

Machine Learning II

Session 1 : Supervised Learning: Discriminant Analysis

Mohamad GHASSANY

EFREI PARIS

Course Overview

Format:

- ▶ Course sessions: 6 sessions of 5 hours each.
- ▶ Sessions are CTP.

Course chapters:

Session 1: Supervised Learning: Discriminant Analysis

Session 2:

Session 3:

Session 4:

Session 5:

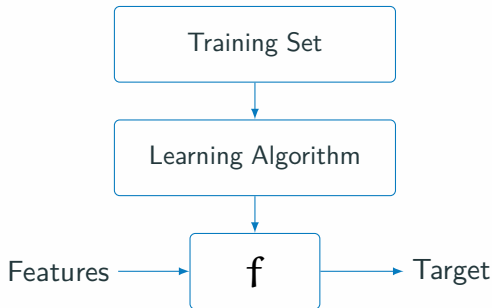
Session 6:

In general, any machine learning problem can be assigned to one of these broad types:



Supervised Learning

The term **supervised learning** refers to the fact that we gave the algorithm a data set in which the “**right answers**” (known as **labels**) were given.



- ▶ Supervised Learning refers to a set of approaches for **estimating f** .
- ▶ f is also called ***hypothesis*** in Machine Learning.

Regression

- ▶ The example of the house price prediction is also called a **regression** problem.
- ▶ A regression problem is when we try to predict a **quantitative (continuous)** value output. Namely the price in the example.

Classification

- ▶ The process for predicting **qualitative (categorical, discrete)** responses is known as classification.
- ▶ Methods: Logistic regression, Support Vector Machines, etc..

Classification

- ▶ Email: Spam / Not Spam?
- ▶ Online Transactions: Fraudulent (Yes/No)?
- ▶ Tumor: Malignant / Benign?
- ▶ Loan Demand (Credit Risk): Safe / Risky

Classification: categorical output

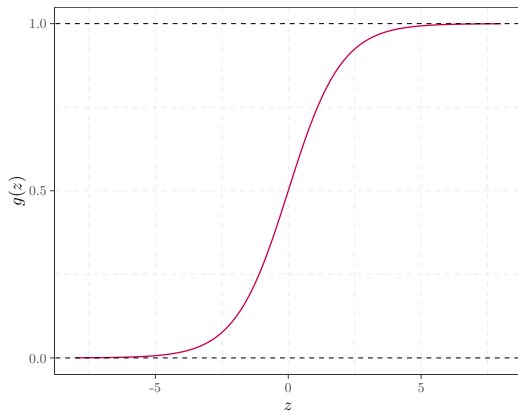
- ▶ $y \in \{0, 1\}$
- ▶ 0: "Negative class"
- ▶ 1: "Positive Class"

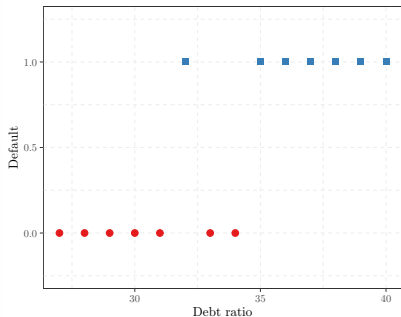
.. and also multiclass classification

Logistic Regression

The logistic function (sigmoid)

$$g(z) = \frac{e^z}{1 + e^z} = \frac{1}{1 + e^{-z}}$$





► $y \in \{0, 1\}$:

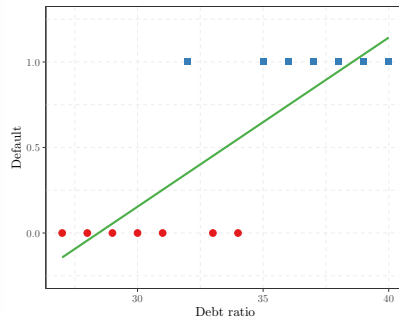
- "0": Negative class (here **no default**)
- "1": Positive class (here **default**)

► $f_{\omega}(x) = \omega'x$ can be > 1 ou < 0 !

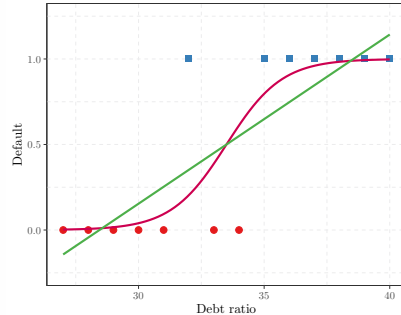
► Ideally $0 \leq f_{\omega}(x) \leq 1$ s.t.:

- If $f_{\omega}(x) \geq 0.5$, predict " $y = 1$ "
- If $f_{\omega}(x) < 0.5$, predict " $y = 0$ "

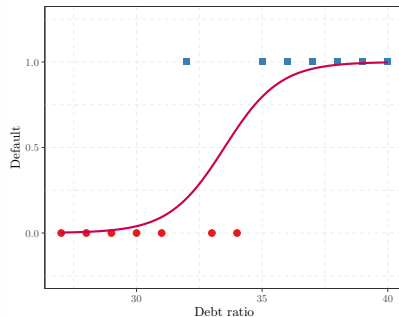
► Let $f_{\omega}(x) = \omega'x$



► Let $f_{\omega}(x) = \cancel{\omega'x} = g(\omega'x) = \frac{1}{1 + e^{-\omega'x}}$



- ▶ $0 \leq g(\omega'x) \leq 1$
- ▶ $f_{\omega}(x) = g(\omega'x)$ = estimated probability that $y = 1$ on input x
- ▶ Probability that $y = 1$, given x , parameterized by ω
- ▶ $g(\omega'x) = p(y = 1 | x) = p(x)$
- ▶ $y \in \{0, 1\}$ so $p(y = 1 | x) + p(y = 0 | x) = 1$



logistic score

$$p(x) = p(y = 1 | x) = \frac{e^{\omega'x}}{1 + e^{\omega'x}} = \frac{1}{1 + e^{-\omega'x}}$$

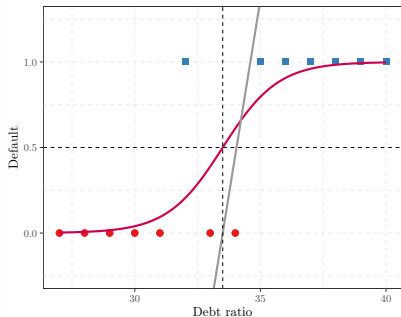
odds (côtes)

$$\frac{p(x)}{1 - p(x)} = e^{\omega'x}$$

log-odds or logit (logarithme des côtes)

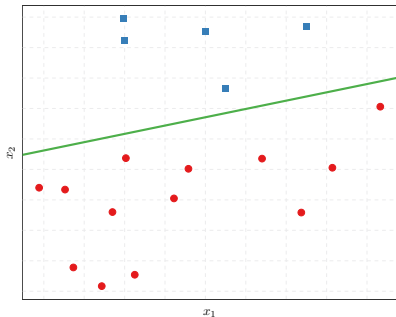
$$\log\left(\frac{p(x)}{1 - p(x)}\right) = \omega'x$$

Logistic Regression: decision boundary



► We predict " $y = 1$ " if $p(x) \geq 0.5$ which means $\omega'x \geq 0$

► $\omega_0 + \omega_1 x \geq 0 \Rightarrow x \geq -\frac{\omega_0}{\omega_1}$



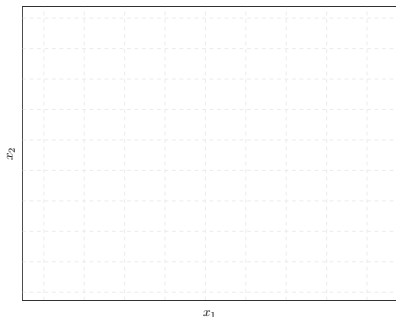
- ▶ $p(x) = p(y = 1 | x) = f_{\omega}(x) = g(\omega'x)$
- ▶ Predict " $y = 1$ " if $p(x) \geq 0.5$ which means $\omega'x \geq 0$

- ▶ $\omega_0 + \omega_1 x_1 + \omega_2 x_2 \geq 0$ So

$$x_2 \geq -\frac{\omega_1}{\omega_2} x_1 - \frac{\omega_0}{\omega_2}$$

Fun

Identify TP, TN, FP, FN on the figure.



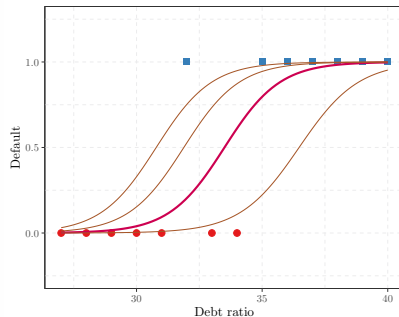
- ▶ Let

$$f_{\omega}(x) = g(\omega_0 + \omega_1 x_1 + \omega_2 x_2 + \omega_3 x_1^2 + \omega_4 x_2^2)$$
- ▶ For example, predict “ $y = 1$ ” if

$$-1 + x_1^2 + x_2^2 \geq 0$$
- ▶ Or, $f_{\omega}(x) = g(\omega_0 + \omega_1 x_1 + \omega_2 x_2 + \omega_3 x_1^2 + \omega_4 x_1^2 x_2 + \omega_5 x_1^2 x_2^2 + \dots)$

Logistic Regression: model estimation

- Parameters to estimate: $\omega = \{\omega_0, \omega_1\}$ if univariate
- $\omega = \{\omega_0, \omega_1, \dots, \omega_p\}$ if multivariate with p features
- How to choose parameters ω ?



¹check: <https://shinyserve.es/shiny/log-maximum-likelihood/>, by Eduardo García Portugués

Cost function of simple linear regression

- ▶ Model: $f_{\omega}(x) = \omega_0 + \omega_1 x = \omega'x$
- ▶ Parameters: ω_0 and ω_1
- ▶ Cost function: $J(\omega_0, \omega_1) = \frac{1}{n} \sum_{i=1}^n \frac{1}{2} (f_{\omega}(x^{(i)}) - y^{(i)})^2$
- ▶ Goal: $\min_{\omega_0, \omega_1} J(\omega_0, \omega_1)$

Non-convex in case of logistic regression !

- ▶ How to choose parameters ω ?
- ▶ $y \in \{0, 1\}$, Let's assume:

$$p(y = 1 \mid x, \omega) = f_{\omega}(x)$$

$$p(y = 0 \mid x, \omega) = 1 - f_{\omega}(x)$$

- ▶ We represent $y \mid x, \omega \sim \mathcal{B}(f_{\omega}(x))$
- ▶ We can write:

$$p(y \mid x, \omega) = (f_{\omega}(x))^y (1 - f_{\omega}(x))^{1-y} \quad y \in \{0, 1\}$$

- ▶ Given the n observations and assuming independance, we estimate ω by maximizing the **likelihood**:

$$\mathcal{L}(\omega) = \prod_{i=1}^n p(y^{(i)} \mid x^{(i)}, \omega)$$

- ▶ The **likelihood**:

$$\begin{aligned}\mathcal{L}(\omega) &= \prod_{i=1}^n p(y^{(i)} | x^{(i)}, \omega) \\ &= \prod_{i=1}^n (f_{\omega}(x^{(i)}))^{y^{(i)}} (1 - f_{\omega}(x^{(i)}))^{1-y^{(i)}}\end{aligned}$$

- ▶ Maximizing the likelihood is same as maximizing its log:

$$\begin{aligned}\ell(\omega) &= \log(\mathcal{L}(\omega)) \\ &= \sum_{i=1}^n y^{(i)} \log f_{\omega}(x^{(i)}) + (1 - y^{(i)}) \log(1 - f_{\omega}(x^{(i)}))\end{aligned}$$

- ▶ Maximizing $\ell(\omega)$ is same as minimizing: $-\frac{1}{n}\ell(\omega)$
- ▶ Let $J(\omega) = -\frac{1}{n}\ell(\omega)$, a **convex cost function** for the logistic regression model (known as *binary cross entropy*).

- ▶ **Goal:** Find ω s.t. $\omega = \operatorname{argmin}_{\omega} J(\omega)$
- ▶ $J(\omega) = -\frac{1}{n} \sum_{i=1}^n y^{(i)} \log f_{\omega}(x^{(i)}) + (1 - y^{(i)}) \log (1 - f_{\omega}(x^{(i)}))$
- ▶ Contrary to the linear regression, this cost function **does not** have an [analytical](#) solution. We need an optimization technique.

GD for logistic regression

- ▶ initialize ω 'randomly'
- ▶ repeat until convergence{

$$\omega_i^{\text{new}} = \omega_i^{\text{old}} - \alpha \frac{\partial J(\omega)}{\partial \omega_i}$$

simultaneously for $i = 0, \dots, p$ }

- ▶ Recall that $g(z) = \frac{e^z}{1 + e^z} = \frac{1}{1 + e^{-z}}$
- ▶ Notice that $g'(z) = g(z)(1 - g(z))$
- ▶ $\frac{\partial J(\omega)}{\partial \omega_i} = (y - f_\omega(x))x_i$

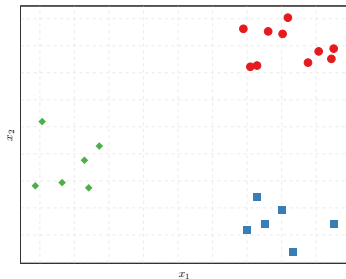
GD for logistic regression

- ▶ initialize ω randomly
- ▶ repeat until convergence{

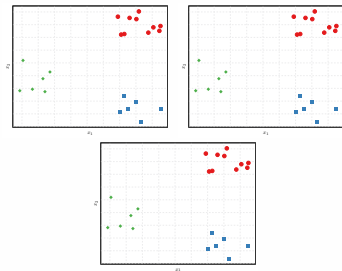
$$\omega_i^{\text{new}} = \omega_i^{\text{old}} - \alpha \frac{1}{n} \sum_{i=1}^n (f_\omega(x^{(i)}) - y^{(i)}) \cdot x_i^{(i)}$$

simultaneously for $i = 0, \dots, p$ }

- Weather: Sunny, Cloudy, Rain, Snow
- Medical diagrams: Not ill, Cold, Flu
- News articles: Sport, Education, Technology, Politics



- $f_{\omega}^{(i)}(x) = P(y = i|x, \omega)$ for $i = 1, 2, 3$
- Train a logistic regression classifier for each class i to predict the probability that $y = i$
- On a new input x , to make a prediction, pick the class i that maximizes $f_{\omega}^{(i)}(x)$



- ▶ Very famous method and maybe the most used
- ▶ Adapted for a binary y
- ▶ Relation with linear regression
- ▶ Linear decision boundary, but can be non linear using other hypothesis
- ▶ Direct calculation of $p(y = 1 | x)$

Disciminant Analysis

- ▶ L'analyse discriminante est une famille méthode de **classification** qui cherche à prédire avec quelle **probabilité** un individu appartient à une classe
- ▶ Au lieu de calculer directement $p(y | x)$, comme dans la régression logistique, on **modélise** $p(x | y)$ et $p(x)$
- ▶ Ensuite, on applique la formule de **Bayes** pour calculer $p(y | x)$
- ▶ Les méthodes présentées dans ce chapitre sont appelées méthodes génératives

Bayes

$$p(y | x) = \frac{p(y)p(x | y)}{p(x)}$$

- ▶ $p(y | x)$ probabilité "a posteriori"
- ▶ $p(y)$ probabilité "a priori"
- ▶ $p(x | y)$ distribution dans les classes
- ▶ $p(x)$ vraisemblance

Pour prédire la classe associée à une observation x :

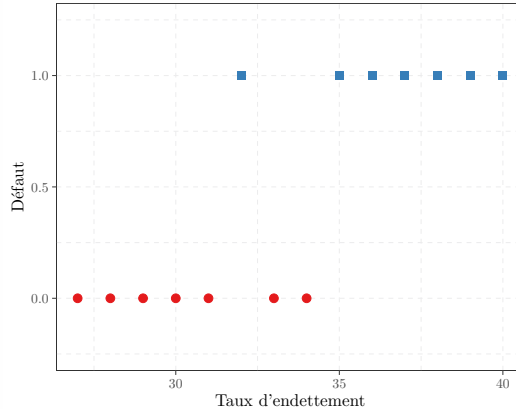
$$\begin{aligned}\arg \max_y p(y|x) &= \arg \max_y \frac{p(x|y)p(y)}{p(x)} \\ &= \arg \max_y p(x|y)p(y)\end{aligned}$$

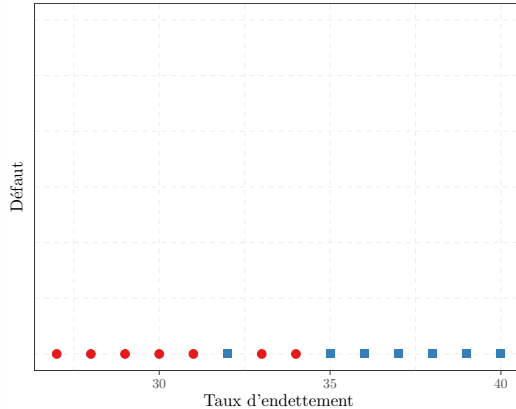
$$p(y | x) = \frac{p(y)p(x | y)}{p(x)}$$

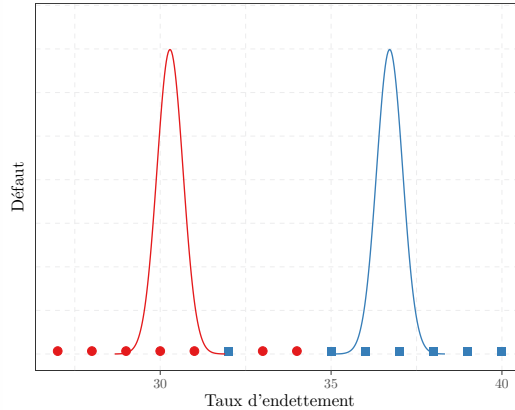
Soit

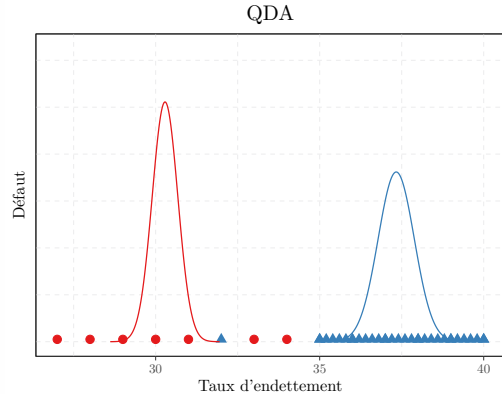
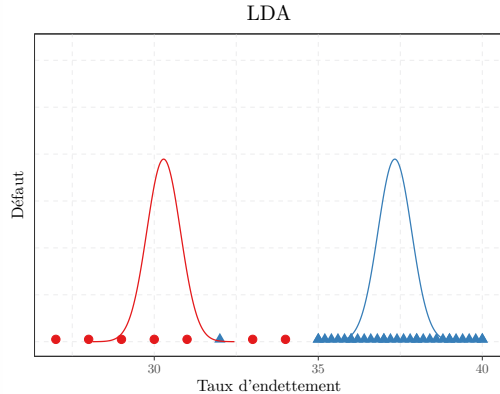
- ▶ $y \sim \mathcal{B}(\phi)$ donc $p(y) = \phi^y(1 - \phi)^{1-y}$
 - ▶ $x | y = 0 \sim \mathcal{N}(\mu_0, \sigma_0^2)$
 - ▶ $x | y = 1 \sim \mathcal{N}(\mu_1, \sigma_1^2)$
-
- ▶ Si $\sigma_0 = \sigma_1 \Rightarrow$ analyse discriminante linéaire (LDA)
 - ▶ Si $\sigma_0 \neq \sigma_1 \Rightarrow$ analyse discriminante quadratique (QDA)

Analyse discriminante: intuition ($p = 1$)









$$p(y | x) = \frac{p(y)p(x | y)}{p(x)}$$

Soit

- ▶ $y \sim \mathcal{B}(\phi)$ donc $p(y) = \phi^y(1 - \phi)^{1-y}$
- ▶ $x | y = 0 \sim \mathcal{N}(\mu_0, \Sigma_0)$
- ▶ $x | y = 1 \sim \mathcal{N}(\mu_1, \Sigma_1)$

$$p(y | x) = \frac{p(y)p(x | y)}{p(x)}$$

Soit

- ▶ $y \sim \mathcal{B}(\phi)$ donc $p(y) = \phi^y(1 - \phi)^{1-y}$
- ▶ $x | y = 0 \sim \mathcal{N}(\mu_0, \Sigma_0)$
- ▶ $x | y = 1 \sim \mathcal{N}(\mu_1, \Sigma_1)$

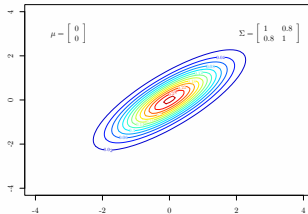
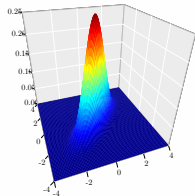
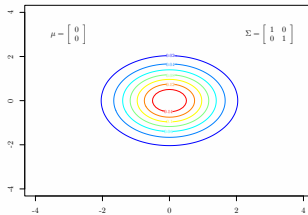
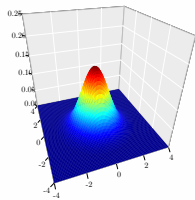
Ici, les gaussiennes sont multi-dimensionnelles. La densité de $\mathcal{N}(\mu, \Sigma)$ se définit comme suit:

$$p(x; \mu, \Sigma) = \frac{1}{(2\pi)^{d/2}|\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(x - \mu)^T \Sigma^{-1}(x - \mu)\right)$$

Où :

- ▶ μ est le vecteur des moyennes, $\mu \in \mathbb{R}^d$
- ▶ Σ est la matrice de covariance, $\Sigma \in \mathbb{R}^{d \times d}$ avec Σ une matrice semi-définie positive

La loi Normale multi-dimensionnelle: exemples ($d = 2$)



- Les paramètres à estimer sont: ϕ , μ , et Σ

$$\phi = \begin{pmatrix} \phi_0 \\ \phi_1 \\ \vdots \\ \phi_k \end{pmatrix} \quad \mu = \begin{pmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_p \end{pmatrix}$$

$$\Sigma = \text{Cov}(x) = \begin{pmatrix} \sigma_1^2 & \text{Cov}(x_1, x_2) & \dots & \text{Cov}(x_1, x_p) \\ \text{Cov}(x_2, x_1) & \sigma_2^2 & \dots & \text{Cov}(x_2, x_p) \\ \vdots & \vdots & \ddots & \vdots \\ \text{Cov}(x_p, x_1) & \text{Cov}(x_p, x_2) & \dots & \sigma_p^2 \end{pmatrix}$$

- ▶ Étant donné les n observations $\{(x^{(i)}, y^{(i)})\}$ et en assumant qu'elles sont indépendantes, on estime ϕ , μ , et Σ qui maximisent la **vraisemblance**:

$$\begin{aligned}
 \mathcal{L}(\phi, \mu_0, \mu_1, \Sigma) &= \prod_{i=1}^n p(x^{(i)}, y^{(i)}; \phi, \mu_0, \mu_1, \Sigma) \\
 &= \prod_{i=1}^n p(y^{(i)}; \phi) p(x^{(i)} | y^{(i)}; \mu_0, \mu_1, \Sigma) \\
 &= \prod_{i=1}^n \phi^{y^{(i)}} (1 - \phi)^{1-y^{(i)}} \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \\
 &\quad \exp\left(-\frac{1}{2} (x^{(i)} - \mu_{y^{(i)}})^T \Sigma^{-1} (x^{(i)} - \mu_{y^{(i)}})\right)
 \end{aligned}$$

- Maximiser la vraisemblance revient à maximiser son log:

$$\begin{aligned}\ell(\phi, \mu, \Sigma) = & \sum_{i=1}^n \left(y^{(i)} \log \phi + (1 - y^{(i)}) \log(1 - \phi) - \frac{1}{2} \log(2\pi) \right. \\ & \left. - \frac{1}{2} \log |\Sigma| - \frac{1}{2} (x^{(i)} - \mu_{y^{(i)}})^T \Sigma^{-1} (x^{(i)} - \mu_{y^{(i)}}) \right)\end{aligned}$$

- Maximiser ℓ possède une solution analytique. On résout le système:

$$\frac{\partial \ell}{\partial \phi} = 0, \quad \frac{\partial \ell}{\partial \mu} = 0, \quad \text{et} \quad \frac{\partial \ell}{\partial \Sigma} = 0$$

- En maximisant la log vraisemblance ℓ , on trouve (pour LDA):

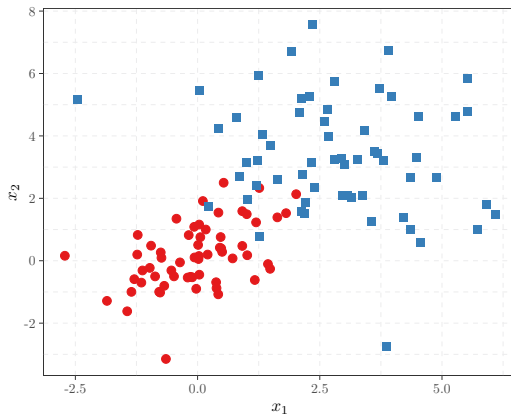
$$\begin{aligned}\phi &= \frac{1}{n} \sum_{i=1}^n 1_{\{y^{(i)}=1\}} \\ \mu_0 &= \frac{\sum_{i=1}^n 1_{\{y^{(i)}=0\}} x^{(i)}}{\sum_{i=1}^n 1_{\{y^{(i)}=0\}}} \text{ et } \mu_1 = \frac{\sum_{i=1}^n 1_{\{y^{(i)}=1\}} x^{(i)}}{\sum_{i=1}^n 1_{\{y^{(i)}=1\}}} \\ \Sigma &= \frac{1}{n} \sum_{i=1}^n \left(x^{(i)} - \mu_{y^{(i)}} \right) \left(x^{(i)} - \mu_{y^{(i)}} \right)^T\end{aligned}$$

- En maximisant la log vraisemblance ℓ , on trouve (pour LDA):

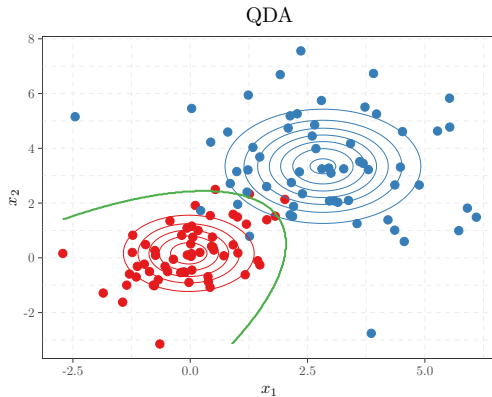
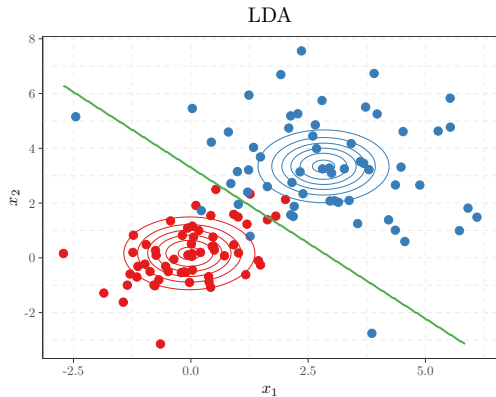
$$\begin{aligned}\phi &= \frac{1}{n} \sum_{i=1}^n 1_{\{y^{(i)}=1\}} \\ \mu_0 &= \frac{\sum_{i=1}^n 1_{\{y^{(i)}=0\}} x^{(i)}}{\sum_{i=1}^n 1_{\{y^{(i)}=0\}}} \text{ et } \mu_1 = \frac{\sum_{i=1}^n 1_{\{y^{(i)}=1\}} x^{(i)}}{\sum_{i=1}^n 1_{\{y^{(i)}=1\}}} \\ \Sigma &= \frac{1}{n} \sum_{i=1}^n \left(x^{(i)} - \mu_{y^{(i)}} \right) \left(x^{(i)} - \mu_{y^{(i)}} \right)^T\end{aligned}$$

Dans le cadre de la QDA:

$$\begin{aligned}\Sigma_0 &= \frac{1}{n_0} \sum_{i=1}^{n_0} (x^{(i)} - \mu_0)(x^{(i)} - \mu_0)^T 1_{\{y^{(i)}=0\}} \\ \Sigma_1 &= \frac{1}{n_1} \sum_{i=1}^{n_1} (x^{(i)} - \mu_1)(x^{(i)} - \mu_1)^T 1_{\{y^{(i)}=1\}}\end{aligned}$$



- Les frontières de séparations sont définies par $p(y = 0|x) = p(y = 1|x)$



- ▶ LDA et QDA méthodes génératives
- ▶ Relation entre LDA et la régression logistique

$$\left. \begin{array}{l} x | y = 0 \sim \mathcal{N}(\mu_0, \Sigma) \\ x | y = 1 \sim \mathcal{N}(\mu_1, \Sigma) \\ y \sim \mathcal{B}(\phi) \end{array} \right\} \implies p(y = 1 | x) = \frac{1}{1 + e^{-\beta'x}}$$

- ▶ Lorsque l'hypothèse de normalité est satisfaite, la méthode générative est préférée