

# Introduction to Machine learning

Sylvain Le Corff

## 1. Mathematical framework

## 2. Discriminant analysis

- The multivariate normal distribution

- Bayes classifier for multivariate normal distributions

## 1. Mathematical framework

## 2. Discriminant analysis

The multivariate normal distribution

Bayes classifier for multivariate normal distributions

## Supervised Learning Framework

- **Input** measurement  $\mathbf{X} \in \mathcal{X}$  (often  $\mathcal{X} \subset \mathbb{R}^d$ ).
- **Output** measurement  $Y \in \mathcal{Y}$ .
- The joint distribution of  $(\mathbf{X}, Y)$  is **unknown**.
- $Y \in \{1, \dots, M\}$  (classification) or  $Y \in \mathbb{R}^m$  (regression).
- A **predictor** is a measurable function  $f : \mathcal{X} \rightarrow \mathcal{Y}$ .

## Training data

- i.i.d. with the same distribution as  $(\mathbf{X}, Y)$ :

$$\mathcal{D}_n = \{(\mathbf{X}_1, Y_1), \dots, (\mathbf{X}_n, Y_n)\}.$$

## Goal

- Construct a **good** predictor  $\hat{f}_n$  from the training data.
- Need to specify the meaning of good.

## Loss function

- $\ell(Y, f(\mathbf{X}))$ : the goodness of the prediction of  $Y$  by  $f(\mathbf{X})$ .
- **Prediction** loss:  $\ell(Y, f(\mathbf{X})) = \mathbf{1}_{Y \neq f(\mathbf{X})}$ .
- **Quadratic** loss:  $\ell(Y, \mathbf{X}) = \|Y - f(\mathbf{X})\|_2^2$ .

## Risk function

- Risk measured as the average loss:

$$\mathcal{R}(f) = \mathbb{E}[\ell(Y, f(\mathbf{X}))].$$

- **Prediction** loss:  $\mathbb{E}[\ell(Y, f(\mathbf{X}))] = \mathbb{P}(Y \neq f(\mathbf{X}))$ .
- **Quadratic** loss:  $\mathbb{E}[\ell(Y, f(\mathbf{X}))] = \mathbb{E}[\|Y - f(\mathbf{X})\|_2^2]$ .
- **Beware**: As  $\hat{f}_n$  depends on  $\mathcal{D}_n$ ,  $\mathcal{R}(\hat{f}_n)$  is a random variable!

## 1. Mathematical framework

## 2. Discriminant analysis

The multivariate normal distribution

Bayes classifier for multivariate normal distributions

## Definition

Let  $\mu \in \mathbb{R}^d$ ,  $\Sigma$  be a positive definite matrix. We write  $X \sim \mathcal{N}(\mu, \Sigma)$  when the Lebesgue density of  $X$  is

$$\begin{aligned} x \in \mathbb{R}^d &\mapsto |2\pi\Sigma|^{-\frac{1}{2}} e^{-\frac{1}{2}(x-\mu)^\top \Sigma^{-1}(x-\mu)} \\ &= \frac{1}{(2\pi)^{\frac{d}{2}} |\Sigma|^{\frac{1}{2}}} e^{-\frac{1}{2}(x-\mu)^\top \Sigma^{-1}(x-\mu)}, \end{aligned}$$

where  $|\Sigma|$  is the determinant of  $\Sigma$ . In addition, we have

$$\mathbb{E}[X] = \mu, \quad \mathbb{V}[X] = \Sigma,$$

where  $\mathbb{V}[X]$  is the covariance matrix of  $X$ .

## Proposition

Let  $\mu^* \in \mathbb{R}^d$ ,  $\Sigma^*$  be a positive definite matrix and  $\{X_1, \dots, X_n\}$  be a sample i.i.d. according to  $\mathcal{N}(\mu^*, \Sigma^*)$ .

Then

$$\hat{\mu}_n = \frac{1}{n} \sum_{i=1}^n X_i$$

and

$$\hat{\Sigma}_n = \frac{1}{n} \sum_{i=1}^n (X_i - \hat{\mu})(X_i - \hat{\mu})^\top$$

are maximum likelihood estimators respectively of  $\mu^*$  and  $\Sigma^*$ .

Proof on blackboard



## Bayes classifier

The **Bayes classifier**  $g^*$  is defined as:

$$g^*(X) = \begin{cases} 1 & \text{if } \mathbb{P}(Y = 1|X) > \mathbb{P}(Y = 0|X), \\ 0 & \text{otherwise.} \end{cases}$$

Equivalently,

$$g^*(X) = \begin{cases} 1 & \text{if } \mathbb{P}(Y = 1|X) > 1/2, \\ 0 & \text{otherwise,} \end{cases}$$

## Lemma

*For any classification rule  $g : \mathbb{R}^d \rightarrow \{0, 1\}$ , one has*

$$\mathcal{R}(g^*) \leq \mathcal{R}(g).$$

In the case where  $Y \in \{1, \dots, M\}$  for  $M > 1$ .

Bayes classifier

$$g^*(X) \in \operatorname{argmax}_{i \in \{1, \dots, M\}} \mathbb{P}(Y = i | X).$$

In practice **we do not know the conditional law of  $Y$  given  $X$** .  
Several solutions to overcome this issue.

## Fully parametric modeling.

Estimate the law of  $(X, Y)$  and use the **Bayes formula** to deduce an estimate of the conditional law of  $Y$ : *LDA/QDA, Naive Bayes...*

## Parametric conditional modeling.

Estimate the conditional law of  $Y$  by a **parametric** law: *linear regression, logistic regression, Feed Forward Neural Networks...*

## Nonparametric conditional modeling.

Estimate the conditional law of  $Y$  by a **non parametric** estimate: *kernel methods, nearest neighbors...*

- ▶  $(X, Y) \in \mathbb{R}^d \times \{1, \dots, M\}$  be a pair of r.v.
- ▶  $Y$  is a label characterizing the class of  $X$ .
- ▶ **Goal:** computing the Bayes classifier when for all  $i \in \{1, \dots, M\}$  the conditional distribution of  $X$  given  $\{Y = i\}$  is Gaussian with positive definite matrix  $\Sigma_i$  and mean  $\mu_i \in \mathbb{R}^d$ .

Recall: a Bayes classifier

For multiclass

$$g^*(X) \in \operatorname{argmax}_{i \in \{1, \dots, M\}} \mathbb{P}(Y = i | X).$$

Assume that for all  $i \in \{1, \dots, M\}$ ,  $\mathbb{P}(Y = i) = \pi_i$ , where  $\pi_i \in (0, 1)$ .

## Proposition

A Bayes classifier  $g^*$  is defined, for all  $x \in \mathbb{R}^d$ , by

$$g^*(x) \in \operatorname{argmax}_{i \in \{1, \dots, M\}} \log(\pi_i) - \frac{1}{2} \log |\Sigma_i| - \frac{1}{2} (x - \mu_i)^\top \Sigma_i^{-1} (x - \mu_i).$$

Proof on blackboard

- ▶ Only two classes ( $M = 2$ )
- ▶ In this case, a Bayes classifier satisfies

$$g^*: X \mapsto \begin{cases} 1 & \text{if } \mathbb{P}(Y = 1|X) > \mathbb{P}(Y = 2|X) \\ 2 & \text{otherwise.} \end{cases}$$

We assume that the covariance is the same in each class.

$$\Sigma_1 = \Sigma_2 = \Sigma.$$

## Proposition

*Define*

$$h: x \in \mathbb{R}^d \mapsto (\mu_1 - \mu_2)^\top \Sigma^{-1} x$$

$$b = \frac{1}{2}(\mu_2^\top \Sigma^{-1} \mu_2 - \mu_1^\top \Sigma^{-1} \mu_1) + \log \left( \frac{\pi_1}{\pi_2} \right).$$

*Then, a Bayes classifier is*

$$g^*: x \in \mathbb{R}^d \mapsto \begin{cases} 1 & \text{if } h(x) + b > 0 \\ 2 & \text{otherwise.} \end{cases}$$

Proof on blackboard, see also <https://sylvainlc.github.io/>

- ▶ Note that the function  $h(x) + b$  is linear in  $x$ .
- ▶ This is a **linear classifier**!

### What happens when $\pi_1 = \pi_2$

- ▶ if  $\pi_1 = \pi_2$ , we have:

$$g^*(x) = 1$$

$$\iff (x - \mu_1)^\top \Sigma^{-1} (x - \mu_1) < (x - \mu_2)^\top \Sigma^{-1} (x - \mu_2),$$

- ▶  $\pi_1 = \pi_2$  if and only if  $x$  is closer to  $\mu_1$  than  $\mu_2$  with respect to the Mahalanobis distance ruled by  $\Sigma$ .



- ▶ Each class is normally distributed
- ▶ But with **different** covariances

## Proposition

*Define*

$$h: x \in \mathbb{R}^d \mapsto \frac{1}{2}x^\top(\Sigma_2^{-1} - \Sigma_1^{-1})x + (\mu_1^\top \Sigma_1^{-1} - \mu_2^\top \Sigma_2^{-1})x$$
$$b = \frac{1}{2}(\mu_2^\top \Sigma_2^{-1} \mu_2 - \mu_1^\top \Sigma_1^{-1} \mu_1) - \frac{1}{2} \log \left( \frac{|\Sigma_1|}{|\Sigma_2|} \right) + \log \left( \frac{\pi_1}{\pi_2} \right).$$

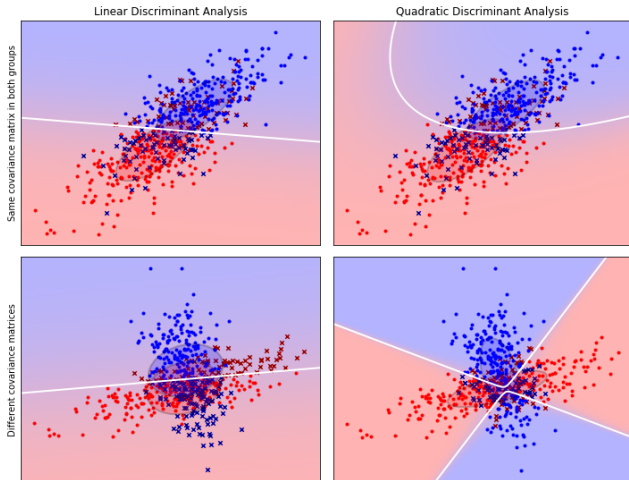
*Then, a Bayes classifier is*

$$g^*: x \in \mathbb{R}^d \mapsto \begin{cases} 1 & \text{if } h(x) + b > 0 \\ 2 & \text{otherwise.} \end{cases}$$

Proof on blackboard, see also <https://sylvainlc.github.io/>

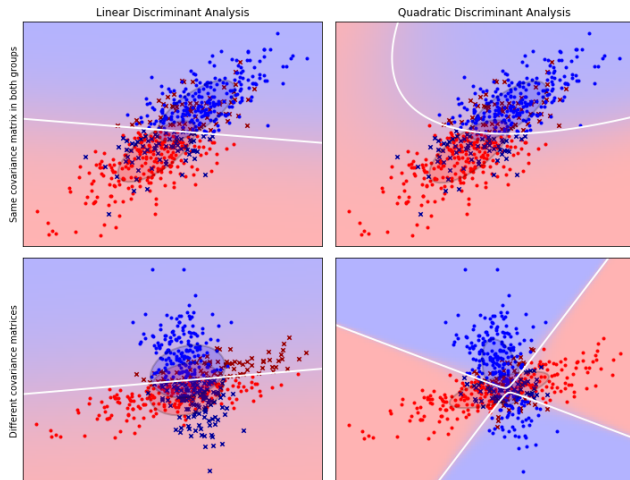
# Comparison of LDA and QDA

18 / 25



**Figure:** (Top) Data are generated with the same covariance matrix in each group. (Bottom) Data are generated with different covariance matrices in the two groups.

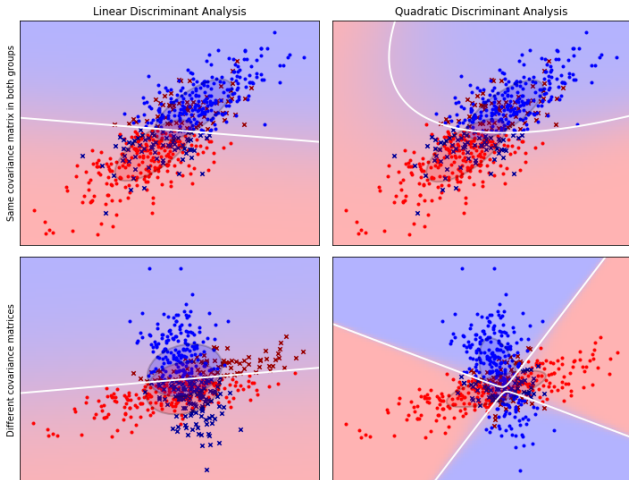
# Comparison of LDA and QDA



**Figure:** (Left) Classification boundary obtained with LDA, assuming the covariance matrix is the same in each group. (Right) Classification obtained with QDA, assuming the covariance matrices are different.

# Comparison of LDA and QDA

20 / 25



**Figure:** Crosses are all false positives i.e. all data wrongly classified by the discriminant analysis. Simulations are inspired by [Discriminant analysis with scikit-learn] and can be found here [?].

Classical algorithm using a **crude modeling for the conditional law of  $X$  given  $Y$**

→ **Feature independence** assumption: all components of  $X$  are independent given  $Y$ .

→ **Simple featurewise model**: binomial if binary, multinomial if finite and Gaussian if continuous.

If all features are continuous, the law of  $X$  given  $Y$  is Gaussian with a **diagonal covariance matrix**!

Very simple learning even in **very high dimension**!

→ **Feature independence** assumption.

For  $k \in \{1, 2\}$ ,  $\mathbb{P}(Y = k) = \pi_k$  and the conditional density of  $X^{(j)}$  given  $\{Y = k\}$  is

$$g_k(x^{(j)}) = (2\pi\sigma_{j,k}^2)^{-1/2} \exp \left\{ -(x^{(j)} - \mu_{j,k})^2 / (2\sigma_{j,k}^2) \right\} .$$

The conditional distribution of  $X$  given  $\{Y = k\}$  is then

$$g_k(x) = (\det(2\pi\Sigma_k))^{-1/2} \exp \left\{ -(x - \mu_k)^T \Sigma_k^{-1} (x - \mu_k) / 2 \right\} ,$$

where  $\Sigma_k = \text{diag}(\sigma_{1,k}^2, \dots, \sigma_{d,k}^2)$  and  $\mu_k = (\mu_{1,k}, \dots, \mu_{d,k})^\top$ .

In a two-classes problem, the optimal classifier is (see **linear discriminant analysis**):

$$f^*(X) = 2\mathbb{1}_{\mathbb{P}(Y=1|X) > \mathbb{P}(Y=-1|X)} - 1.$$

→ When the parameters are unknown, they may be replaced by their **maximum likelihood estimates**. This yields, for  $k \in \{1, 2\}$ ,

$$\hat{\pi}_k^n = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{Y_i=k},$$

$$\hat{\mu}_k^n = \frac{1}{\sum_{i=1}^n \mathbb{1}_{Y_i=k}} \sum_{i=1}^n \mathbb{1}_{Y_i=k} X_i,$$

$$\hat{\Sigma}_k^n = \text{diag} \left( \frac{1}{n} \sum_{i=1}^n (X_i - \hat{\mu}_k^n)(X_i - \hat{\mu}_k^n)^T \mathbb{1}_{Y_i=k} \right).$$

The loglikelihood of the observations is given by

$$\begin{aligned}\ell_n(\theta) = & -\frac{nd}{2} \log(2\pi) - \frac{n}{2} \log \det(\Sigma) + n_1 \log \pi_1 + (n - n_2) \log(1 - \pi_1) \\ & - \frac{1}{2} \sum_{i=1}^n \mathbb{1}_{Y_i=1} (X_i - \mu_1)^\top \Sigma^{-1} (X_i - \mu_1) \\ & - \frac{1}{2} \sum_{i=1}^n \mathbb{1}_{Y_i=-1} (X_i - \mu_{-1})^\top \Sigma^{-1} (X_i - \mu_{-1}),\end{aligned}$$

where  $n_1 = \sum_{i=1}^n \mathbb{1}_{Y_i=1}$ . This yields, for  $k \in \{1, 2\}$ , the following MLE estimates:

$$\begin{aligned}\hat{\pi}_k &= \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{Y_i=k}, \quad \hat{\mu}_k = \frac{1}{\sum_{i=1}^n \mathbb{1}_{Y_i=k}} \sum_{i=1}^n \mathbb{1}_{Y_i=k} X_i, \\ \hat{\Sigma} &= \frac{1}{n} \sum_{i=1}^n (X_i - \hat{\mu}_{Y_i})(X_i - \hat{\mu}_{Y_i})^\top.\end{aligned}$$



We assume that the covariance is the same in each class.

$$\Sigma_1 = \Sigma_2 = \Sigma.$$

## Proposition

*Define*

$$\begin{aligned}\hat{h}: x \in \mathbb{R}^d &\mapsto (\hat{\mu}_1 - \hat{\mu}_2)^\top \hat{\Sigma}^{-1} x \\ \hat{b} &= \frac{1}{2}(\hat{\mu}_2^\top \hat{\Sigma}^{-1} \hat{\mu}_2 - \hat{\mu}_1^\top \hat{\Sigma}^{-1} \hat{\mu}_1) + \log \left( \frac{\hat{\pi}_1}{\hat{\pi}_2} \right).\end{aligned}$$

*Then, a "Plug-in" Bayes classifier is*

$$\hat{g}: x \in \mathbb{R}^d \mapsto \begin{cases} 1 & \text{if } \hat{h}(x) + \hat{b} > 0 \\ 2 & \text{otherwise.} \end{cases}$$

Proof on blackboard, see also <https://sylvainlc.github.io/>