

## Introduction à l'apprentissage supervisé

Y. Pradat (*auteur*), L. Verlingue & D. Gautheret (*responsables cours*)

CentraleSupélec (labo MICS) & Institut Gustave Roussy (IGR)

23 janvier 2020



- 1 Notations et méthodes
- 2 Régression linéaire
- 3 Régression logistique
- 4 Modèles à risques proportionnels de Cox
- 5 Conclusion

# Notations

## Notations

- $n, p \in \mathbb{N}^*$  nombre d'observations (unités, individus) et nombre de variables (facteurs, cofacteurs);
- $X$  (resp  $\mathbf{X}$ ) variable aléatoire scalaire (resp vectorielle);
- $\mathbf{X}_{1:n}$  vecteur (resp matrice) des variables aléatoires  $X_1, \dots, X_n$  (resp  $\mathbf{X}_1, \dots, \mathbf{X}_n$ );
- $x$  (resp  $\mathbf{x}$ ) observation de  $X$  (resp  $\mathbf{X}$ );
- $\mathbf{x}_{1:n}$  observation de  $\mathbf{X}_{1:n}$ .

## Objectif :

- 1 Proposer un modèle mathématique;
- 2 Estimer les paramètres du modèle.

# Definition vraisemblance

## Définition 1. Modèle statistique

*Un modèle statistique est une collection de lois (ou densités) candidates paramétrée par  $\theta \in \Theta$*

$$\mathcal{M}_{\Theta} = \{p_{\theta} | \theta \in \Theta\} \quad (1)$$

## Définition 2. Vraisemblance

*Soient  $\mathbf{x}_{1:n}$  un échantillon d'observations de  $\mathbf{X}_{1:n} \sim p_{\theta^*}$ . La vraisemblance du paramètre  $\theta$  pour cet échantillon est*

$$\mathcal{L}(\theta; \mathbf{x}_{1:n}) = \prod_{i=1}^n p_{\theta}(x_i) \quad (2)$$

# Definition vraisemblance

## Définition 1. Modèle statistique

Un modèle statistique est une collection de lois (ou densités) candidates paramétrée par  $\theta \in \Theta$

$$\mathcal{M}_\Theta = \{p_\theta | \theta \in \Theta\} \quad (1)$$

## Définition 2. Vraisemblance

Soient  $\mathbf{x}_{1:n}$  un échantillon d'observations de  $\mathbf{X}_{1:n} \sim p_{\theta^*}$ . La vraisemblance du paramètre  $\theta$  pour cet échantillon est

$$\mathcal{L}(\theta; \mathbf{x}_{1:n}) = \prod_{i=1}^n p_\theta(x_i) \quad (2)$$

## Exemple 1

Soit  $\mathbf{x}_{1:n}$  un  $n$ -échantillon de  $\mathbf{X}_{1:n} \sim \mathcal{N}(\mu^*, \sigma^{2*})$ . Alors, la vraisemblance de  $(\mu, \sigma^2)$  est

$$\mathcal{L}(\mu, \sigma^2; \mathbf{x}_{1:n}) = \prod_{i=1}^n \frac{1}{\sqrt{\sigma^2 2\pi}} e^{-\frac{1}{2\sigma^2}(x_i - \mu)^2} \quad (3)$$

# Estimation par MV

## Définition 3. Estimateur du MV

Pour un échantillon  $\mathbf{x}_{1:n}$  de  $\mathbf{X}_{1:n} \sim p_{\theta^*}$ , l'équation

$$\hat{\theta}^*(\mathbf{x}_{1:n}) \in \operatorname{argmax}_{\theta \in \Theta} \mathcal{L}(\theta; \mathbf{x}_{1:n}) \quad (4)$$

définit un estimateur du maximum de vraisemblance de  $\theta^*$ .

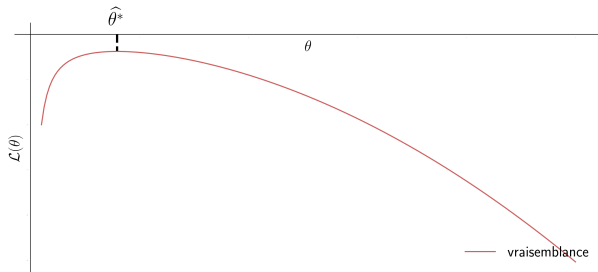
# Estimation par MV

## Définition 3. Estimateur du MV

Pour un échantillon  $\mathbf{x}_{1:n}$  de  $\mathbf{X}_{1:n} \sim p_{\theta^*}$ , l'équation

$$\hat{\theta}^*(\mathbf{x}_{1:n}) \in \operatorname{argmax}_{\theta \in \Theta} \mathcal{L}(\theta; \mathbf{x}_{1:n}) \quad (4)$$

définit un estimateur du maximum de vraisemblance de  $\theta^*$ .



# Exemple estimation par MV

## Définition 3. Estimateur du MV

Pour un échantillon  $\mathbf{x}_{1:n}$  de  $\mathbf{X}_{1:n} \sim p_{\theta^*}$ , l'équation

$$\hat{\theta}^*(\mathbf{x}_{1:n}) \in \operatorname{argmax}_{\theta \in \Theta} \mathcal{L}(\theta; \mathbf{x}_{1:n}) \quad (5)$$

définit un estimateur du maximum de vraisemblance de  $\theta^*$ .

## Exemple 2

Soit  $\mathbf{x}_{1:100}$  un 100-échantillon de  $\mathbf{X}_{1:100} \sim \mathcal{B}(\theta^*)$  (100 lancers de pièces identiques).  
L'estimateur par MV de  $\theta$  maximise

$$\prod_{i=1}^{100} p_{X_i}(x_i; \theta) \quad (6)$$



# Exemple estimation par MV

## Définition 3. Estimateur du MV

Pour un échantillon  $\mathbf{x}_{1:n}$  de  $\mathbf{X}_{1:n} \sim p_{\theta^*}$ , l'équation

$$\hat{\theta}^*(\mathbf{x}_{1:n}) \in \operatorname{argmax}_{\theta \in \Theta} \mathcal{L}(\theta; \mathbf{x}_{1:n}) \quad (5)$$

définit un estimateur du maximum de vraisemblance de  $\theta^*$ .

## Exemple 2

Soit  $\mathbf{x}_{1:100}$  un 100-échantillon de  $\mathbf{X}_{1:100} \sim \mathcal{B}(\theta^*)$  (100 lancers de pièces identiques).  
L'estimateur par MV de  $\theta$  maximise

$$\prod_{i=1}^{100} p_{X_i}(x_i; \theta) \quad (6)$$

Pour une loi de Bernoulli,  $p(x; \theta) = \theta^x (1 - \theta)^{1-x}$  de sorte que

$$\hat{\theta}^*(\mathbf{x}_{1:100}) = \frac{1}{100} \sum_{i=1}^{100} x_i \quad (7)$$

- 1 Notations et méthodes
- 2 Régression linéaire
- 3 Régression logistique
- 4 Modèles à risques proportionnels de Cox
- 5 Conclusion

## Construction modèle

Le modèle Soit  $\mathbf{x}_{1:n} \in (\mathbb{R}^p)^n$  un  $n$ -échantillon de  $\mathbf{X}_{1:n}$  et soient, pour  $\mathbf{x}_{1:n}$  fixé, les variables aléatoires  $Y_1, \dots, Y_n$  données par

$$\forall i = 1, \dots, n, \quad Y_i = \mathbf{x}_i^\top \beta + \mathcal{E}_i \quad (8)$$

avec  $\mathcal{E}_i \sim \mathcal{N}(0, \sigma^2)$  (résidus).

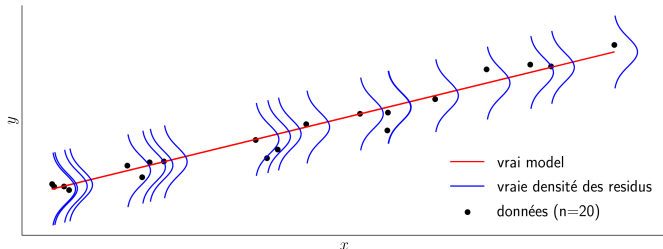
# Construction modèle

Le modèle Soit  $\mathbf{x}_{1:n} \in (\mathbb{R}^p)^n$  un  $n$ -échantillon de  $\mathbf{X}_{1:n}$  et soient, pour  $\mathbf{x}_{1:n}$  fixé, les variables aléatoires  $Y_1, \dots, Y_n$  données par

$$\forall i = 1, \dots, n, \quad Y_i = \mathbf{x}_i^\top \beta + \mathcal{E}_i \quad (8)$$

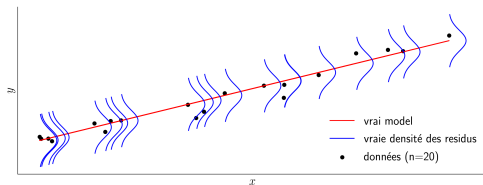
avec  $\mathcal{E}_i \sim \mathcal{N}(0, \sigma^2)$  (résidus).

## Explication



# Construction modèle

## Explication



- Les paramètres du modèle sont  $(\beta, \sigma^2)$ .
- C'est un modèle statistique  $\mathcal{M}_{\beta, \sigma^2}$  sur les densités conditionnelles  $p_{Y_i | \mathbf{X}_i = x_i}$ .
- **Hypothèses du modèle**
  - 1 Relation linéaire entre  $\mathbf{x}_i$  et  $Y_i$  ;
  - 2 Résidus indépendants et gaussiens ;
  - 3 Homoscedasticité.

# Estimation paramètre $\beta$

Par maximum de vraisemblance,

$$\mathcal{L}(\beta, \sigma^2; \mathbf{x}_{1:n}, \mathbf{y}_{1:n}) = \prod_{i=1}^n p_{Y_i|X_i=\mathbf{x}_i}(y_i) \quad (9)$$

$$= \prod_{i=1}^n \frac{1}{\sqrt{\sigma^2 2\pi}} e^{-\frac{1}{2\sigma^2} (y_i - \mathbf{x}_i^\top \beta)^2} \quad (10)$$

Exercice : Calculer l'estimateur  $\hat{\beta}(\mathbf{x}_{1:n}, \mathbf{y}_{1:n})$  de  $\beta$  par MV en minimisant  $\ell(\beta, \sigma^2) = -\log \mathcal{L}(\beta, \sigma^2)$ .

# Estimation paramètre $\beta$

Par maximum de vraisemblance,

$$\mathcal{L}(\beta, \sigma^2; \mathbf{x}_{1:n}, \mathbf{y}_{1:n}) = \prod_{i=1}^n p_{Y_i|X_i=\mathbf{x}_i}(y_i) \quad (9)$$

$$= \prod_{i=1}^n \frac{1}{\sqrt{\sigma^2 2\pi}} e^{-\frac{1}{2\sigma^2} (y_i - \mathbf{x}_i^\top \beta)^2} \quad (10)$$

Exercice : Calculer l'estimateur  $\hat{\beta}(\mathbf{x}_{1:n}, \mathbf{y}_{1:n})$  de  $\beta$  par MV en minimisant  $\ell(\beta, \sigma^2) = -\log \mathcal{L}(\beta, \sigma^2)$ .

Tous calculs faits, si  $\mathbf{x}_{1:n}^\top \mathbf{x}_{1:n}$  inversible,

$$\hat{\beta}(\mathbf{x}_{1:n}, \mathbf{y}_{1:n}) = (\mathbf{x}_{1:n}^\top \mathbf{x}_{1:n})^{-1} \mathbf{x}_{1:n}^\top \mathbf{y}_{1:n} \quad (11)$$

- 1 Notations et méthodes
- 2 Régression linéaire
- 3 Régression logistique
- 4 Modèles à risques proportionnels de Cox
- 5 Conclusion



# Construction modèle

Le modèle Soit  $\mathbf{x}_{1:n} \in (\mathbb{R}^p)^n$  un  $n$ -échantillon de  $\mathbf{X}_{1:n}$  et soient, pour  $\mathbf{x}_{1:n}$  fixé, les variables aléatoires  $Y_1, \dots, Y_n$  données par

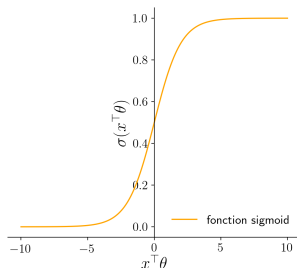
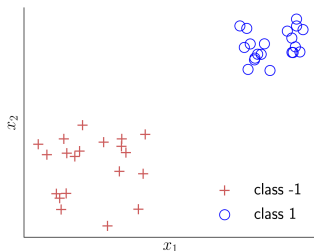
$$\forall i = 1, \dots, n, \quad Y_i \sim \mathcal{B}(\sigma(\mathbf{x}_i^\top \theta)) \quad (12)$$

# Construction modèle

Le modèle Soit  $\mathbf{x}_{1:n} \in (\mathbb{R}^p)^n$  un  $n$ -échantillon de  $\mathbf{X}_{1:n}$  et soient, pour  $\mathbf{x}_{1:n}$  fixé, les variables aléatoires  $Y_1, \dots, Y_n$  données par

$$\forall i = 1, \dots, n, \quad Y_i \sim \mathcal{B}(\sigma(\mathbf{x}_i^\top \theta)) \quad (12)$$

## Explication



# Estimation paramètre $\theta$

## Explications

- Les paramètres du modèle sont  $\theta$ .
- C'est un modèle statistique  $\mathcal{M}_\theta$  sur les densités conditionnelles  $p_{Y_i|X_i=x_i}$ .
- **Hypothèses du modèle**
  - 1 Relation linéaire entre  $\text{logit} p_{Y_i|X_i=x_i}(1)$  et  $\mathbf{x}_i$ ;

Estimation Par maximum de vraisemblance,

$$\mathcal{L}(\theta; \mathbf{x}_{1:n}, \mathbf{y}_{1:n}) = \prod_{i=1}^n p_{Y_i|X_i=x_i}(y_i) \quad (13)$$

$$= \prod_{i=1}^n \sigma(\theta^\top \mathbf{x}_i)^{y_i} (1 - \sigma(\theta^\top \mathbf{x}_i))^{1-y_i} \quad (14)$$

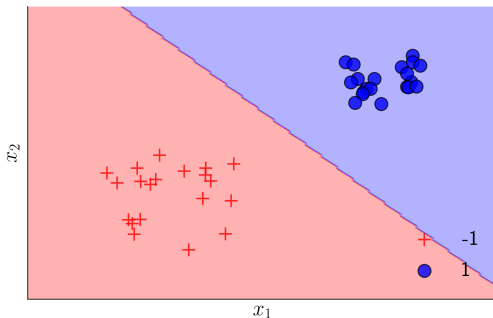
$$\ell(\theta; \mathbf{x}_{1:n}, \mathbf{y}_{1:n}) = - \sum_{i=1}^n y_i \log \sigma(\theta^\top \mathbf{x}_i) + (1 - y_i) \log(1 - \sigma(\theta^\top \mathbf{x}_i)) \quad (15)$$

## Estimation paramètre $\theta$

Comme  $\sigma'(z) = \sigma(z)(1 - \sigma(z))$  on obtient que :

$$\nabla \ell(\theta) = - \sum_{i=1}^n \mathbf{x}_i (y_i - \sigma(\mathbf{x}_i^\top \theta)) \quad (16)$$

qui nous sert à minimiser  $\ell(\theta)$  *numériquement*.



- 1 Notations et méthodes
- 2 Régression linéaire
- 3 Régression logistique
- 4 Modèles à risques proportionnels de Cox
- 5 Conclusion

# Modélisation de survie

## Notations

- ①  $T_D, T_C$  variable aléatoires positives, temps à l'évènement et temps à la censure respectivement ;
- ②  $\Delta$  variable aléatoire binaire indiquant l'occurrence de l'évènement ;
- ③  $T = \min(T_D, T_C)$  variable aléatoire observée ;
- ④  $\mathbf{Z}$  vecteur aléatoire des covariables ;

## Définition 4. Fonction de survie

*La fonction de survie est donnée par*

$$S: \begin{array}{ll} \mathbb{R}_+ & \rightarrow [0, 1] \\ t & \mapsto \mathbb{P}(T_D \geq t), \end{array}$$

# Modélisation de survie

## Définition 5. Taux de risque

*La taux de risque (instantané) est la fonction donnée par*

$$\lambda: \begin{array}{ccc} \mathbb{R}_+ & \rightarrow & \mathbb{R}_+ \\ t & \mapsto & \lim_{h \rightarrow 0} \mathbb{P}(T_D \leq t + h | T_D \geq t). \end{array}$$

# Modélisation de survie

## Définition 5. Taux de risque

La *taux de risque (instantané)* est la fonction donnée par

$$\lambda: \begin{array}{ccc} \mathbb{R}_+ & \rightarrow & \mathbb{R}_+ \\ t & \mapsto & \lim_{h \rightarrow 0} \mathbb{P}(T_D \leq t + h | T_D \geq t). \end{array}$$

Soit :  $f_D$  (resp  $F_D$ ) densité (resp f.r.) de  $T_D$ . Alors

$$\lambda(t) = \frac{f_D(t)}{1 - F_D(t)} \quad (17)$$

$$= \frac{-S'(t)}{S(t)} \quad (18)$$

$$\text{d'où} \quad \int_0^t \lambda(s) ds = -\log(S(t)) \quad (\text{car } S(0)=1) \quad (19)$$



## Construction modèle

Le modèle Soient  $\delta_{1:n}$ ,  $\mathbf{z}_{1:n}$   $n$ -échantillon de  $\Delta_{1:n}$ ,  $\mathbf{Z}_{1:n}$ . Le modèle de Cox modélise  $T_{D,1:n}$ , à  $\mathbf{z}_{1:n}$  fixés, via le taux de risque selon la relation

$$\lambda_i(t) = \lambda_0(t)e^{\mathbf{z}_i^\top \boldsymbol{\beta}} \quad (20)$$

# Construction modèle

Le modèle Soient  $\delta_{1:n}, \mathbf{z}_{1:n}$   $n$ -échantillon de  $\Delta_{1:n}, \mathbf{Z}_{1:n}$ . Le modèle de Cox modélise  $T_{D,1:n}$ , à  $\mathbf{z}_{1:n}$  fixés, via le taux de risque selon la relation

$$\lambda_i(t) = \lambda_0(t)e^{\mathbf{z}_i^\top \beta} \quad (20)$$

## Explications Hypothèses

- *Proportionnalité* Exercice : Pourquoi “risques proportionnels” ?

# Construction modèle

Le modèle Soient  $\delta_{1:n}, \mathbf{z}_{1:n}$   $n$ -échantillon de  $\Delta_{1:n}, \mathbf{Z}_{1:n}$ . Le modèle de Cox modélise  $T_{D,1:n}$ , à  $\mathbf{z}_{1:n}$  fixés, via le taux de risque selon la relation

$$\lambda_i(t) = \lambda_0(t)e^{\mathbf{z}_i^\top \beta} \quad (20)$$

## Explications Hypothèses

- *Proportionnalité* Exercice : Pourquoi “risques proportionnels” ?  
Supposons  $\mathbf{z} = \mathbf{1}$  pour le groupe traité et  $\mathbf{z} = \mathbf{0}$  pour le groupe de contrôle. Alors,

$$\forall t \geq 0, \quad \frac{\lambda(t, \mathbf{1})}{\lambda(t, \mathbf{0})} = \frac{\lambda_0(t)e^\beta}{\lambda_0(t)} = e^\beta \quad (21)$$

# Construction modèle

Le modèle Soient  $\delta_{1:n}$ ,  $\mathbf{z}_{1:n}$   $n$ -échantillon de  $\Delta_{1:n}$ ,  $\mathbf{Z}_{1:n}$ . Le modèle de Cox modélise  $T_{D,1:n}$ , à  $\mathbf{z}_{1:n}$  fixés, via le taux de risque selon la relation

$$\lambda_i(t) = \lambda_0(t)e^{\mathbf{z}_i^\top \beta} \quad (20)$$

## Explications Hypothèses

- *Proportionnalité* Exercice : Pourquoi “risques proportionnels” ?  
Supposons  $\mathbf{z} = 1$  pour le groupe traité et  $\mathbf{z} = 0$  pour le groupe de contrôle. Alors,

$$\forall t \geq 0, \quad \frac{\lambda(t, 1)}{\lambda(t, 0)} = \frac{\lambda_0(t)e^\beta}{\lambda_0(t)} = e^\beta \quad (21)$$

- *Linéarité* entre  $\log(\lambda_i)$  et les covariables  $\mathbf{z}_i$ .

# Vraisemblance totale

Par maximum de vraisemblance

- $\mathbb{P}(T_i = t_i | \Delta_i = 0) = \mathbb{P}(T_{D,i} \geq t_i) = \mathcal{S}_i(t_i)$
- $\mathbb{P}(T_i = t_i | \Delta_i = 1) = \mathbb{P}(T_{D,i} = t_i) = \lambda_i(t_i) \mathcal{S}_i(t_i)$

# Vraisemblance totale

Par maximum de vraisemblance

- $\mathbb{P}(T_i = t_i | \Delta_i = 0) = \mathbb{P}(T_{D,i} \geq t_i) = S_i(t_i)$
- $\mathbb{P}(T_i = t_i | \Delta_i = 1) = \mathbb{P}(T_{D,i} = t_i) = \lambda_i(t_i) S_i(t_i)$

Alors,

$$\begin{aligned}
 \mathcal{L}(\beta; \mathbf{z}_{1:n}, \mathbf{t}_{1:n}, \delta_{1:n}) &= \prod_{i=1}^n \lambda(t_i)^{\delta_i} S_i(t_i) \\
 &= \prod_{i=1}^n \left[ \frac{\lambda_i(t_i)}{\sum_{j \in R(t_i)} \lambda_j(t_j)} \right]^{\delta_i} \left[ \sum_{j \in R(t_i)} \lambda_j(t_i) \right]^{\delta_i} S_i(t_i) \\
 &= \prod_{i=1}^n \left[ \frac{\lambda_0(t_i) e^{\mathbf{z}_i^\top \beta}}{\sum_{j \in R(t_i)} \lambda_0(t_i) e^{\mathbf{z}_j^\top \beta}} \right]^{\delta_i} \left[ \sum_{j \in R(t_i)} \lambda_j(t_j) \right]^{\delta_i} S_i(t_i)
 \end{aligned}$$

# Vraisemblance partielle

On estime  $\beta$  par

$$\hat{\beta}(\mathbf{z}_{1:n}, \mathbf{t}_{1:n}, \delta_{1:n}) \in \operatorname{argmax}_{\beta} \mathcal{L}_{\text{partiel}}(\beta) = \prod_{i=1}^n \left[ \frac{\lambda_0(t_i) e^{\mathbf{z}_i^\top \beta}}{\sum_{j \in R(t_i)} \lambda_0(t_i) e^{\mathbf{z}_j^\top \beta}} \right]^{\delta_i} \quad (22)$$

# Exemple

individual	$X_i$	$\delta_i$	$Z_i$
1	9	1	4
2	8	0	5
3	6	1	7
4	10	1	3



## Exemple

individual	$X_i$	$\delta_i$	$Z_i$
1	9	1	4
2	8	0	5
3	6	1	7
4	10	1	3

ordered failure				Likelihood contribution	
$j$	time	$X_i$	$\mathcal{R}(X_i)$	$i_j$	$\left[e^{\beta Z_i} / \sum_{j \in \mathcal{R}(X_i)} e^{\beta Z_j}\right]^{\delta_i}$
1	6		{1,2,3,4}	3	$e^{7\beta} / [e^{4\beta} + e^{5\beta} + e^{7\beta} + e^{3\beta}]$
2	8		{1,2,4}	2	1
3	9		{1,4}	1	$e^{4\beta} / [e^{4\beta} + e^{3\beta}]$
4	10		{4}	4	$e^{3\beta} / e^{3\beta} = 1$

Nous avons vu 3 modèles différents :

- ➊ Régression linéaire, paramètres  $\beta, \sigma^2$ . Formule théorique pour  $\beta$  (et  $\sigma^2$ ).

Nous avons vu 3 modèles différents :

- ➊ Régression linéaire, paramètres  $\beta, \sigma^2$ . Formule théorique pour  $\beta$  (et  $\sigma^2$ ).
- ➋ Régression logistique, paramètres  $\theta$ . Estimation *numérique* (descente de gradient).

Nous avons vu 3 modèles différents :

- ➊ Régression linéaire, paramètres  $\beta, \sigma^2$ . Formule théorique pour  $\beta$  (et  $\sigma^2$ ).
- ➋ Régression logistique, paramètres  $\theta$ . Estimation *numérique* (descente de gradient).
- ➌ Modèle de Cox, paramètres  $\beta, \lambda_0(t)$ . Estimation *numérique* de  $\beta$  (descente de gradient).