\operatorname{argmin}} \left[-\log \pi_c + \frac{1}{2} \log |\boldsymbol{\Sigma}_c| + \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_c)^T \boldsymbol{\Sigma}_c^{-1} (\mathbf{x} - \boldsymbol{\mu}_c) \right]$$

which can be thought as a **nearest centroid classifier**

- in fact, with an uniform prior and $\boldsymbol{\Sigma}_c = \boldsymbol{\Sigma}$

$$\hat{y}(\mathbf{x}) = \underset{c}{\operatorname{argmin}} (\mathbf{x} - \boldsymbol{\mu}_c)^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}_c) = \underset{c}{\operatorname{argmin}} \|\mathbf{x} - \boldsymbol{\mu}_c\|_{\boldsymbol{\Sigma}}^2$$

- in this case, we select the class c whose center $\boldsymbol{\mu}_c$ is closest to \mathbf{x} (using the **Mahalanobis distance** $\|\mathbf{x} - \boldsymbol{\mu}_c\|_{\boldsymbol{\Sigma}}$)

Mahalanobis Distance

- the covariance matrix Σ can be diagonalized since it is a symmetric real matrix

$$\Sigma = \mathbf{U}\mathbf{D}\mathbf{U}^T = \sum_{i=1}^D \lambda_i \mathbf{u}_i \mathbf{u}_i^T$$

where $\mathbf{U} = [\mathbf{u}_1, \dots, \mathbf{u}_D]$ is an orthonormal matrix of eigenvectors (i.e. $\mathbf{U}^T \mathbf{U} = \mathbf{I}$) and λ_i are the corresponding eigenvalues ($\lambda_i \geq 0$ since $\Sigma \geq 0$)

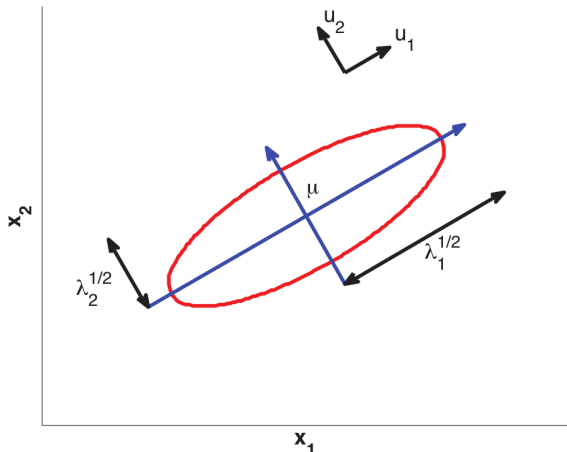
- one has immediately $\Sigma^{-1} = \mathbf{U}\mathbf{D}^{-1}\mathbf{U}^T = \sum_{i=1}^D \frac{1}{\lambda_i} \mathbf{u}_i \mathbf{u}_i^T$
- the **Mahalanobis distance** is defined as $\|\mathbf{x} - \boldsymbol{\mu}\|_{\Sigma} \triangleq \left((\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right)^{1/2}$
- one can rewrite

$$\begin{aligned} (\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}) &= (\mathbf{x} - \boldsymbol{\mu})^T \left(\sum_{i=1}^D \frac{1}{\lambda_i} \mathbf{u}_i \mathbf{u}_i^T \right) (\mathbf{x} - \boldsymbol{\mu}) = \\ &= \sum_{i=1}^D \frac{1}{\lambda_i} (\mathbf{x} - \boldsymbol{\mu})^T \mathbf{u}_i \mathbf{u}_i^T (\mathbf{x} - \boldsymbol{\mu}) = \sum_{i=1}^D \frac{y_i^2}{\lambda_i} \end{aligned}$$

where $y_i \triangleq \mathbf{u}_i^T (\mathbf{x} - \boldsymbol{\mu})$ (or equivalently $\mathbf{y} \triangleq \mathbf{U}^T (\mathbf{x} - \boldsymbol{\mu})$)

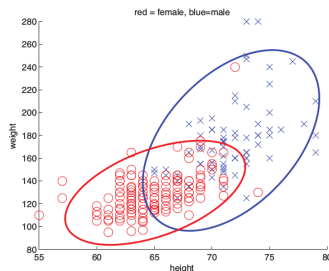
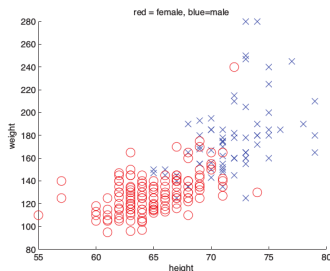
Mahalanobis Distance

- $\Sigma = \mathbf{U}\mathbf{D}\mathbf{U}^T = \sum_{i=1}^D \lambda_i \mathbf{u}_i \mathbf{u}_i^T$
- $(\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}) = \sum_{i=1}^D \frac{y_i^2}{\lambda_i}$ (where $\mathbf{y} \triangleq \mathbf{U}^T (\mathbf{x} - \boldsymbol{\mu})$)
- (1) center w.r.t. $\boldsymbol{\mu}$ (2) rotate by \mathbf{U}^T (3) get a norm weighted by the $\frac{1}{\lambda_i}$



Gaussian Discriminant Analysis

GDA



- *left*: height/weight data for the two classes male/female
- *right*: visualization of 2D Gaussian fit to each class
- we can see that the features are correlated (tall people tend to weigh more)

1 Basics

- Multivariate Gaussian

2 MLE for an MVN

- Theorem

3 Gaussian Discriminant Analysis

- Generative Classifiers
- Gaussian Discriminant Analysis (GDA)
- **Quadratic Discriminant Analysis (QDA)**
- Linear Discriminant Analysis (LDA)
- MLE for Gaussian Discriminant Analysis
- Diagonal LDA
- Bayesian Procedure

Quadratic Discriminant Analysis

QDA

- the **complete class posterior** with Gaussian densities is

$$p(y = c | \mathbf{x}, \theta) = \frac{\pi_c |2\pi \Sigma_c|^{-1/2} \exp[-\frac{1}{2}(\mathbf{x} - \mu_c)^T \Sigma_c^{-1}(\mathbf{x} - \mu_c)]}{\sum_{c'} \pi_{c'} |2\pi \Sigma_{c'}|^{-1/2} \exp[-\frac{1}{2}(\mathbf{x} - \mu_{c'})^T \Sigma_{c'}^{-1}(\mathbf{x} - \mu_{c'})]}$$

- the **quadratic decision boundaries** can be found by imposing

$$p(y = c' | \mathbf{x}, \theta) = p(y = c'' | \mathbf{x}, \theta)$$

or equivalently

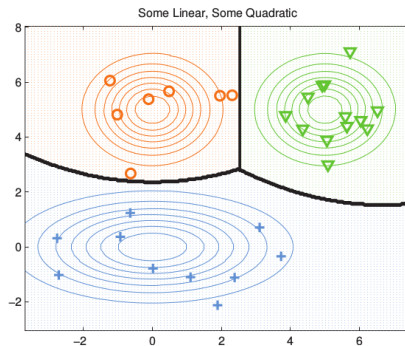
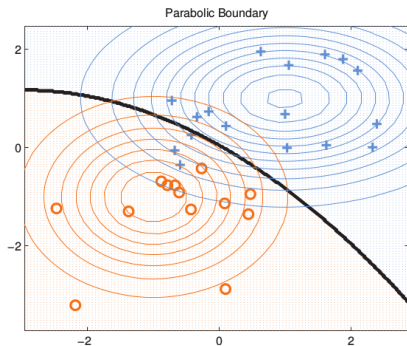
$$\log p(y = c' | \mathbf{x}, \theta) = \log p(y = c'' | \mathbf{x}, \theta)$$

for each pair of "adjacent" classes (c' , c''), which results in the quadratic equation

$$-\frac{1}{2}(\mathbf{x} - \mu_{c'})^T \Sigma_{c'}^{-1}(\mathbf{x} - \mu_{c'}) = -\frac{1}{2}(\mathbf{x} - \mu_{c''})^T \Sigma_{c''}^{-1}(\mathbf{x} - \mu_{c''}) + \text{constant}$$

Quadratic Discriminant Analysis

QDA



- *left*: dataset with 2 classes
- *right*: dataset with 3 classes

1 Basics

- Multivariate Gaussian

2 MLE for an MVN

- Theorem

3 Gaussian Discriminant Analysis

- Generative Classifiers
- Gaussian Discriminant Analysis (GDA)
- Quadratic Discriminant Analysis (QDA)
- **Linear Discriminant Analysis (LDA)**
- MLE for Gaussian Discriminant Analysis
- Diagonal LDA
- Bayesian Procedure

Linear Discriminant Analysis

LDA

- we now consider the GDA in the special case $\Sigma_c = \Sigma$ for $c \in \{1, \dots, C\}$
- in this case we have

$$\begin{aligned} p(y = c | \mathbf{x}, \theta) &\propto \pi_c \exp \left[-\frac{1}{2} (\mathbf{x} - \mu_c)^T \Sigma^{-1} (\mathbf{x} - \mu_c) \right] = \\ &= \exp \left[-\frac{1}{2} \mathbf{x}^T \Sigma^{-1} \mathbf{x} \right] \exp \left[\mu_c^T \Sigma^{-1} \mathbf{x} - \frac{1}{2} \mu_c^T \Sigma^{-1} \mu_c + \log \pi_c \right] \end{aligned}$$

- note that the quadratic term $-\frac{1}{2} \mathbf{x}^T \Sigma^{-1} \mathbf{x}$ is independent of c and it will cancel out in the numerator and denominator of the complete class posterior equation
- we define

$$\begin{aligned} \gamma_c &\triangleq -\frac{1}{2} \mu_c^T \Sigma^{-1} \mu_c + \log \pi_c \\ \beta_c &\triangleq \Sigma^{-1} \mu_c \end{aligned}$$

- we can rewrite

$$p(y = c | \mathbf{x}, \theta) = \frac{e^{\beta_c^T \mathbf{x} + \gamma_c}}{\sum_{c'} e^{\beta_{c'}^T \mathbf{x} + \gamma_{c'}}}$$

Linear Discriminant Analysis

LDA

- we have

$$p(y = c | \mathbf{x}, \boldsymbol{\theta}) = \frac{e^{\beta_c^T \mathbf{x} + \gamma_c}}{\sum_{c'} e^{\beta_{c'}^T \mathbf{x} + \gamma_{c'}}} \triangleq S(\boldsymbol{\eta})_c$$

where $\boldsymbol{\eta} \triangleq [\beta_1^T \mathbf{x} + \gamma_1, \dots, \beta_C^T \mathbf{x} + \gamma_C]^T \in \mathbb{R}^C$ and the function $S(\boldsymbol{\eta})$ is the **softmax function** defined as

$$S(\boldsymbol{\eta}) \triangleq \left[\frac{e^{\eta_1}}{\sum_{c'} e^{\eta_{c'}}}, \dots, \frac{e^{\eta_C}}{\sum_{c'} e^{\eta_{c'}}} \right]^T$$

and $S(\boldsymbol{\eta})_c \in \mathbb{R}$ is just its c -th component

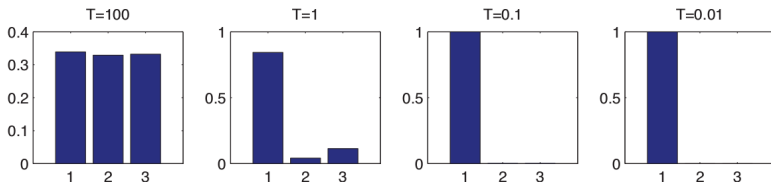
- the softmax function $S(\boldsymbol{\eta})$ is so-called since it acts a bit like the max function. To see this, divide each component η_c by a **temperature** T , then

$$S(\boldsymbol{\eta}/T)_c = \begin{cases} 1 & \text{if } c = \operatorname{argmax}_{c'} \eta_{c'} \\ 0 & \text{otherwise} \end{cases} \quad \text{as } T \rightarrow 0$$

- in other words, at low temperature $S(\boldsymbol{\eta}/T)_c$ returns the most probable state, whereas at high temperatures $S(\boldsymbol{\eta}/T)_c$ returns one of the states with a uniform probability (cfr. **Boltzmann distribution** in physics)

Linear Discriminant Analysis

Softmax



- softmax distribution $S(\eta/T)$, where $\eta = [3, 0, 1]^T$, at different temperatures T
- when the temperature is high (left), the distribution is uniform, whereas when the temperature is low (right), the distribution is “spiky”, with all its mass on the largest element

Linear Discriminant Analysis

LDA

- in order to find the decision boundaries we impose

$$p(y = c|\mathbf{x}, \boldsymbol{\theta}) = p(y = c'|\mathbf{x}, \boldsymbol{\theta})$$

which entails

$$e^{\boldsymbol{\beta}_c^T \mathbf{x} + \gamma_c} = e^{\boldsymbol{\beta}_{c'}^T \mathbf{x} + \gamma_{c'}}$$

- in this case, taking the logs returns

$$\boldsymbol{\beta}_c^T \mathbf{x} + \gamma_c = \boldsymbol{\beta}_{c'}^T \mathbf{x} + \gamma_{c'}$$

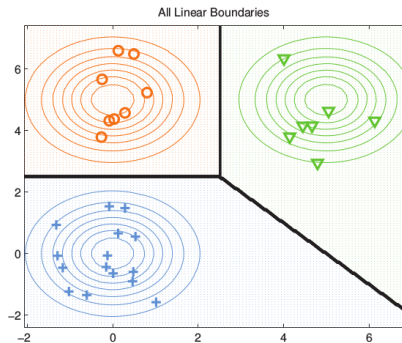
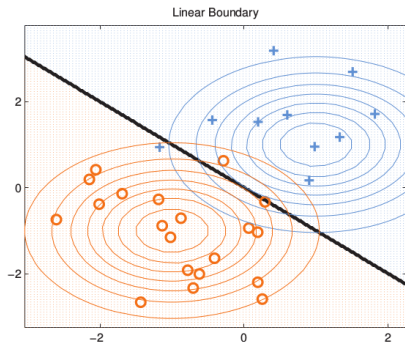
which in turn corresponds to a **linear decision boundary**¹

$$(\boldsymbol{\beta}_c - \boldsymbol{\beta}_{c'})^T \mathbf{x} = -(\gamma_c - \gamma_{c'})$$

¹in D dimensions this corresponds to an hyperplane, in 3D to a plane, in 2D to a straight line

Linear Discriminant Analysis

LDA



- *left*: dataset with 2 classes
- *right*: dataset with 3 classes

Linear Discriminant Analysis

two-class LDA

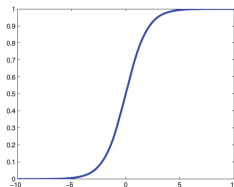
- let us consider an LDA with just two classes (i.e. $y \in \{0, 1\}$)
- in this case

$$p(y = 1|\mathbf{x}, \boldsymbol{\theta}) = \frac{e^{\beta_1^T \mathbf{x} + \gamma_1}}{e^{\beta_1^T \mathbf{x} + \gamma_1} + e^{\beta_0^T \mathbf{x} + \gamma_0}} = \frac{1}{1 + e^{(\beta_0 - \beta_1)^T \mathbf{x} + (\gamma_0 - \gamma_1)}}$$

that is

$$p(y = 1|\mathbf{x}, \boldsymbol{\theta}) = \text{sigm}((\beta_0 - \beta_1)^T \mathbf{x} + (\gamma_0 - \gamma_1))$$

where $\text{sigm}(\eta) \triangleq \frac{1}{1 + \exp(-\eta)}$ is the **sigmoid function** (aka logistic function)



Linear Discriminant Analysis

two-class LDA

- the linear decision boundary is

$$(\beta_0 - \beta_1)^T \mathbf{x} + (\gamma_0 - \gamma_1) = 0$$

- if we define

$$\mathbf{w} \triangleq \beta_1 - \beta_0 = \Sigma^{-1}(\mu_1 - \mu_0)$$

$$\mathbf{x}_0 \triangleq \frac{1}{2}(\mu_1 + \mu_0) - (\mu_1 - \mu_0) \frac{\log(\pi_1/\pi_0)}{(\mu_1 - \mu_0)^T \Sigma^{-1}(\mu_1 - \mu_0)}$$

we obtain $\mathbf{w}^T \mathbf{x}_0 = -(\gamma_1 - \gamma_0)$

- the linear decision boundary can be rewritten as

$$\mathbf{w}^T (\mathbf{x} - \mathbf{x}_0) = 0$$

- in fact we have

$$p(y = 1 | \mathbf{x}, \theta) = \text{sigm}(\mathbf{w}^T (\mathbf{x} - \mathbf{x}_0))$$

Linear Discriminant Analysis

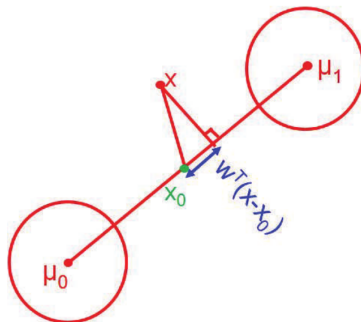
two-class LDA

- we have

$$\mathbf{w} \triangleq \beta_1 - \beta_0 = \Sigma^{-1}(\mu_1 - \mu_0)$$

$$\mathbf{x}_0 \triangleq \frac{1}{2}(\mu_1 + \mu_0) - (\mu_1 - \mu_0) \frac{\log(\pi_1/\pi_0)}{(\mu_1 - \mu_0)^T \Sigma^{-1}(\mu_1 - \mu_0)}$$

- the linear decision boundary is $\mathbf{w}^T(\mathbf{x} - \mathbf{x}_0) = 0$
- in the case $\Sigma_1 = \Sigma_2 = \mathbf{I}$ and $\pi_1 = \pi_0$, one has $\mathbf{w} = \mu_1 - \mu_0$ and $\mathbf{x}_0 = \frac{1}{2}(\mu_1 + \mu_0)$



- 1 Basics
 - Multivariate Gaussian
- 2 MLE for an MVN
 - Theorem
- 3 Gaussian Discriminant Analysis
 - Generative Classifiers
 - Gaussian Discriminant Analysis (GDA)
 - Quadratic Discriminant Analysis (QDA)
 - Linear Discriminant Analysis (LDA)
 - **MLE for Gaussian Discriminant Analysis**
 - Diagonal LDA
 - Bayesian Procedure

MLE for GDA

- how to fit the GDA model?
- the simplest way is to use MLE
- let's assume iid samples, then it is $p(\mathcal{D}|\theta) = \prod_{i=1}^N p(\mathbf{x}_i, y_i|\theta)$
- one has

$$p(\mathbf{x}_i, y_i|\theta) = p(\mathbf{x}_i|y_i, \theta)p(y_i|\pi)$$
$$p(\mathbf{x}_i|y_i, \theta) = \prod_c \mathcal{N}(\mathbf{x}_i|\mu_c, \Sigma_c)^{\mathbb{I}(y_i=c)} \quad p(y_i|\pi) = \prod_c \pi_c^{\mathbb{I}(y_i=c)}$$

where θ is a compound parameter vector containing the parameters π , μ_c and Σ_c

- the log-likelihood function is

$$\log p(\mathcal{D}|\theta) = \left[\sum_{i=1}^N \sum_{c=1}^C \mathbb{I}(y_i = c) \log \pi_c \right] + \sum_{c=1}^C \left[\sum_{i:y_i=c} \log \mathcal{N}(\mathbf{x}_i|\mu_c, \Sigma_c) \right]$$

which is the sum of $C + 1$ distinct terms: the first depending on π and the other C terms depending both on μ_c and Σ_c

- we can estimate each parameter by optimizing the log-likelihood separately w.r.t. it

- the log-likelihood function is

$$\log p(\mathcal{D}|\boldsymbol{\theta}) = \left[\sum_{i=1}^N \sum_{c=1}^C \mathbb{I}(y_i = c) \log \pi_c \right] + \sum_{c=1}^C \left[\sum_{i:y_i=c} \log \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_c, \boldsymbol{\Sigma}_c) \right]$$

- for the class prior, as with the NBC model, we have

$$\hat{\pi}_c = \frac{N_c}{N}$$

- for the class conditional densities, we partition the data based on its class label, and compute the MLE for each Gaussian term

$$\hat{\boldsymbol{\mu}}_c = \frac{1}{N_c} \sum_{i:y_i=c}^{N_c} \mathbf{x}_i$$

$$\hat{\boldsymbol{\Sigma}}_c = \frac{1}{N_c} \sum_{i:y_i=c}^{N_c} (\mathbf{x}_i - \hat{\boldsymbol{\mu}}_c)(\mathbf{x}_i - \hat{\boldsymbol{\mu}}_c)^T$$

- once the model is fit and the parameters are estimated we can make predictions by using a plug-in approximation

$$p(y = c | \mathbf{x}, \hat{\boldsymbol{\theta}}) \propto \hat{\pi}_c |2\pi \hat{\boldsymbol{\Sigma}}_c|^{-1/2} \exp[-\frac{1}{2}(\mathbf{x} - \hat{\boldsymbol{\mu}}_c)^T \hat{\boldsymbol{\Sigma}}_c^{-1}(\mathbf{x} - \hat{\boldsymbol{\mu}}_c)]$$

Overfitting for GDA

- the MLE is fast and simple, however it can badly overfit in high dimensions
- in particular, $\hat{\Sigma}_c = \frac{1}{N_c} \sum_{i:y_i=c}^{N_c} (\mathbf{x}_i - \hat{\mu}_c)(\mathbf{x}_i - \hat{\mu}_c)^T \in \mathbb{R}^{D \times D}$ is singular for $N_c < D$
- even when $N_c > D$, the MLE can be ill-conditioned (close to singular)
- possible simple strategies to solve this issue (they reduce the number of parameters)
 - use NBC model/assumption (i.e. Σ_c are diagonal)
 - use LDA (i.e. $\Sigma_c = \Sigma$)
 - use diagonal LDA (i.e. $\Sigma_c = \Sigma = \text{diag}(\sigma_1^2, \dots, \sigma_D^2)$) (following subsection)
 - use Bayesian approach: estimate full covariance by imposing a prior and then integrating out (following subsection)

1 Basics

- Multivariate Gaussian

2 MLE for an MVN

- Theorem

3 Gaussian Discriminant Analysis

- Generative Classifiers
- Gaussian Discriminant Analysis (GDA)
- Quadratic Discriminant Analysis (QDA)
- Linear Discriminant Analysis (LDA)
- MLE for Gaussian Discriminant Analysis
- **Diagonal LDA**
- Bayesian Procedure

Diagonal LDA

- the diagonal LDA assumes $\Sigma_c = \Sigma = \text{diag}(\sigma_1^2, \dots, \sigma_D^2)$ for $c \in \{1, \dots, C\}$
- one has

$$p(\mathbf{x}_i, y_i = c | \boldsymbol{\theta}) = p(\mathbf{x}_i | y_i = c, \boldsymbol{\theta}_c) p(y_i = c | \boldsymbol{\pi}) = \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_c, \Sigma) \pi_c = \prod_{j=1}^D \mathcal{N}(x_{ij} | \mu_{cj}, \sigma_j^2)$$

and taking the logs

$$\log p(\mathbf{x}_i, y_i = c | \boldsymbol{\theta}) = - \sum_{j=1}^D \frac{(x_{ij} - \mu_{cj})^2}{2\sigma_j^2} + \log \pi_c$$

- typically the estimates of the parameters are

$$\hat{\mu}_{cj} = \frac{1}{N_c} \sum_{i: y_i = c} x_{ij}$$

$$\hat{\sigma}_j^2 = \frac{1}{N - C} \sum_{c=1}^C \sum_{i: y_i = c} (x_{ij} - \hat{\mu}_{cj})^2 \quad (\text{pooled empirical variance})$$

- in high-dimensional settings, this model can work much better than LDA and RDA

1 Basics

- Multivariate Gaussian

2 MLE for an MVN

- Theorem

3 Gaussian Discriminant Analysis

- Generative Classifiers
- Gaussian Discriminant Analysis (GDA)
- Quadratic Discriminant Analysis (QDA)
- Linear Discriminant Analysis (LDA)
- MLE for Gaussian Discriminant Analysis
- Diagonal LDA
- Bayesian Procedure

Bayesian Procedure

- we now follow the full Bayesian procedure to fit the GDA model
- let's restart from the expression of the **posterior predictive** PDF

$$p(y = c|\mathbf{x}, \mathcal{D}) = \frac{p(y = c, \mathbf{x}|\mathcal{D})}{p(\mathbf{x}|\mathcal{D})} = \frac{p(\mathbf{x}|y = c, \mathcal{D})p(y = c|\mathcal{D})}{p(\mathbf{x}|\mathcal{D})}$$

- since we are interested in computing

$$c^* = \underset{c}{\operatorname{argmax}} \quad p(y = c|\mathbf{x}, \mathcal{D})$$

we can neglect the constant $p(\mathbf{x}|\mathcal{D})$ and use the following simpler expression

$$p(y = c|\mathbf{x}, \mathcal{D}) \propto p(\mathbf{x}|y = c, \mathcal{D})p(y = c|\mathcal{D})$$

- note that we didn't use the model parameters in the previous equation
- now we use the **Bayesian procedure** in which we integrate out the unknown parameters
- for simplicity we now consider a vector parameter π for the PMF $p(y = c|\mathcal{D})$ and a vector parameter θ_c for the PDF $p(\mathbf{x}|y = c, \mathcal{D})$

Bayesian Procedure

- as for the PMF $p(y = c|\mathcal{D})$ we can integrate out π as follows

$$p(y = c|\mathcal{D}) = \int p(y = c, \pi|\mathcal{D})d\pi$$

- we know that $y \sim \text{Cat}(\pi)$ i.e. $p(y|\pi) = \prod_c \pi_c^{\mathbb{I}(y=c)}$
- we can decompose $p(y = c, \pi|\mathcal{D})$ as follows

$$p(y = c, \pi|\mathcal{D}) = p(y = c|\pi, \mathcal{D})p(\pi|\mathcal{D}) = p(y = c|\pi)p(\pi|\mathcal{D}) = \pi_c p(\pi|\mathcal{D})$$

where $p(\pi|\mathcal{D})$ is the posterior w.r.t. π

- using the previous equation in integral above we have

$$p(y = c|\mathcal{D}) = \int p(y = c, \pi|\mathcal{D})d\pi = \int \pi_c p(\pi|\mathcal{D})d\pi = \mathbb{E}[\pi_c|\mathcal{D}] = \frac{N_c + \alpha_c}{N + \alpha_0}$$

which is the posterior mean computed for the **Dirichlet-multinomial** model (cfr lecture 4 slides)

Bayesian Procedure

- as for the PDF $p(\mathbf{x}|y = c, \mathcal{D})$ we can integrate out θ_c as follows

$$p(\mathbf{x}|y = c, \mathcal{D}) = \int p(\mathbf{x}, \theta_c|y = c, \mathcal{D})d\theta_c = \int p(\mathbf{x}, \theta_c|\mathcal{D}_c)d\theta_c$$

where for simplicity we introduce $\mathcal{D}_c \triangleq \{(\mathbf{x}_i, y_i) \in \mathcal{D} | y_i = c\}$

- we know that $p(\mathbf{x}|\theta_c) = \mathcal{N}(\mathbf{x}|\mu_c, \Sigma_c)$ where $\theta_c = (\mu_c, \Sigma_c)$
- we can use the following decomposition

$$p(\mathbf{x}, \theta_c|\mathcal{D}_c) = p(\mathbf{x}|\theta_c, \mathcal{D}_c)p(\theta_c|\mathcal{D}_c) = p(\mathbf{x}|\theta_c)p(\theta_c|\mathcal{D}_c)$$

where $p(\theta_c|\mathcal{D}_c)$ is the posterior w.r.t. θ_c

- hence one has

$$\begin{aligned} p(\mathbf{x}|y = c, \mathcal{D}) &= \int p(\mathbf{x}, \theta_c|\mathcal{D}_c)d\theta_c = \int p(\mathbf{x}|\theta_c)p(\theta_c|\mathcal{D}_c)d\theta_c = \\ &= \int \int \mathcal{N}(\mathbf{x}|\mu_c, \Sigma_c)p(\mu_c, \Sigma_c|\mathcal{D}_c)d\mu_c d\Sigma_c \end{aligned}$$

- one has

$$p(\mathbf{x}|y = c, \mathcal{D}) = \int \int \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_c, \boldsymbol{\Sigma}_c) p(\boldsymbol{\mu}_c, \boldsymbol{\Sigma}_c|\mathcal{D}_c) d\boldsymbol{\mu}_c d\boldsymbol{\Sigma}_c$$

- the posterior is (see sect. 4.6.3.3 of the book)

$$p(\boldsymbol{\mu}_c, \boldsymbol{\Sigma}_c|\mathcal{D}_c) = \text{NIW}(\mathbf{m}_c, \boldsymbol{\Sigma}_c|\mathbf{m}_N^c, \kappa_N^c, \nu_N^c, \mathbf{S}_N^c)$$

- then (see sect. 4.6.3.6)

$$p(\mathbf{x}|y = c, \mathcal{D}) = \int \int \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_c, \boldsymbol{\Sigma}_c) \text{NIW}(\boldsymbol{\mu}_c, \boldsymbol{\Sigma}_c|\mathbf{m}_N^c, \kappa_N^c, \nu_N^c, \mathbf{S}_N^c) d\boldsymbol{\mu}_c d\boldsymbol{\Sigma}_c =$$

$$p(\mathbf{x}|y = c, \mathcal{D}) = \mathcal{T}(\mathbf{x}|\mathbf{m}_N^c, \frac{\kappa_N^c + 1}{\kappa_N^c(\nu_N^c - D + 1)} \mathbf{S}_N^c, \nu_N^c - D + 1)$$

Bayesian Procedure

- let's summarize what we obtained by applying the Bayesian procedure
- we first found

$$p(y = c|\mathcal{D}) = \mathbb{E}[\pi_c|\mathcal{D}] = \frac{N_c + \alpha_c}{N + \alpha_0}$$

and then

$$p(\mathbf{x}|y = c, \mathcal{D}) = \mathcal{T}(\mathbf{x}|\mathbf{m}_N^c, \frac{\kappa_N^c + 1}{\kappa_N^c(\nu_N^c - D + 1)} \mathbf{S}_N^c, \nu_N^c - D + 1)$$

- then combining everything in the starting posterior predictive we have

$$\begin{aligned} p(y = c|\mathbf{x}, \mathcal{D}) &\propto p(\mathbf{x}|y = c, \mathcal{D})p(y = c|\mathcal{D}) = \\ &= \mathbb{E}[\pi_c|\mathcal{D}]\mathcal{T}(\mathbf{x}|\mathbf{m}_N^c, \frac{\kappa_N^c + 1}{\kappa_N^c(\nu_N^c - D + 1)} \mathbf{S}_N^c, \nu_N^c - D + 1) \end{aligned}$$

- Kevin Murphy's book