

Lecture 09: Linear Discriminant Analysis

Introduction to Machine Learning [25737]

Sajjad Amini

Sharif University of Technology

- 1 Approach Definitions
- 2 Gaussian Discriminant Analysis
- 3 Connection Between LDA and MLR
- 4 Naive Bayes Classifier
- 5 Generative vs. Discriminative

Except explicitly cited, the reference for the material in slides is:

- Murphy, K. P. (2022). *Probabilistic machine learning: an introduction*. MIT press.

Section 1

Approach Definitions

Discriminant Analysis

Assume we consider the following model for classification:

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

where:

$$p(y = c|\boldsymbol{\theta})$$

Prior probability over labels

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

Class conditional density

Using special options for class conditional density, we can show that:

$$p(y|\mathbf{x}; \boldsymbol{\theta}) = \mathbf{w}^T \mathbf{x} + \text{constant}$$

The resulting model is known as linear discriminant analysis.

Section 2

Gaussian Discriminant Analysis

Class Conditional Density

For Gaussian discriminant analysis, the class conditional density is:

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

The above selection result in the following posterior over class labels:

$$p(y = c|\mathbf{x}, \boldsymbol{\theta}) \propto \pi_c \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_c, \boldsymbol{\Sigma}_c)$$

where

$$\pi_c = p(y = c|\boldsymbol{\theta})$$

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

Decision Boundary

Quadratic Decision Boundary

Consider the log posterior probability as:

$$\log p(y = c | \mathbf{x}; \boldsymbol{\theta}) = \log \pi_c - \frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_c)^T \boldsymbol{\Sigma}_c^{-1}(\mathbf{x} - \boldsymbol{\mu}_c) + \text{const}$$

This method is called *Quadratic Discriminant Analysis* (QDA) because the decision boundary is a quadratic function.

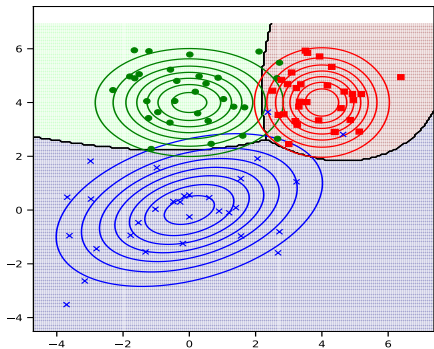
Linear Decision Boundary

If we assume $\boldsymbol{\Sigma} = \boldsymbol{\Sigma}_c$, then:

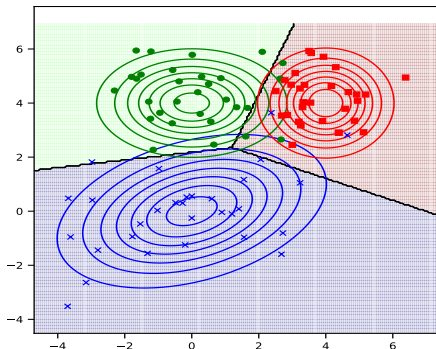
$$\begin{aligned} \log p(y = c | \mathbf{x}; \boldsymbol{\theta}) &= \log \pi_c - \frac{1}{2} \boldsymbol{\mu}_c^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_c + \mathbf{x}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_c + \text{const} - \frac{1}{2} \mathbf{x}^T \boldsymbol{\Sigma}^{-1} \mathbf{x} \\ &= \gamma_c + \mathbf{x}^T \boldsymbol{\beta}_c + \kappa \end{aligned}$$

This method is called *Linear Discriminant Analysis* (LDA) because the decision boundary is a linear function.

Quadratic vs Linear Discriminant Analysis



(a) QDA



(b) LDA

Figure: Decision boundary comparison

Section 3

Connection Between LDA and MLR

Similarity

As we can see, LDA can be formulated as:

$$p(y = c | \mathbf{x}, \boldsymbol{\theta}) = \frac{\exp(\boldsymbol{\beta}_c^T \mathbf{x} + \gamma_c)}{\sum_{c'=1}^C \exp(\boldsymbol{\beta}_{c'}^T \mathbf{x} + \gamma_{c'})} = \frac{\exp(\mathbf{w}_c^T [1; \mathbf{x}])}{\sum_{c'=1}^C \exp(\mathbf{w}_{c'}^T [1; \mathbf{x}])}$$

Thus the posterior form is similar to MLR.

Difference

- In LDA, we first estimate prior probability over labels and class conditional density, and derive $\{\mathbf{w}_c\}_{c=1}^C$ from them.
- In MLR, we estimate $\{\mathbf{w}_c\}_{c=1}^C$ directly to maximize conditional likelihood $p(y | \mathbf{x}, \boldsymbol{\theta})$

MLE

The likelihood function can be formulated as:

$$\begin{aligned}
 p(\mathcal{D}|\boldsymbol{\theta}) &= p(\{(\mathbf{x}_n, y_n)\}_{n=1}^N|\boldsymbol{\theta}) \stackrel{(1)}{=} \prod_{n=1}^N p(\mathbf{x}_n, y_n|\boldsymbol{\theta}) \\
 &\stackrel{(2)}{=} \prod_{n=1}^N p(y_n|\boldsymbol{\theta})p(\mathbf{x}_n|y_n, \boldsymbol{\theta}) = \prod_{n=1}^N \text{Cat}(y_n|\boldsymbol{\pi}) \prod_{c=1}^C \mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_c, \boldsymbol{\Sigma}_c)^{\mathbb{I}(y_n=c)}
 \end{aligned}$$

where we use independency of training samples and probability chain rule for equality (1) and (2), respectively. Note that the parameter vector is:

$$\boldsymbol{\theta} = [\boldsymbol{\pi}; \boldsymbol{\mu}_1; \dots; \boldsymbol{\mu}_C; \text{vec}(\boldsymbol{\Sigma}_1); \dots; \text{vec}(\boldsymbol{\Sigma}_C)]$$

MLE (Continue)

The likelihood function and its log version are:

$$p(\mathcal{D}|\boldsymbol{\theta}) = \prod_{n=1}^N \text{Cat}(y_n|\boldsymbol{\pi}) \prod_{c=1}^C \mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_c, \boldsymbol{\Sigma}_c)^{\mathbb{I}(y_n=c)}$$

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

Using differentiation, we can calculate the model parameters as:

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

$$\hat{\boldsymbol{\mu}}_c = \frac{1}{N_c} \sum_{n:y_n=c} \mathbf{x}_n$$

$$\hat{\boldsymbol{\Sigma}}_c = \frac{1}{N_c} \sum_{n:y_n=c} (\mathbf{x}_n - \hat{\boldsymbol{\mu}}_c)(\mathbf{x}_n - \hat{\boldsymbol{\mu}}_c)^T$$

Tied Covariance Matrices

Tied Covariance Matrices

Tied covariance matrices is the situation where we force all covariance matrices to be equal as:

$$\Sigma_c = \Sigma, \quad c = 1, \dots, C$$

MLE estimation for tied covariance matrix is:

$$\hat{\Sigma} = \frac{1}{N} \sum_{c=1}^C \sum_{n: y_n=c} (\mathbf{x}_n - \hat{\boldsymbol{\mu}}_c)(\mathbf{x}_n - \hat{\boldsymbol{\mu}}_c)^T$$

LDA and Tied Covariance Matrix

When the covariance matrix is tied, QDA simplifies to LDA.

Diagonal LDA

We can simplify tied covariance matrix further by assuming it to be diagonal, so: $\Sigma_c = \mathbf{D}$, $c = 1, \dots, C$

Nearest Centroid Classifier

Nearest Centroid Classifier

Assume the prior probability over classes is uniform, so:

$$\pi_c = \frac{1}{C}, \quad c = 1, \dots, C$$

If the covariance matrices are tied, then:

$$\begin{aligned}\hat{y}(\mathbf{x}) &= \operatorname{argmax}_c \log p(y = c | \mathbf{x}, \boldsymbol{\theta}) = \operatorname{argmin}_c (\mathbf{x} - \boldsymbol{\mu}_c)^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}_c) \\ &= \operatorname{argmin}_c \Delta_{\boldsymbol{\Sigma}}^2(\mathbf{x}, \boldsymbol{\mu}_c)\end{aligned}$$

Thus the class whose mean has minimum Mahalanobis distance to the query point \mathbf{x} is selected as the label.

Section 4

Naive Bayes Classifier

Naive Bayes Classifier

Main Assumption

The input features are mutually independent given the class label. In other words:

$$p(\mathbf{x}|y = c, \boldsymbol{\theta}) = \prod_{d=1}^D p(x_d|y = c, \boldsymbol{\theta}_{dc})$$

where $\boldsymbol{\theta}_{dc}$ is model parameter vector for conditional density for class c and feature d . The posterior over class label is:

$$p(y = c|\mathbf{x}, \boldsymbol{\theta}) = \frac{p(y = c|\boldsymbol{\pi}) \prod_{d=1}^D p(x_d|y = c, \boldsymbol{\theta}_{dc})}{\sum_{c'} p(y = c'|\boldsymbol{\pi}) \prod_{d=1}^D p(x_d|y = c', \boldsymbol{\theta}_{dc'})}$$

Pros and cons

- The naive model may not hold in many real world application.
- Naive Bayes model is relatively immune to overfitting.

Binary Features

In this case $x_d \in \{0, 1\}$ and thus the class conditional density is:

$$p(\mathbf{x}|y = c, \boldsymbol{\theta}) = \prod_{d=1}^D \text{Ber}(x_d|\theta_{dc})$$

where θ_{dc} shows the probability that $x_d = 1$ in class c . This model is known as *multivariate Bernoulli naive Bayes*.

Categorical Features

In this case $x_d \in \{0, 1, \dots, K\}$ and thus the class conditional density is:

$$p(\mathbf{x}|y = c, \boldsymbol{\theta}) = \prod_{d=1}^D \text{Cat}(x_d|\boldsymbol{\theta}_{dc})$$

where θ_{dck} shows the probability that $x_d = k$ in class c .

Real-values Features

In this case $x_d \in \mathbb{R}$ and thus we can use univariate Gaussian for each dimension in each class. Thus the class conditional density is:

$$p(\mathbf{x}|y = c, \boldsymbol{\theta}) = \prod_{d=1}^D \mathcal{N}(x_d|\mu_{dc}, \sigma_{dc}^2)$$

where μ_{dc} and σ_{dc}^2 shows the mean and variance of feature d in class c .

MLE

The likelihood for the dataset \mathcal{D} is:

$$\begin{aligned} p(\mathcal{D}|\boldsymbol{\theta}) &= \prod_{n=1}^N p(y_n|\boldsymbol{\theta})p(\mathbf{x}_n|y_n, \boldsymbol{\theta}) \stackrel{(1)}{=} \prod_{n=1}^N p(y_n|\boldsymbol{\theta}) \prod_{d=1}^D p(x_{nd}|y_n, \boldsymbol{\theta}_d) \\ &= \prod_{n=1}^N \text{Cat}(y_n|\boldsymbol{\pi}) \prod_{d=1}^D \prod_{c=1}^C p(x_{nd}|\boldsymbol{\theta}_{dc})^{\mathbb{I}(y_n=c)} \end{aligned}$$

where we use Naive Bayes assumption for equality (1).

MLE (Continue)

$$\begin{aligned}\log p(\mathcal{D}|\boldsymbol{\theta}) &= \log \prod_{n=1}^N \left(\left[\prod_{c=1}^C \pi_c^{\mathbb{I}(y_n=c)} \right] \left[\prod_{d=1}^D \prod_{c=1}^C p(x_{nd}|\boldsymbol{\theta}_{dc})^{\mathbb{I}(y_n=c)} \right] \right) \\ &= \sum_{n=1}^N \left(\log \left[\prod_{c=1}^C \pi_c^{\mathbb{I}(y_n=c)} \right] + \log \left[\prod_{d=1}^D \prod_{c=1}^C p(x_{nd}|\boldsymbol{\theta}_{dc})^{\mathbb{I}(y_n=c)} \right] \right) \\ &= \sum_{n=1}^N \left(\left[\sum_{c=1}^C \mathbb{I}(y_n = c) \log \pi_c \right] + \left[\sum_{d=1}^D \sum_{c=1}^C \mathbb{I}(y_n = c) \log p(x_{nd}|\boldsymbol{\theta}_{dc}) \right] \right) \\ &= \left[\sum_{n=1}^N \sum_{c=1}^C \mathbb{I}(y_n = c) \log \pi_c \right] + \left[\sum_{n=1}^N \sum_{d=1}^D \sum_{c=1}^C \mathbb{I}(y_n = c) \log p(x_{nd}|\boldsymbol{\theta}_{dc}) \right]\end{aligned}$$

MLE for π

Irrespective of class conditional density, the MLE for π is the vector of empirical counts as $\hat{\pi}_c = \frac{N_c}{N}$.

MLE for θ_{dc}

- Binary features: $\hat{\theta}_{dc} = \frac{N_{dc}}{N_c}$
- Categorical features: $\hat{\theta}_{dck} = \frac{N_{dck}}{N_c}$, $k = 1, \dots, K$
- Real-valued features (Univariate Gaussian):

$$\hat{\mu}_{dc} = \frac{1}{N_c} \sum_{n=1}^N \mathbb{I}(y_n = c) x_{nd}$$

$$\hat{\sigma}_{dc}^2 = \frac{1}{N_c} \sum_{n=1}^N \mathbb{I}(y_n = c) (x_{nd} - \hat{\mu}_{dc})^2$$

Section 5

Generative vs. Discriminative

MLR vs DA

In MLR, we have $p(y|\mathbf{x}; \boldsymbol{\theta}) = \text{Cat}(y|\mathcal{S}(\mathbf{W}^T \mathbf{x} + \mathbf{b}))$ and the likelihood is:

$$p(\mathcal{D}|\boldsymbol{\theta}) = p(\{y_n\}_{n=1}^N | \{\mathbf{x}_n\}_{n=1}^N, \boldsymbol{\theta}) = \prod_{n=1}^N p(y_n | \mathbf{x}_n, \boldsymbol{\theta})$$

In Discriminant analysis, we have $p(y = c|\mathbf{x}; \boldsymbol{\theta}) = \frac{p(\mathbf{x}|y=c;\boldsymbol{\theta})p(y=c|\boldsymbol{\theta})}{\sum_{c'} p(\mathbf{x}|y=c';\boldsymbol{\theta})p(y=c'|\boldsymbol{\theta})}$ and the likelihood is:

$$p(\mathcal{D}|\boldsymbol{\theta}) = p(\{(\mathbf{x}_n, y_n)\}_{n=1}^N | \boldsymbol{\theta}) = \prod_{n=1}^N p(\mathbf{x}_n, y_n | \boldsymbol{\theta}) = \prod_{n=1}^N p(\mathbf{x}_n | y_n, \boldsymbol{\theta}) p(y_n | \boldsymbol{\theta})$$

Discriminative

By training MLR:

- You have access to $p(y|\mathbf{x}, \boldsymbol{\theta})$ which can be used to generate label for a query input \mathbf{x} (discriminate the label of \mathbf{x}).
- You can't generate samples from specific class $y = k$.

Generative

By training DA:

- You have access to $p(y|\mathbf{x}, \boldsymbol{\theta})$ which can be used to generate label for a query input \mathbf{x} (discriminate the label of \mathbf{x}).
- You have access to $p(\mathbf{x}|y, \boldsymbol{\theta})$ that can be used to generate samples from specific class $y = k$.