

# Introduction to data science & artificial intelligence (INF7100)

Arthur Charpentier

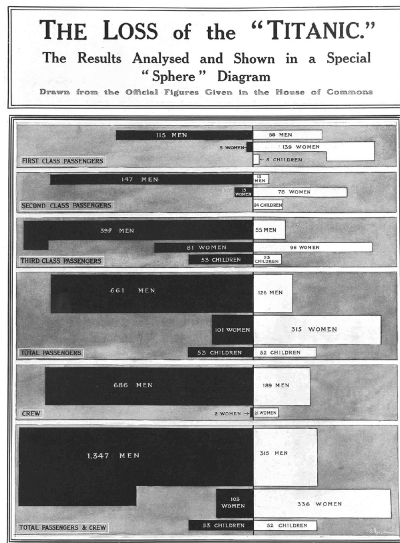
#351 Classification

été 2020

# Titanic, who survived?

$y \in \{0, 1\}$ , see

[freakonometrics.hypotheses.org](https://freakonometrics.hypotheses.org)



SPECIALLY DRAWN FOR "THE SPHERE" BY G. DECK  
The Black Indicates Passengers and Crew NOT SAVED, the White Indicates the SAVED

# Fisher's Discriminant Analysis

Suppose that we only want to predict the class,  $\hat{y} \in \{0, 1\}$  (or more generally  $\hat{y} \in \{a_1, a_2, \dots, a_J\}$ )

$$m^*(\mathbf{x}) = \operatorname{argmin}_{y \in \{0,1\}} \{\mathbb{P}[Y = y | \mathbf{X} = \mathbf{x}]\}$$

i.e.

$$m^*(\mathbf{x}) = \operatorname{argmin}_{y \in \{0,1\}} \left\{ \frac{\mathbb{P}[\mathbf{X} = \mathbf{x} | Y = y]}{\mathbb{P}[\mathbf{X} = \mathbf{x}]} \right\}$$

(where  $\mathbb{P}[\mathbf{X} = \mathbf{x}]$  is  $f(\mathbf{x})$  is the continuous case).

If  $y$  takes two values – i.e.  $\{0, 1\}$

$$m^*(\mathbf{x}) = \begin{cases} 1 & \text{if } \mathbb{E}(Y | \mathbf{X} = \mathbf{x}) > \frac{1}{2} \\ 0 & \text{otherwise} \end{cases}$$

# Fisher's Discriminant Analysis

The set

$$\mathcal{D}_S = \left\{ \mathbf{x} : \mathbb{E}(Y|\mathbf{X} = \mathbf{x}) = \frac{1}{2} \right\}$$

is called **decision border**.

Assume that  $\mathbf{X}|Y = 0 \sim \mathcal{N}(\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0)$  and  $\mathbf{X}|Y = 1 \sim \mathcal{N}(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1)$ , then, if  $r_y^2$  is Mahalanobis distance from  $\mathbf{x}$  to  $\boldsymbol{\mu}_y$

$$r_y^2 = [\mathbf{x} - \boldsymbol{\mu}_y]^\top \boldsymbol{\Sigma}_y^{-1} [\mathbf{x} - \boldsymbol{\mu}_y] \text{ for } y \in \{0, 1\},$$

$$m^*(\mathbf{x}) = \begin{cases} 1 & \text{if } r_1^2 < r_0^2 + 2 \log \frac{\mathbb{P}(Y = 1)}{\mathbb{P}(Y = 0)} + \log \frac{|\boldsymbol{\Sigma}_0|}{|\boldsymbol{\Sigma}_1|} \\ 0 & \text{otherwise} \end{cases}$$

# Fisher's Discriminant Analysis

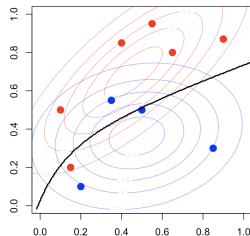
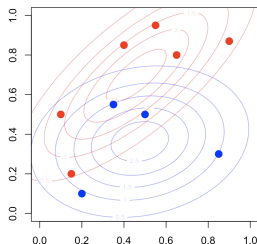
Let  $\delta_y$  be defined (for  $y \in \{0, 1\}$ ) as

$$\delta_y(\mathbf{x}) = -\frac{1}{2} \log |\boldsymbol{\Sigma}_y| - \frac{1}{2} [\mathbf{x} - \boldsymbol{\mu}_y]^\top \boldsymbol{\Sigma}_y^{-1} [\mathbf{x} - \boldsymbol{\mu}_y] + \log \mathbb{P}(Y = y)$$

so that the decision frontier is

$$\{\mathbf{x} \text{ such that } \delta_0(\mathbf{x}) = \delta_1(\mathbf{x})\}$$

which is quadratic in  $\mathbf{x}$

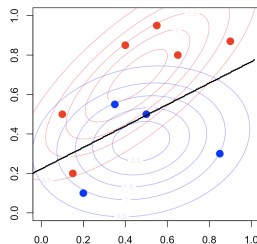
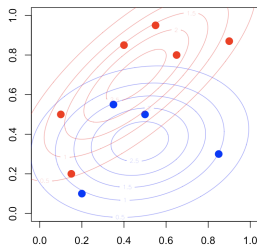


# Fisher's Discriminant Analysis

Fisher (1936) added the assumption  $\Sigma_0 = \Sigma_1$ . Then

$$\delta_y(\mathbf{x}) = \mathbf{x}^\top \Sigma^{-1} \mu_y - \frac{1}{2} \mu_y^\top \Sigma^{-1} \mu_y + \log \mathbb{P}(Y = y)$$

and the decision border is linear in  $\mathbf{x}$



# Fisher's Discriminant Analysis

If  $\mathbf{X}|Y = 0 \sim \mathcal{N}(\boldsymbol{\mu}_0, \boldsymbol{\Sigma})$  and  $\mathbf{X}|Y = 1 \sim \mathcal{N}(\boldsymbol{\mu}_1, \boldsymbol{\Sigma})$  then

$$\log \frac{\mathbb{P}(Y = 1|\mathbf{X} = \mathbf{x})}{\mathbb{P}(Y = 0|\mathbf{X} = \mathbf{x})}$$

is equal to

$$\mathbf{x}^\top \boldsymbol{\Sigma}^{-1}[\boldsymbol{\mu}_y] - \frac{1}{2}[\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0]^\top \boldsymbol{\Sigma}^{-1}[\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0] + \log \frac{\mathbb{P}(Y = 1)}{\mathbb{P}(Y = 0)}$$

which is linear in  $\mathbf{x}$ , i.e.

$$\log \frac{\mathbb{P}(Y = 1|\mathbf{X} = \mathbf{x})}{\mathbb{P}(Y = 0|\mathbf{X} = \mathbf{x})} = \mathbf{x}^\top \boldsymbol{\beta}$$

which will be close to the linear regression...

# Logistic Regression

$$p_i = \mathbb{E}(Y_i | \mathbf{X}_i = \mathbf{x}_i) \in [0, 1] \neq \mathbf{x}_i^\top \boldsymbol{\beta}$$

→ use the **odds ratio**

$$\text{odds}_i = \frac{\mathbb{P}[Y_i = 1]}{\mathbb{P}[Y_i = 0]} = \frac{p_i}{1 - p_i} \in [0, \infty].$$

i.e. if we take the logarithm

$$\log(\text{odds}_i) = \log\left(\frac{p_i}{1 - p_i}\right) \in \mathbb{R}.$$

That transformation is called **logit**,

$$\text{logit}(p_i) = \log\left(\frac{p_i}{1 - p_i}\right) = \mathbf{x}_i^\top \boldsymbol{\beta}$$

or

$$p_i = \text{logit}^{-1}(\mathbf{x}_i^\top \boldsymbol{\beta}) = \frac{\exp[\mathbf{x}_i^\top \boldsymbol{\beta}]}{1 + \exp[\mathbf{x}_i^\top \boldsymbol{\beta}]}.$$



# Logistic Regression

Data :  $\{(\mathbf{x}_i, y_i) = (x_{1,i}, x_{2,i}, y_i), i = 1, \dots, n\}$

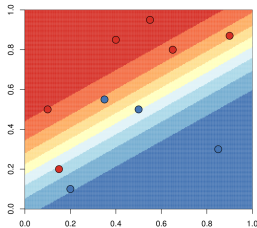
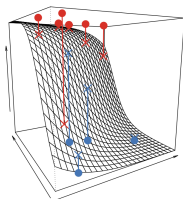
$$\mathbb{P}(Y = 1 | \mathbf{X} = \mathbf{x}) = \frac{\exp[\mathbf{x}^T \boldsymbol{\beta}]}{1 + \exp[\mathbf{x}^T \boldsymbol{\beta}]}$$

Inference using maximum likelihood techniques

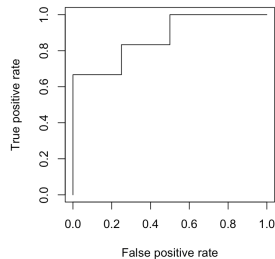
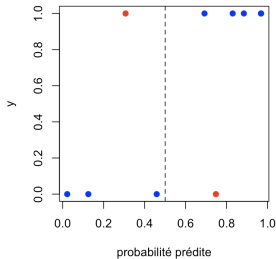
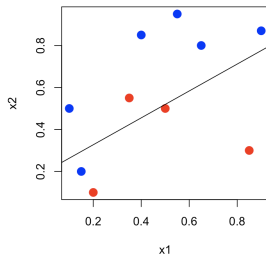
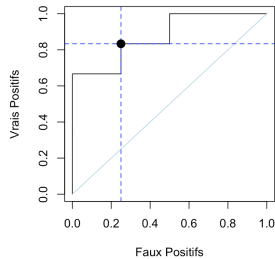
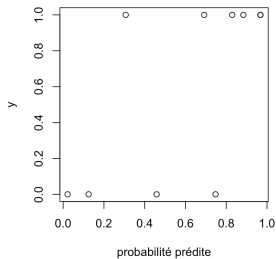
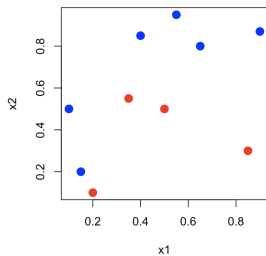
$$\hat{\boldsymbol{\beta}} = \operatorname{argmin} \left\{ \sum_{i=1}^n \log[\mathbb{P}(Y = y_i | \mathbf{X} = \mathbf{x}_i)] \right\}$$

and the score model is then

$$s(\mathbf{x}) = \frac{\exp[\mathbf{x}^T \hat{\boldsymbol{\beta}}]}{1 + \exp[\mathbf{x}^T \hat{\boldsymbol{\beta}}]}$$



# ROC Curve



## ROC Curve

$$FPR = \frac{\mathbb{P}[y = 0, \hat{y} = 1]}{\mathbb{P}[y = 0]} \text{ et } TPR = \frac{\mathbb{P}[y = 1, \hat{y} = 1]}{\mathbb{P}[y = 1]}$$