

Actuariat de l'Assurance Non-Vie # 2

A. Charpentier (Université de Rennes 1)



ENSAE ParisTech, 2016/2017.

<http://freakonometrics.hypotheses.org>

Modélisation économétrique d'une variable dichotomique

Utile en fraude, marketing, crédit, catastrophes, décès ou pour modéliser la fréquence sur des branches spécifiques (max 1 sinistre par an).

Références: Frees (2010), chapitre 11 (p 305-342), Greene (2012), sections 17.2 et 17.3 (p 863-715), de Jong & Heller (2008), chapitre 7.

Remarque: la régression logistique est un cas particulier des modèles GLM, avec une loi **binomiale** et une fonction de lien **logit**.

Modèles introduits par Gaddum (1933), [Bliss \(1935\)](#), ou [Berkson \(1944\)](#).

Modélisation économétrique d'une variable dichotomique

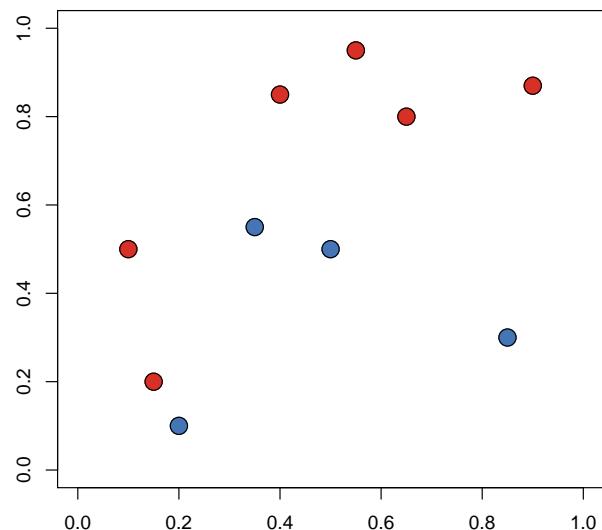
Ici $y_i \in \{0, 1\}$, mais pour les méthodes de machine learning $y_i \in \{-1, +1\}$, ou pour les dessins $y_i \in \{\bullet, \circ\}$.

La prévision se fera en deux étapes

- estimation d'une fonction de **score**,

$$s(\mathbf{x}) = \mathbb{P}(Y = 1 | \mathbf{X} = \mathbf{x}) \in [0, 1]$$

- construction d'un **classifieur** $s(\mathbf{x}) \rightarrow \hat{y} \in \{0, 1\}$.



Bases de données

- Repésentation d'un assuré par un avocat (accident corporel)

```

1 > avocat <- read.table("http://freakonometrics.free.fr/AutoBI.csv",
2   header=TRUE,sep=",")
3 > avocat$CLMSEX <- factor(avocat$CLMSEX, labels=c("M","F"))
4 > avocat$MARITAL <- factor(avocat$MARITAL, labels=c("M","C","V","D"))
5 > avocat <- avocat[,-c(1,5)]

```

- Survenance d'accident sur des polices d'assurance de camping car

```

5 > camping <- read.table("http://freakonometrics.free.fr/baseN.txt",
6   header=TRUE,sep=";")

```

- Décès par infarctus suite à une admission aux urgences

```

6 > myocarde <- read.table("http://freakonometrics.free.fr/myocarde.csv",
7   " ,head=TRUE,sep=";")

```

(cf. code en ligne sur le blog, détails donnés dans les slides)

La loi binomiale $\mathcal{B}(\pi)$

Soient Y_i des variables aléatoires prenant les valeurs $y_i \in \{0, 1\}$, avec probabilité $1 - \pi_i$ et π_i (respectivement), i.e.

$$\mathbb{P}(Y_i = y_i) = \pi_i^{y_i} [1 - \pi_i]^{1-y_i}, \text{ pour } y_i \in \{0, 1\},$$

avec $\pi_i \in [0, 1]$. Aussi, $\mathbb{P}(Y_i = 1) = \pi_i$ et $\mathbb{P}(Y_i = 0) = 1 - \pi_i$.

Alors $\mathbb{E}(Y_i) = \pi_i$ et $\text{Var}(Y_i) = \pi_i[1 - \pi_i] < \mathbb{E}(Y_i)$.

En statistique, on suppose que $\pi_i = \pi, \forall i \in \{1, 2, \dots, n\}$,

et on cherche à **estimer** π

Maximum de vraisemblance et méthode des moments

La fonction de vraisemblance, pour un échantillon y_1, \dots, y_n s'écrit

$$\mathcal{L}(\pi; \mathbf{y}) = \prod_{i=1}^n \mathbb{P}(Y_i = y_i) = \prod_{i=1}^n \pi^{y_i} [1 - \pi]^{1-y_i}$$

et la log-vraisemblance est

$$\log \mathcal{L}(\pi; \mathbf{y}) = \sum_{i=1}^n y_i \log[\pi] + (1 - y_i) \log[1 - \pi]$$

La condition du premier ordre est

$$\frac{\partial \log \mathcal{L}(\pi; \mathbf{y})}{\partial \pi} = \sum_{i=1}^n \frac{y_i}{\pi} - \frac{1 - y_i}{1 - \pi} = 0$$

Intervalle de confiance pour π

On a un intervalle de confiance asymptotique car on sait que l'estimateur de vraisemblance est asymptotiquement efficace,

$$\sqrt{n}(\hat{\pi} - \pi) \xrightarrow{\mathcal{L}} \mathcal{N}(0, I^{-1}(\pi))$$

où $I(\pi)$ est l'information de Fisher, i.e.

$$I(\pi) = \frac{1}{\pi[1-\pi]}$$

d'où un intervalle de confiance approché (à 95%) pour π de la forme

$$\left[\hat{\pi} \pm \frac{1.96}{\sqrt{n}} \sqrt{\hat{\pi}[1-\hat{\pi}]} \right].$$

Intervalle de confiance pour π

On peut aussi construire un intervalle de confiance, par le théorème central limite, car $\hat{\pi} = \bar{X}$. On sait que

$$\sqrt{n} \frac{\bar{X} - \mathbb{E}(X)}{\sqrt{\text{Var}(X)}} \xrightarrow{\mathcal{L}} \mathcal{N}(0, 1),$$

avec ici $\bar{X} = \hat{\pi}$, $\mathbb{E}(X) = \pi$ et $\text{Var}(X) = \pi(1 - \pi)$, i.e. un intervalle de confiance est obtenu par l'approximation

$$\sqrt{n} \frac{\hat{\pi} - \pi}{\sqrt{\hat{\pi}[1 - \hat{\pi}]}} \xrightarrow{\mathcal{L}} \mathcal{N}(0, 1),$$

d'où un intervalle de confiance (à 95%) pour π de la forme

$$\left[\hat{\pi} \pm \frac{1.96}{\sqrt{n}} \sqrt{\hat{\pi}[1 - \hat{\pi}]} \right].$$

Mise en oeuvre pratique

Considérons l'échantillon suivant, $\{y_1, \dots, y_n\}$

```

1 > set.seed(1)
2 > n=20
3 > (Y=sample(0:1, size=n, replace=TRUE))
4 [1] 0 0 1 1 0 1 1 1 1 0 0 0 1 0 1 0 1 1 0 1

```

Comme $Y_i \sim \mathcal{B}(\pi)$, avec $\pi = \mathbb{E}(Y)$. On note $\hat{\pi} = \bar{y}$, soit ici

```

1 > mean(x)
2 [1] 0.55

```

On peut faire un test de $H_0 : \pi = \pi_\star$ contre $H_1 : \pi \neq \pi_\star$ (par exemple 50%)

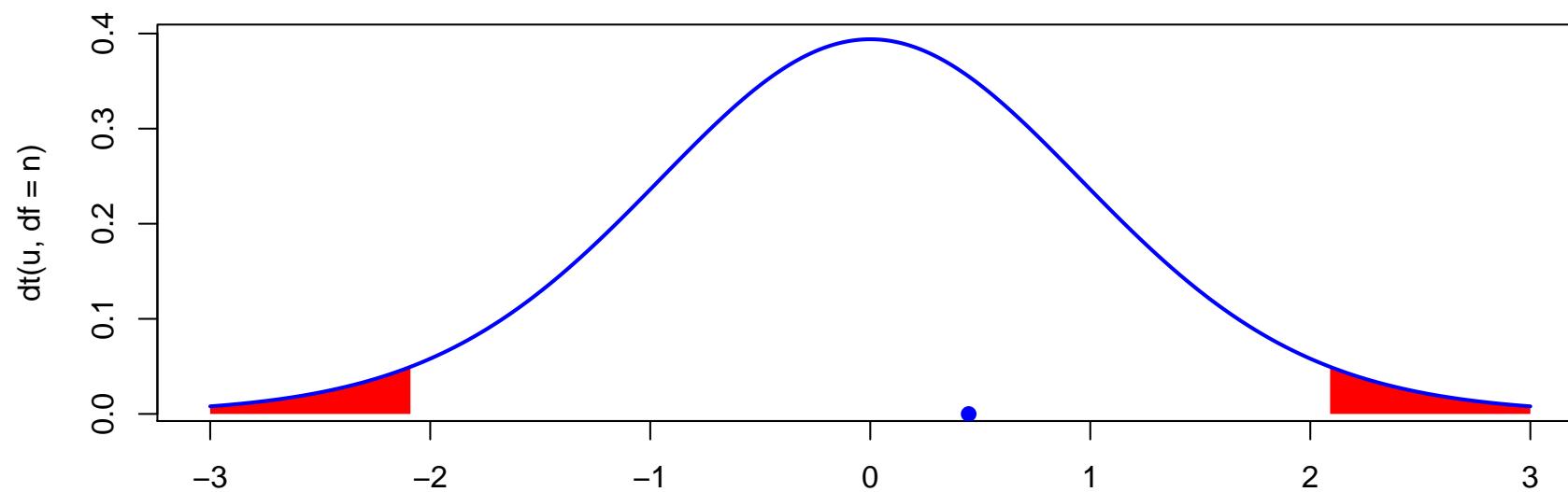
On peut utiliser le test de Student,

$$T = \sqrt{n} \frac{\hat{\pi} - \pi_\star}{\sqrt{\pi_\star(1 - \pi_\star)}}$$

qui suit, sous H_0 une loi de Student à n degrés de liberté.

Mise en oeuvre pratique

```
1 > (T=sqrt(n)*(pn-p0)/(sqrt(p0*(1-p0))))  
2 [1] 0.4472136  
3 > abs(T)<qt(1-alpha/2,df=n)  
4 [1] TRUE
```

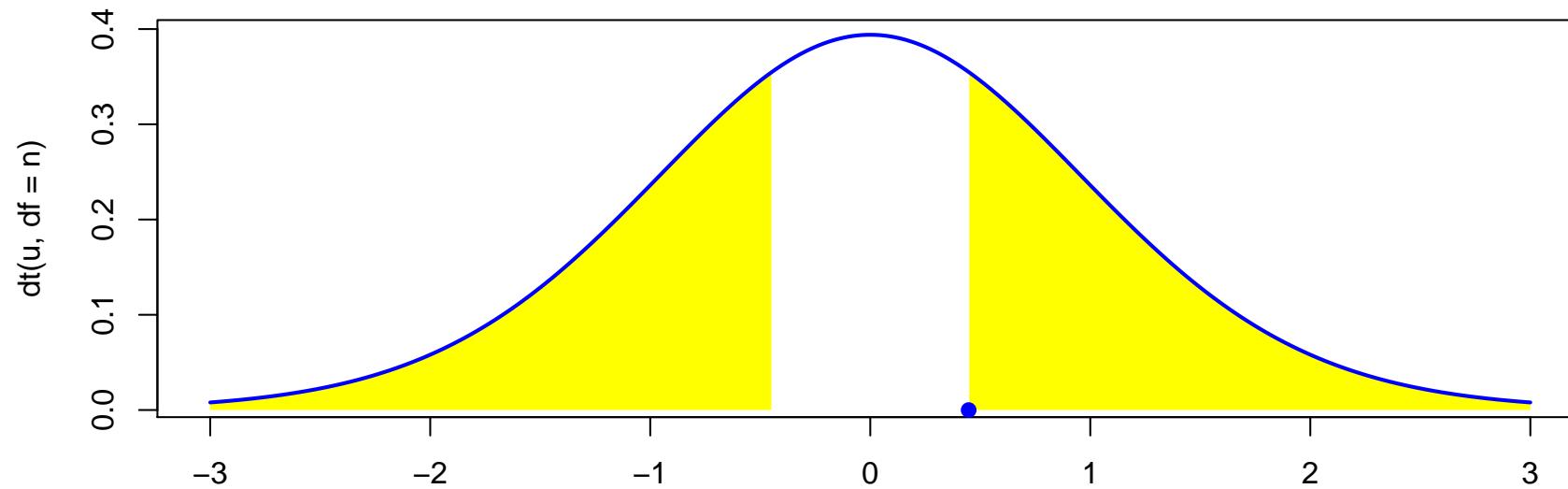


Mise en oeuvre pratique

On est ici dans la région d'acceptation du test.

On peut aussi calculer la p -value, $\mathbb{P}(|T| > |t_{obs}|)$,

```
1 > 2*(1-pt(abs(T),df=n))  
2 [1] 0.6595265
```



Mise en oeuvre pratique

Test de Wald l'idée est d'étudier la différence entre $\hat{\pi}$ et π_\star . Sous H_0 ,

$$T = n \frac{(\hat{\pi} - \pi_\star)^2}{I^{-1}(\pi_\star)} \xrightarrow{\mathcal{L}} \chi^2(1)$$

Test du rapport de vraisemblance l'idée est d'étudier la différence entre $\log \mathcal{L}(\hat{\theta})$ et $\log \mathcal{L}(\theta_\star)$. Sous H_0 ,

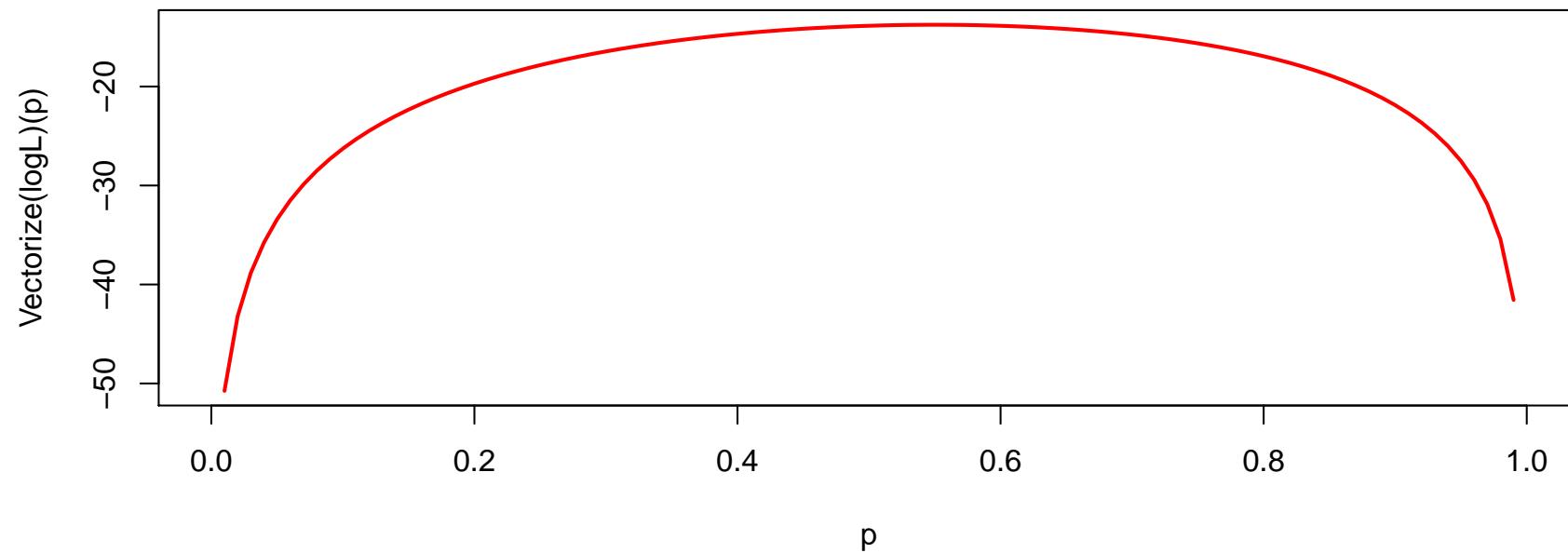
$$T = 2 \log \left(\frac{\log \mathcal{L}(\theta_\star)}{\log \mathcal{L}(\hat{\theta})} \right) \xrightarrow{\mathcal{L}} \chi^2(1)$$

Test du score l'idée est d'étudier la différence entre $\frac{\partial \log \mathcal{L}(\pi_\star)}{\partial \pi}$ et 0. Sous H_0 ,

$$T = \left(\frac{1}{n} \sum_{i=1}^n \frac{\partial \log f_{\pi_\star}(x_i)}{\partial \pi} \right)^2 \xrightarrow{\mathcal{L}} \chi^2(1)$$

Mise en oeuvre pratique

```
1 > p=seq(0,1,by=.01)
2 > logL=function(p){sum(log(dbinom(X,size=1,prob=p)))}
3 > plot(p,Vectorize(logL)(p),type="l",col="red",lwd=2)
```



Mise en oeuvre pratique

On peut trouver numériquement le maximum de $\log \mathcal{L}$,

```
1 > neglogL=function(p){-sum(log(dbinom(X,size=1,prob=p)))}
2 > pml=optim(fn=neglogL,par=p0,method="BFGS")
3 > pml
4 $par
5 [1] 0.5499996
6
7$value
8 [1] 13.76278
```

i.e. on retrouve ici $\hat{\pi} = \bar{y}$.

Mise en oeuvre pratique

Maintenant, on peut tester $H_0 : \pi = \pi_\star = 50\%$ versus $H_1 : \pi \neq 50\%$. Pour le test de Wald, on peut commencer par calculer $nI(\theta_\star)$ ou ,

```

1 > nx=sum(X==1)
2 > f = expression(nx*log(p)+(n-nx)*log(1-p))
3 > Df = D(f, "p")
4 > Df2 = D(Df, "p")
5 > p=p0=0.5
6 > (IF=-eval(Df2))
7 [1] 80

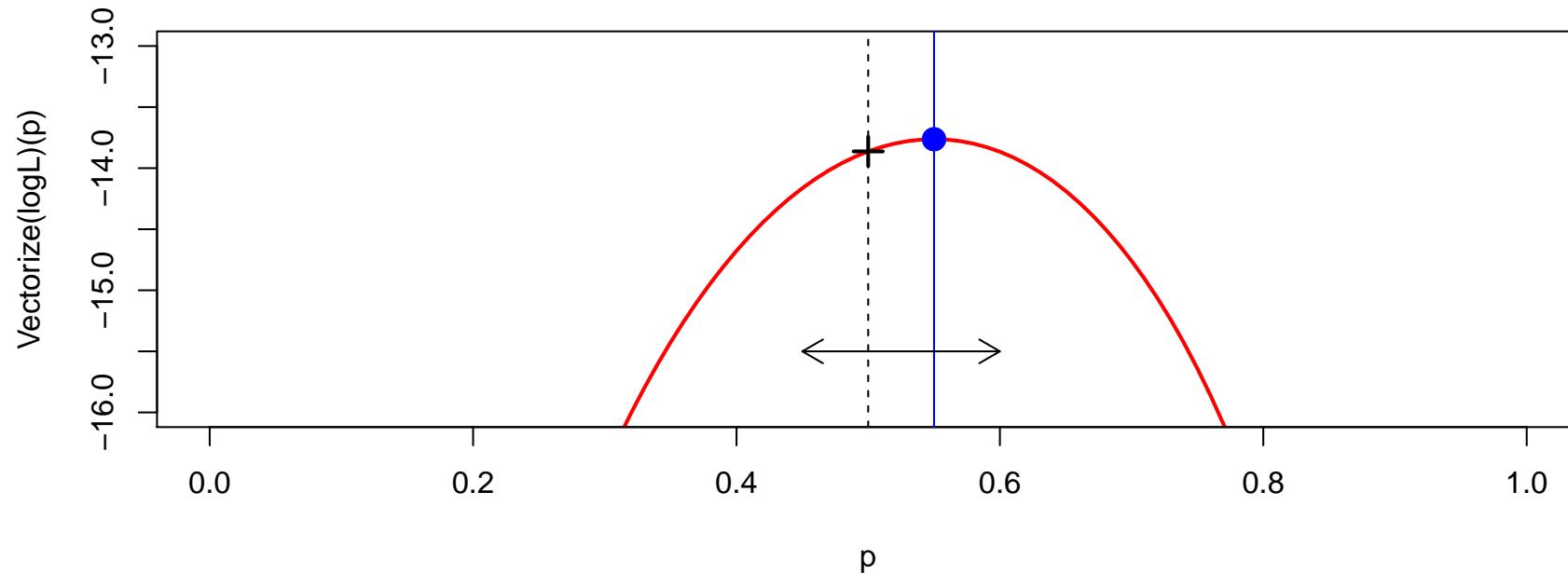
```

Mise en oeuvre pratique

On peut d'ailleurs comparer cette valeur avec la valeur théorique, car

$$I(\pi)^{-1} = \pi(1 - \pi)$$

```
1 > 1 / (p0 * (1-p0) / n)
2 [1] 80
```



Mise en oeuvre pratique

La statistique du test de Wald est ici

```
1 > pml=optim(fn=neglogL,par=p0,method="BFGS")$par
2 > (T=(pml-p0)^2*IF)
3 [1] 0.199997
```

que l'on va chercher à comparer avec le quantile de la loi du χ^2 ,

```
1 > T<qchisq(1-alpha,df=1)
2 [1] TRUE
```

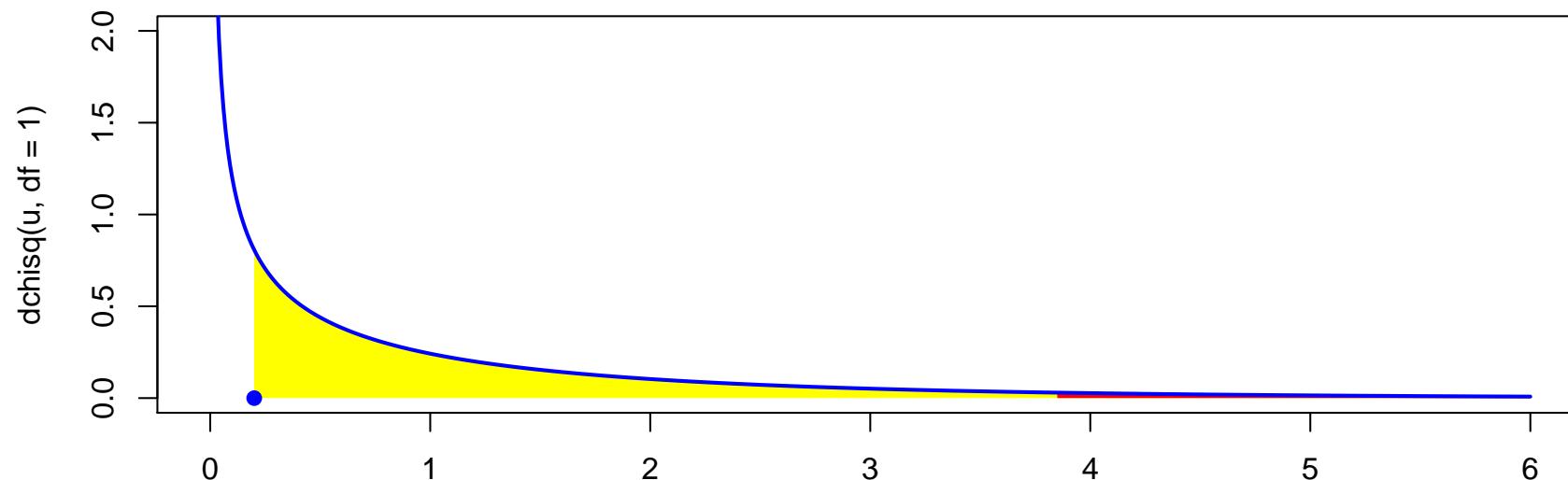
i.e. on est dans la région d'acceptation du test.

Mise en oeuvre pratique

De manière duale, on peut calculer la p -value du test,

```
1 > 1-pchisq(T, df=1)
2 [1] 0.6547233
```

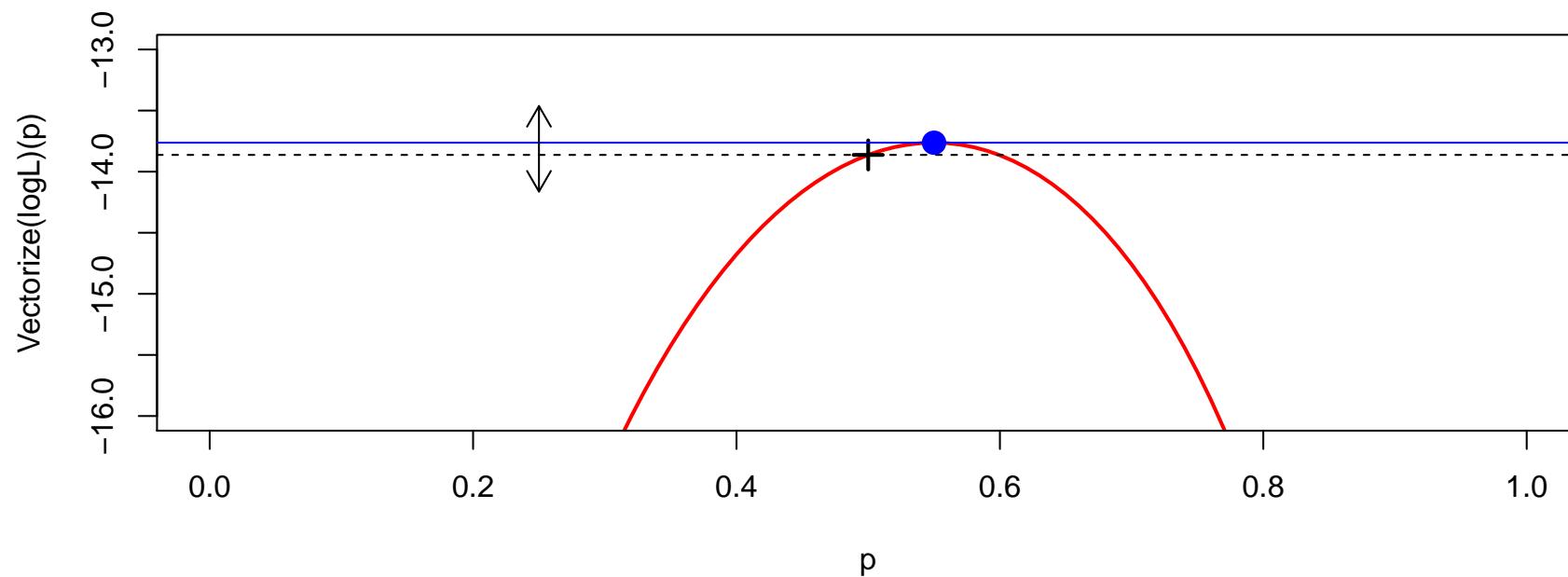
Donc on va accepter H_0 .



Mise en oeuvre pratique

Pour le test du rapport de vraisemblance, T est ici

```
1 > (T=2*(logL(pml)-logL(p0)))
2 [1] 0.2003347
```



Mise en oeuvre pratique

Là encore, on est dans la région d'acceptation,

```
1 > T<qchisq(1-alpha, df=1)
2 [1] TRUE
```

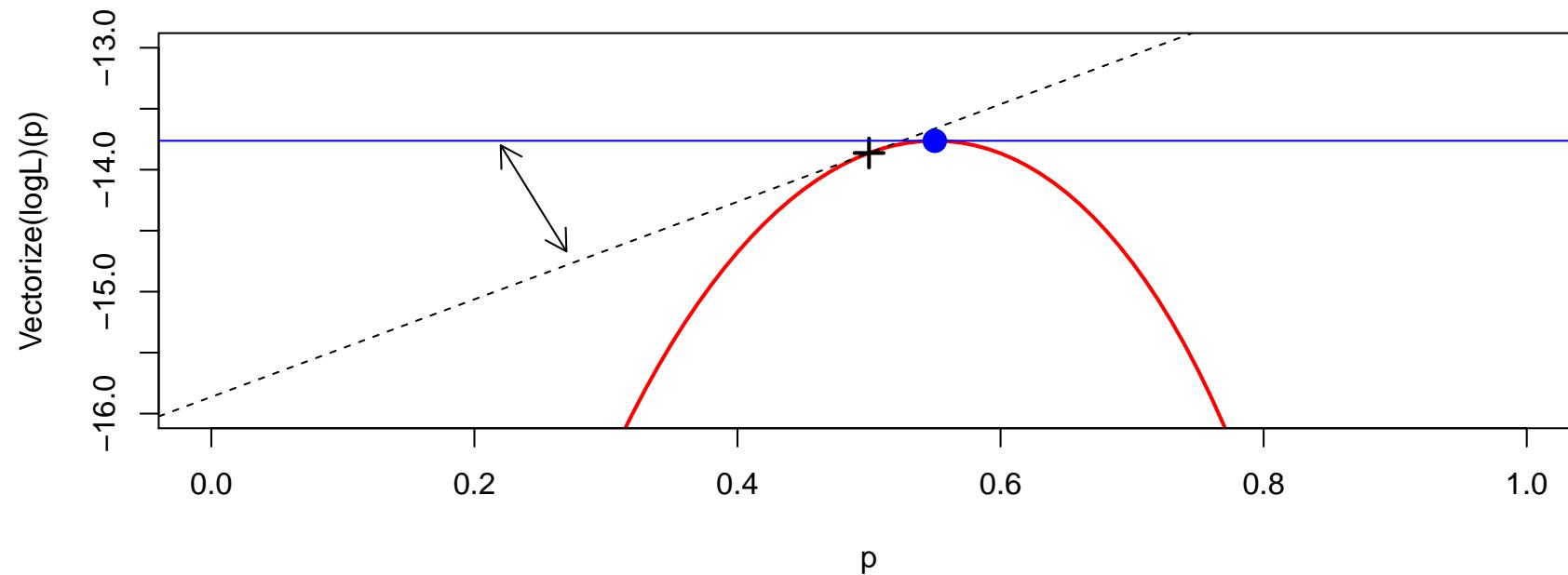
Enfin, on peut faire le test du score. Le score se calcule facilement,

```
1 > nx=sum(X==1)
2 > f = expression(nx*log(p)+(n-nx)*log(1-p))
3 > Df = D(f, "p")
4 > p=p0
5 > score=eval(Df)
```

La statistique de test est alors

```
1 > (T=score^2/IF)
2 [1] 0.2
```

Mise en oeuvre pratique



Là encore, on est dans la région d'acceptation,

```
1 > T< qchisq(1-alpha, df=1)
2 [1] TRUE
```

Prise en compte de l'hétérogénéité

Supposons que π_i soit fonction de variables explicatives \mathbf{X}_i , e.g.

$$\pi_i = \mathbb{E}(Y_i | \mathbf{X}_i).$$

Le modèle classique est le [modèle linéaire](#), $\mathbb{E}(Y_i | \mathbf{X}_i) = \mathbf{X}'_i \boldsymbol{\beta}$

On peut faire une estimation par moindres carrés pour estimer $\boldsymbol{\beta}$,

$$\hat{\boldsymbol{\beta}} = \operatorname{argmin} \left\{ \sum_{i=1}^n [Y_i - \mathbf{X}'_i \boldsymbol{\beta}]^2, \boldsymbol{\beta} \in \mathbb{R}^k \right\}$$

Problème: $\pi_i \in [0, 1]$ et $\mathbf{X}'_i \boldsymbol{\beta} \in \mathbb{R}$.

Problème: modèle non homoscédastique, en effet

$$\varepsilon_i = \begin{cases} 1 - \mathbf{X}'_i \boldsymbol{\beta} & \text{avec probabilité } \pi_i \\ -\mathbf{X}'_i \boldsymbol{\beta} & \text{avec probabilité } 1 - \pi_i \end{cases}$$

i.e. $\operatorname{Var}(\varepsilon) = \mathbf{X}'_i \boldsymbol{\beta} \cdot [1 - x \mathbf{X}'_i \boldsymbol{\beta}]$.

Utilisation de la cote (odds)

$\pi_i \in [0, 1]$ alors que $\mathbf{X}'_i \boldsymbol{\beta} \in \mathbb{R}$, donc on ne peut pas supposer

$$\pi_i = \mathbb{E}(Y_i | \mathbf{X}_i) = \mathbf{X}'_i \boldsymbol{\beta}$$

→ utilisation de la **côte**

$$\text{odds}_i = \frac{\pi_i}{1 - \pi_i} \in [0, \infty].$$

ou, en passant au logarithme

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

On appelle **logit** cette transformation,

$$\text{logit}(\pi_i) = \log\left(\frac{\pi_i}{1 - \pi_i}\right)$$

Utilisation de la cote (odds)

On va supposer que

$$\text{logit}(\pi_i) = \mathbf{X}'_i \boldsymbol{\beta}$$

ou

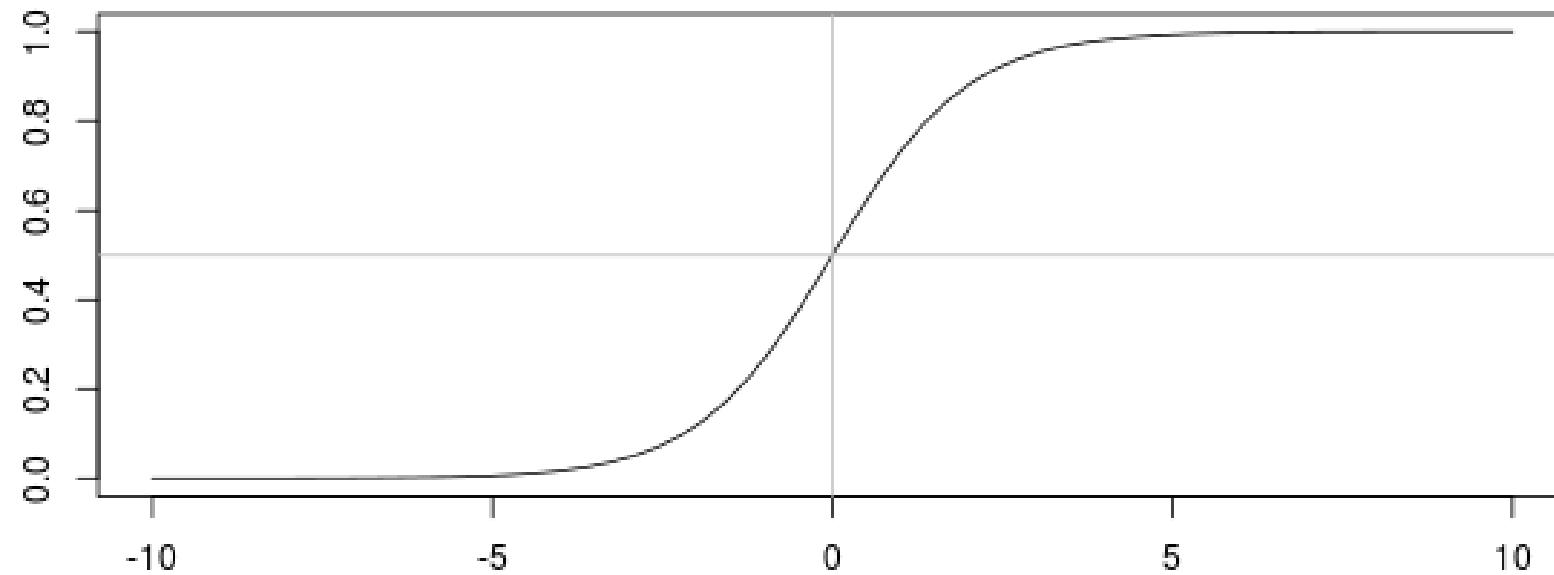
$$\pi_i = \text{logit}^{-1}(\mathbf{X}'_i \boldsymbol{\beta}) = \frac{\exp[\mathbf{X}'_i \boldsymbol{\beta}]}{1 + \exp[\mathbf{X}'_i \boldsymbol{\beta}]}.$$

Remarque: si π_i est faible ($\pi_i \sim 0$) alors $\text{logit}^{-1} \sim \exp$, i.e.

$$\pi_i \sim \exp(\mathbf{X}'_i \boldsymbol{\beta})$$

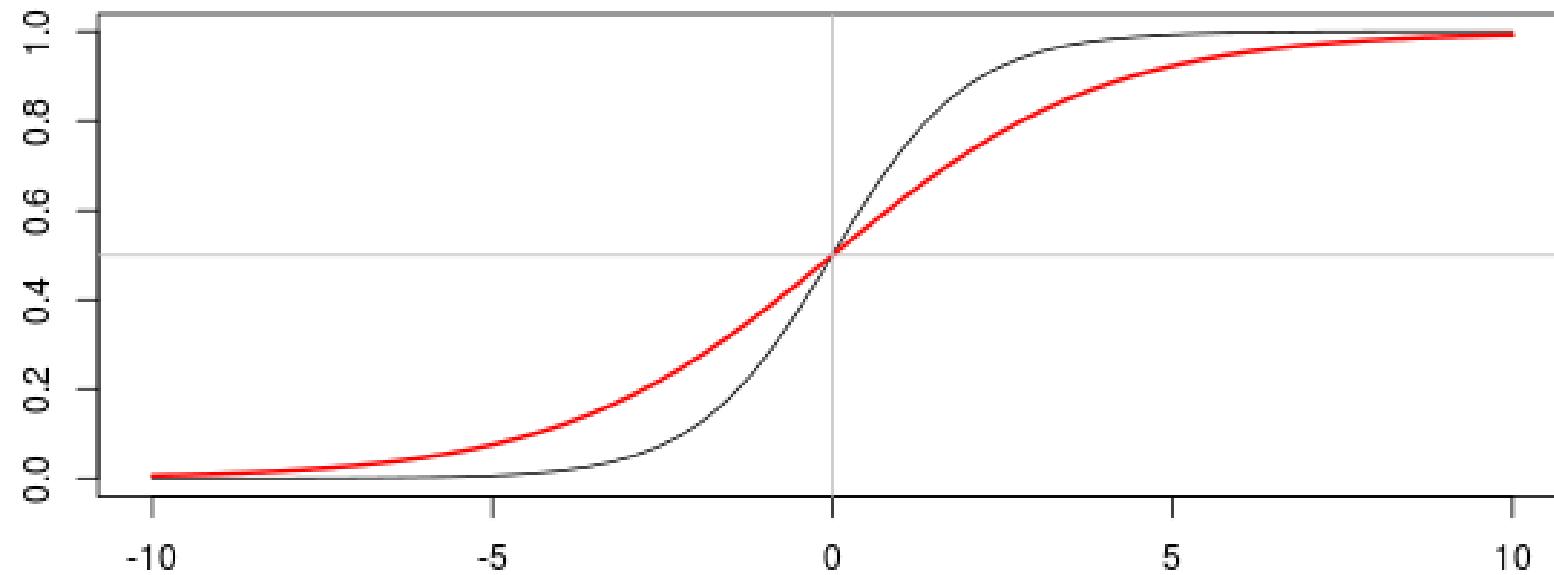
La fonction logistique

$$F(x) = \frac{\exp[x]}{1 + \exp[x]}, \forall x \in \mathbb{R}.$$



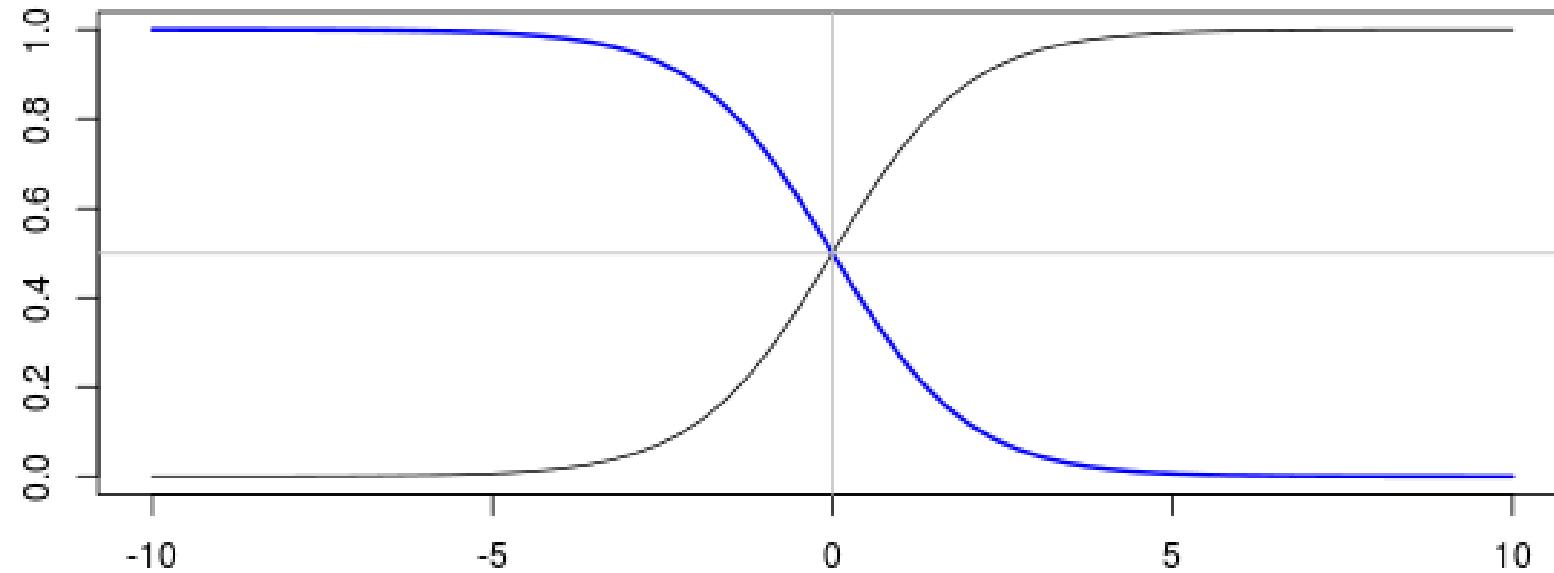
La fonction logistique

$$F(x) = \frac{\exp[2x]}{1 + \exp[2x]}, \forall x \in \mathbb{R}.$$



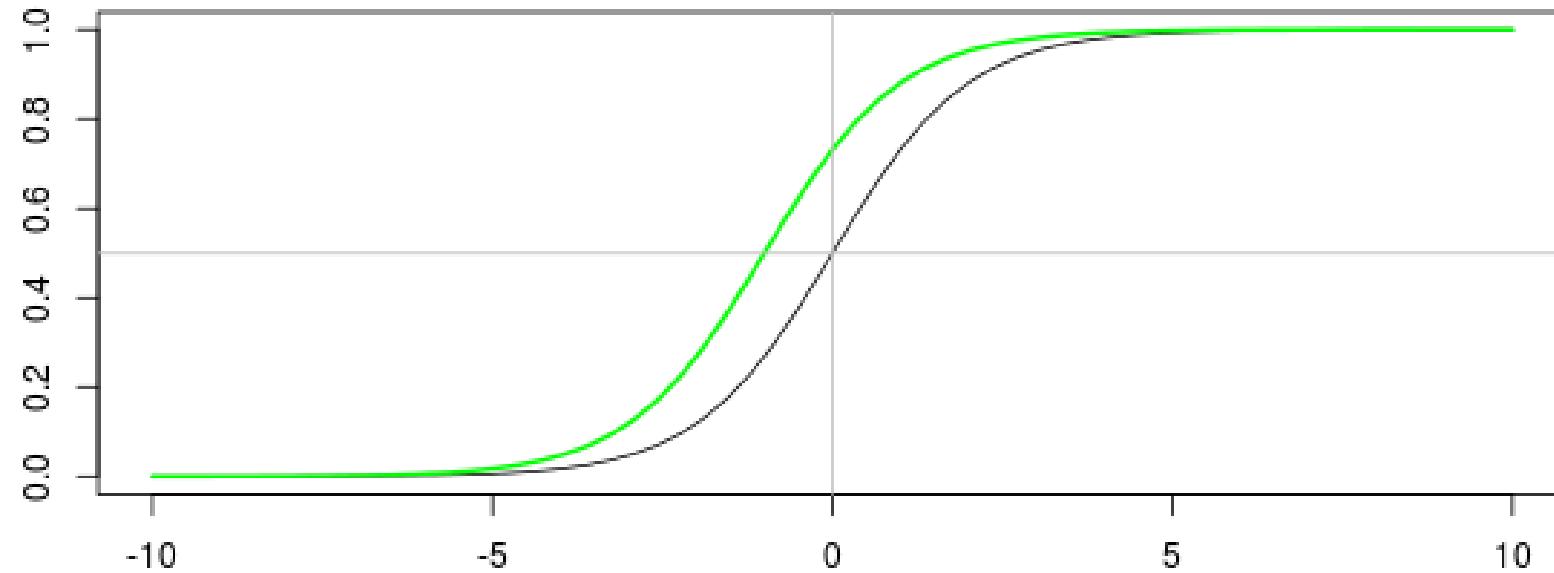
La fonction logistique

$$F(x) = \frac{\exp[-x]}{1 + \exp[-x]}, \forall x \in \mathbb{R}.$$



La fonction logistique

$$F(x) = \frac{\exp[1+x]}{1 + \exp[1+x]}, \forall x \in \mathbb{R}.$$



La régression logistique

La log-vraisemblance est ici

$$\log \mathcal{L}(\boldsymbol{\beta}) = \sum_{i=1}^n y_i \log(\pi_i) + (1-y_i) \log(1-\pi_i) = \sum_{i=1}^n y_i \log(\pi_i(\boldsymbol{\beta})) + (1-y_i) \log(1-\pi_i(\boldsymbol{\beta}))$$

La résolution se fait en écrivant les conditions du premier ordre. Pour cela, on écrit le gradient $\nabla \log \mathcal{L}(\boldsymbol{\beta}) = [\partial \log \mathcal{L}(\boldsymbol{\beta}) / \partial \beta_k]$, où

$$\frac{\partial \log \mathcal{L}(\boldsymbol{\beta})}{\partial \beta_k} = \sum_{i=1}^n \frac{y_i}{\pi_i(\boldsymbol{\beta})} \frac{\partial \pi_i(\boldsymbol{\beta})}{\partial \beta_k} - \frac{1-y_i}{\pi_i(\boldsymbol{\beta})} \frac{\partial \pi_i(\boldsymbol{\beta})}{\partial \beta_k}$$

or compte tenu de la forme de $\pi_i(\boldsymbol{\beta})$,

$$\frac{\partial \pi_i(\boldsymbol{\beta})}{\partial \beta_k} = \pi_i(\boldsymbol{\beta})[1 - \pi_i(\boldsymbol{\beta})]X_{k,i}$$

on obtient

$$\frac{\partial \log \mathcal{L}(\boldsymbol{\beta})}{\partial \beta_k} = \sum_{i=1}^n X_{k,i}[y_i - \pi_i(\boldsymbol{\beta})]$$

Newton Raphson - ‘Descente’ de gradient

Remarque Pour rechercher le **zéro** d'une fonction (réelle), i.e. $f(x) = 0$: on se donne x_0 , et on pose

$$x_k = x_{k-1} - \frac{f(x_{k-1})}{f'(x_{k-1})}$$

Remarque Pour rechercher le **maximum** d'un fonction dérivable, on recherche le zéro de la dérivée, i.e. $h'(x) = 0$: on se donne x_0 , et on pose

$$x_k = x_{k-1} - \frac{f'(x_{k-1})}{f''(x_{k-1})}$$

Il faut trouver numériquement le maximum de $\log \mathcal{L}(\beta)$.

L'algorithme standard est

1. partir d'une valeur initiale β_0
2. poser $\beta_k = \beta_{k-1} - H(\beta_{k-1})^{-1} \nabla \log \mathcal{L}(\beta_{k-1})$

où $\nabla \log \mathcal{L}(\boldsymbol{\beta})$ est le gradient, et $H(\boldsymbol{\beta})$ la matrice Hessienne (on parle parfois de Score de Fisher).

Ici, le terme générique de la matrice Hessienne est

$$\frac{\partial^2 \log \mathcal{L}(\boldsymbol{\beta})}{\partial \beta_k \partial \beta_\ell} = \sum_{i=1}^n X_{k,i} X_{\ell,i} [y_i - \pi_i(\boldsymbol{\beta})]$$

Posons

$$\boldsymbol{\Omega} = [\omega_{i,j}] = \text{diag}(\hat{\pi}_i(1 - \hat{\pi}_i))$$

alors le gradient est

$$\nabla \log \mathcal{L}(\boldsymbol{\beta}) = \frac{\partial \log \mathcal{L}(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}} = \mathbf{X}^\top (\mathbf{y} - \boldsymbol{\pi})$$

alors que la matrice Hessienne estimation

$$H(\boldsymbol{\beta}) = \frac{\partial^2 \log \mathcal{L}(\boldsymbol{\beta})}{\partial \boldsymbol{\beta} \partial \boldsymbol{\beta}^\top} = -\mathbf{X}^\top \boldsymbol{\Omega} \mathbf{X}$$

Aussi, l'algorithme de Newton Raphson s'écrit

$$\boldsymbol{\beta}_k = (\mathbf{X}^\top \boldsymbol{\Omega} \mathbf{X})^{-1} \mathbf{X}^\top \boldsymbol{\Omega} \mathbf{Z} \text{ où } \mathbf{Z} = \mathbf{X} \boldsymbol{\beta}_{k-1} + \mathbf{X}^\top \boldsymbol{\Omega}^{-1} \mathbf{y} - \boldsymbol{\pi}),$$

On reconnaît une régression pondérée itérée,

$$\boldsymbol{\beta}_k = \operatorname{argmin} \{(\mathbf{Z} - \mathbf{X} \boldsymbol{\beta})^\top \boldsymbol{\Omega} (\mathbf{Z} - \mathbf{X} \boldsymbol{\beta})\}$$

D'un point de vue computationnel, l'estimation se fait de la manière suivante:

```
1 X=cbind(rep(1,nrow(avocat)),avocat$LOSS)
2 Y=(avocat$ATTORNEY=2)*1
```

on part d'une valeur initiale (e.g. une estimation classique de modèle linéaire)

```
1 beta=lm(Y~0+X)$coefficients
```

On fait ensuite une boucle,

```
1 for(s in 1:20){
2 pi=exp(X%*%beta)/(1+exp(X%*%beta))
```

```

3 gradient=t(X) *%*(Y-pi)
4 omega=matrix(0,nrow(X),nrow(X));diag(omega)=(pi*(1-pi))
5 hesienne=-t(X) *%* omega *%* X
6 beta=beta-solve(hesienne) *%* gradient}

```

Ici, on obtient les valeurs suivantes

	[,1]	[,2]
1		
2	[1 ,] 0.50034948	0.001821554
3	[2 ,] -0.05689403	0.015434935
4	[3 ,] -0.19941087	0.055570172
5	[4 ,] -0.45629828	0.140490153
6	[5 ,] -0.73627462	0.253154397
7	[6 ,] -0.85830164	0.309712845
8	[7 ,] -0.87192186	0.316615292
9	[8 ,] -0.87207800	0.316697581

i.e. on converge (et vite).

Comportement asymptotique (et intervalle de confiance)

On peut montrer que $\widehat{\beta} \xrightarrow{\mathbb{P}} \beta$ et

$$\sqrt{n}(\widehat{\beta} - \beta) \xrightarrow{\mathcal{L}} \mathcal{N}(\mathbf{0}, I(\beta)^{-1}).$$

Numériquement, la encore, on peut approcher $I(\beta)^{-1}$ qui est la variance (asymptotique) de notre estimateur. Or $I(\beta) = H(\beta)$, donc les écart-types de $\widehat{\beta}$ sont

```

1 > pi=exp(X*%beta)/(1+exp(X*%beta))
2 + omega=matrix(0,nrow(X),nrow(X));diag(omega)=(pi*(1-pi))
3 + hesienne=-t(X)*%omega%*%X
4 + sqrt(-diag(solve(hesienne)))
5 [1] 0.09128379 0.02902411

```

On retrouve toutes ces valeurs en utilisant

```

1 > logistique <- glm((ATTORNEY==1)~LOSS, data=avocat, family=binomial(
2   link="logit"))
3
4 Coefficients:
5             Estimate Std. Error z value Pr(>|z|)
6 (Intercept) -0.87208    0.09128 -9.553  <2e-16 ***
7 LOSS         0.31670    0.02902 10.912  <2e-16 ***
8 ---
9 Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
10 (Dispersion parameter for binomial family taken to be 1)
11
12 Null deviance: 1857.6 on 1340 degrees of freedom
13 Residual deviance: 1609.4 on 1338 degrees of freedom
14 AIC: 1613.4
15
16 Number of Fisher Scoring iterations: 8

```

Test d'hypothèse et significativité

Considérons un test $H_0 : \beta_k = 0$ contre $H_1 : \beta_k \neq 0$ (classique dans le modèle linéaire).

Compte tenu du comportement asymptotiquement Gaussien des estimateurs, on a automatiquement un test de significativité des coefficients: sous H_0 ,

$$\frac{\widehat{\beta}_k}{\sqrt{\widehat{\text{Var}}(\widehat{\beta}_k)}} \sim \mathcal{N}(0, 1).$$

Remarque: ce test correspond à un test de Wald (d'où la notation \mathbf{z} sous \mathbf{R})

Régression linéaire sur une variable catégorielle

Supposons que l'on régresse sur X , une variable binaire, i.e $X \in \{0, 1\}$, avec

$$\mathbb{P}(Y = 1|X) = \beta_0 + \beta_1 X$$

qui prend les valeurs β_0 si $x = 0$ (modalité de référence) et $\beta_0 + \beta_1$ si $x = 1$.

La log-vraisemblance devient

$$\begin{aligned} \log \mathcal{L}(\beta_0, \beta_1) &= \sum_{i, x_i=0} y_i \log[\beta_0] + \sum_{i, x_i=1} y_i \log[\beta_0 + \beta_1] \\ &\quad + \sum_{i, x_i=0} (1 - y_i) \log[1 - \beta_0] + \sum_{i, x_i=1} (1 - y_i) \log[1 - \beta_0 + \beta_1] \end{aligned}$$

et les conditions du premier ordre donnent

$$\widehat{\beta}_0 = \frac{\sum_{x_i=0} y_i}{\sum_{x_i=0} 1} \text{ et } \widehat{\beta}_1 = \frac{\sum_{x_i=1} y_i}{\sum_{x_i=0} 1} - \frac{\sum_{x_i=0} y_i}{\sum_{x_i=0} 1}$$

Régression linéaire sur une variable catégorielle

Sur nos données, on peut considérer le sexe comme variable explicative.

```

1 > xtabs(~ ATTORNEY+CLMSEX , data=avocat)
2
3 CLMSEX
4
5 ATTORNEY   M     F
6
7       1 325 352
8       2 261 390

```

Si on fait la régression linéaire

```

1 > reglm = lm((ATTORNEY==1) ~ CLMSEX , data=avocat)
2 > summary(reglm)
3
4 Coefficients:
5
6             Estimate Std. Error t value Pr(>|t|)    
7 (Intercept) 0.55461    0.02060 26.922 < 2e-16 ***
8 CLMSEXF     -0.08021    0.02756 -2.911  0.00367 ** 

```

Régression linéaire sur une variable catégorielle

Ces paramètres correspondent aux fréquences empiriques,

```
1 > attach(avocat)
2 > sum((ATTORNEY==1)&(CLMSEX=="M") ,na.rm=TRUE) / sum(CLMSEX=="M" ,na.rm=
   TRUE)
3 [1] 0.5546075
4 > sum((ATTORNEY==1)&(CLMSEX=="F") ,na.rm=TRUE) / sum(CLMSEX=="F" ,na.rm=
   TRUE)
5 [1] 0.4743935
```

Régression logistique sur une variable catégorielle

Avec un modèle logistique, on obtient les mêmes prédictions,

```

1 > reglogit = glm((ATTORNEY==2) ~ CLMSEX, data=avocat, family=binomial)
2 > summary(reglogit)
3
4 Coefficients:
5             Estimate Std. Error z value Pr(>|z|)
6 (Intercept) -0.21930   0.08312 -2.639  0.00833 ** 
7 CLMSEXF      0.32182   0.11097  2.900  0.00373 ** 
8 ---
9 Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
10
11 > exp(reglogit$coefficients[1])/(1+exp(reglogit$coefficients[1]))
12 [1] 0.4453925
13 > exp(sum(reglogit$coefficients[1:2]))/(1+exp(sum(reglogit$coefficients[1:2])))
14 [1] 0.5256065

```

Modèle alternatif, le modèle probit

Soit H une fonction de répartition d'une variable réelle, $H : \mathbb{R} \mapsto [0, 1]$. On suppose ici

$$\pi_i = H(\mathbf{X}'_i \boldsymbol{\beta}) \text{ ou } \mathbf{X}'_i \boldsymbol{\beta} = H^{-1}(\pi_i)$$

- si H est la fonction de répartition de la loi logistique, on retrouve le modèle logit,
- si H est la fonction de répartition de la loi $\mathcal{N}(0, 1)$, on obtient le modèle probit,

Remarque: probit signifie probability unit.

On suppose qu'il existe une variable latente Y_i^* non observée, telle que

$$\pi_i = \mathbb{P}(Y_i = 1) = \mathbb{P}(Y_i^* > 0)$$

On suppose que $Y_i^* = \mathbf{X}'_i \boldsymbol{\beta} + \epsilon_i$, où ϵ_i est une erreur de loi H . Alors

$$\pi_i = \mathbb{P}(Y_i = 1) = \mathbb{P}(Y_i^* > 0) = \mathbb{P}(\epsilon_i > -\mathbf{X}'_i \boldsymbol{\beta}) = 1 - H(-\mathbf{X}'_i \boldsymbol{\beta})$$

alors que

$$1 - \pi_i = \mathbb{P}(Y_i = 0) = \mathbb{P}(Y_i^* \leq 0) = \mathbb{P}(\epsilon_i \leq -\mathbf{X}'_i \boldsymbol{\beta}) = H(-\mathbf{X}'_i \boldsymbol{\beta})$$

Si $H = \Phi$, alors H est symétrique par rapport à 0, et

$$\pi_i = \mathbb{P}(Y_i = 1) = 1 - \Phi(-\mathbf{X}'_i \boldsymbol{\beta}) = \Phi(\mathbf{X}'_i \boldsymbol{\beta}),$$

ou encore $\eta_i = \mathbf{X}'_i \boldsymbol{\beta} = \Phi^{-1}(\pi_i)$.

Si $H(x) = \frac{e^x}{1 + e^x}$, on a une loi logistique, qui est aussi symétrique par rapport à 0, et

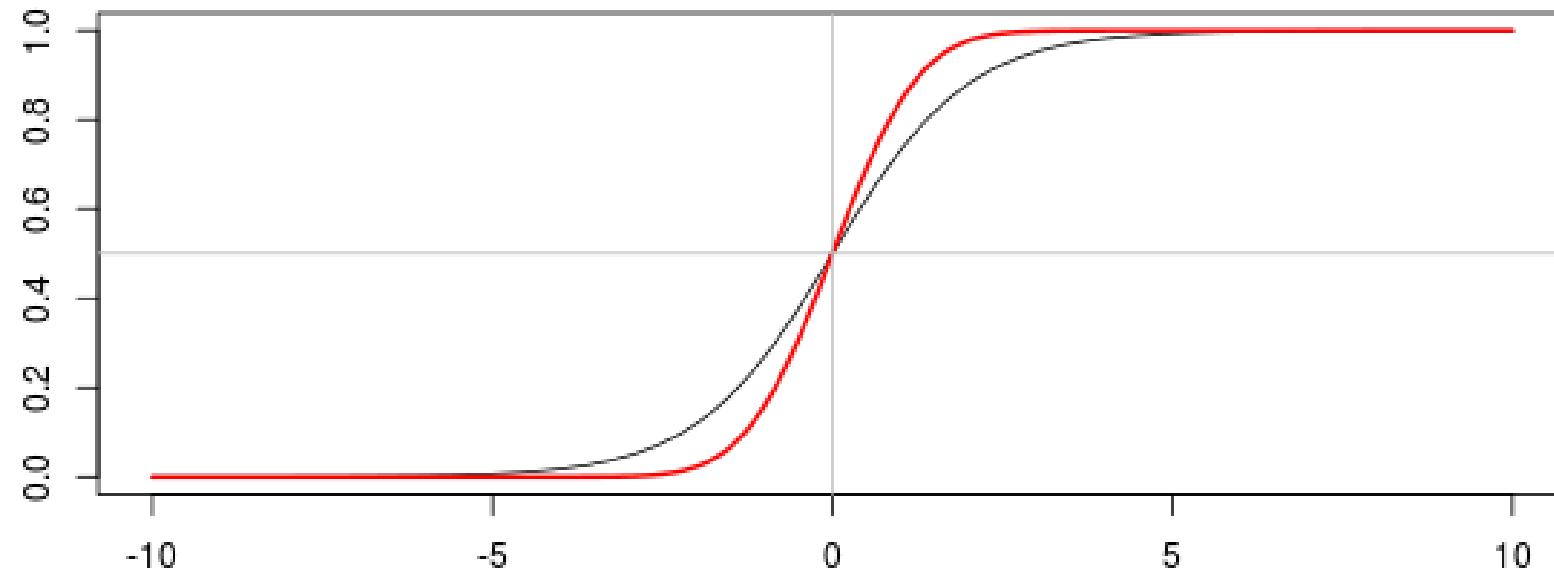
$$\pi_i = \mathbb{P}(Y_i = 1) = 1 - H(-\mathbf{X}'_i \boldsymbol{\beta}) = H(\mathbf{X}'_i \boldsymbol{\beta}),$$

ou encore $\eta_i = \mathbf{X}'_i \boldsymbol{\beta} = H^{-1}(\pi_i)$.

Remarque: pourquoi prendre un seuil nul ? sinon le modèle ne serait pas identifiable. De même si on ne fixe pas la variance des résidus. On prend donc (arbitrairement) $s = 0$ et $\sigma^2 = 1$.

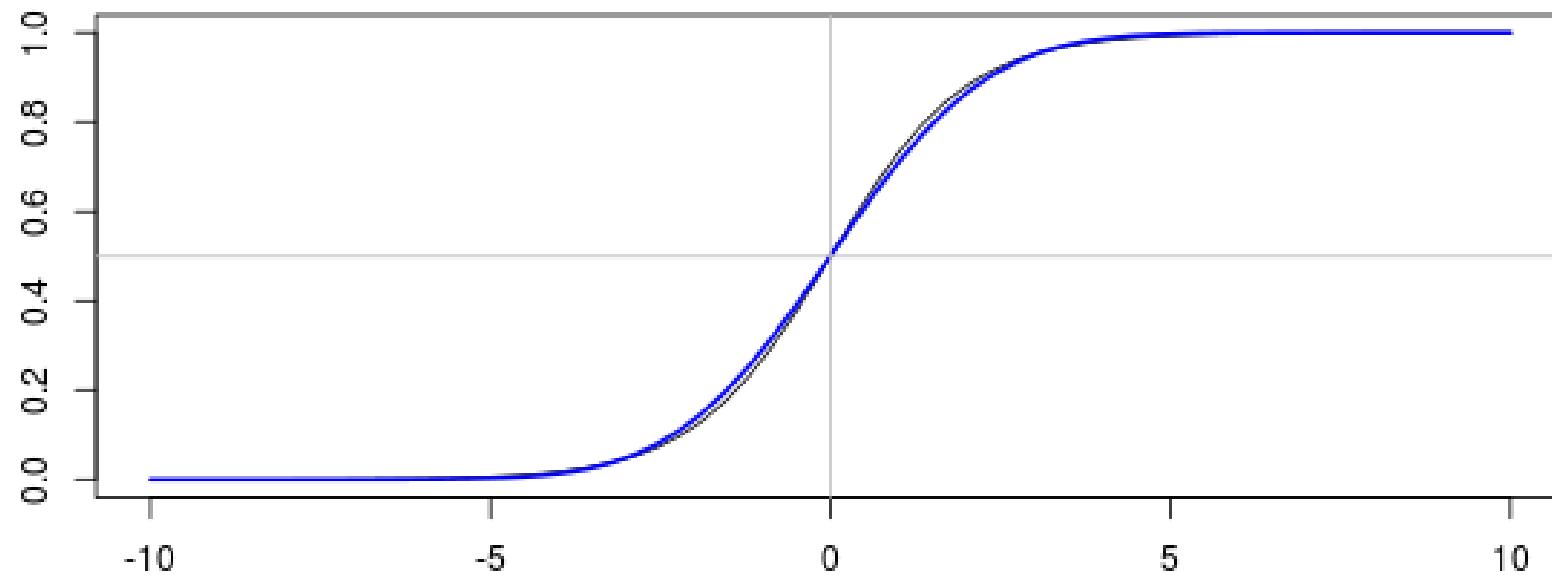
La fonction logistique

$$\Phi(x) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^x \exp\left(-\frac{z^2}{2}\right) dz, \forall x \in \mathbb{R}.$$



La fonction logistique

$$\Phi(x) = \frac{\sqrt{3}}{\pi\sqrt{2\pi}} \int_{-\infty}^x \exp\left(-\frac{3z^2}{2\pi^2}\right) dz, \forall x \in \mathbb{R}.$$



Qualité de l'ajustement

Pour juger de la validité du modèle, on peut considérer plusieurs quantités.

On peut définir la **déviance**, qui compare la log-vraisemblance du modèle et celle du modèle parfait (cf. section sur les GLMs),

$$D = -2 \sum_{i=1}^n [Y_i \log[\hat{\pi}_i] + (1 - Y_i) \log[1 - \hat{\pi}_i]]$$

Plus la variance est faible, meilleur est le modèle .

On peut aussi considérer la statistique de Pearson

$$\chi^2 = \sum_{i=1}^n \frac{(y_i - \hat{\pi}_i)^2}{\hat{\pi}_i[1 - \hat{\pi}_i]}$$

Étude des résidus

A l'aide de la relation précédante, on peut définir

$$\hat{\varepsilon}_i = \frac{Y_i - \hat{\pi}_i}{\sqrt{\hat{\pi}_i[1 - \hat{\pi}_i]}}$$

On parlera de **résidus de Pearson**.

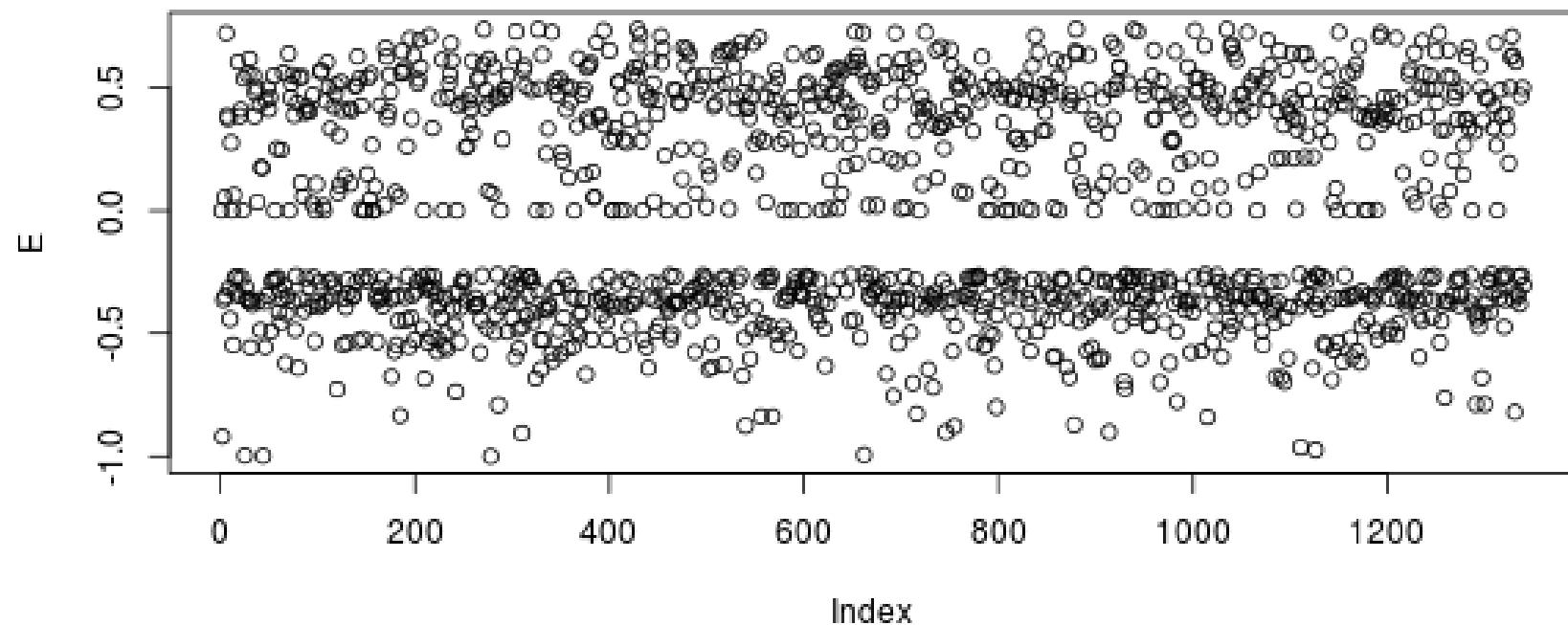
Il existe aussi les **résidus de déviance**.

$$\hat{\varepsilon}_i = \pm [Y_i \log[\hat{\pi}_i] + (1 - Y_i) \log[1 - \hat{\pi}_i]]^{\frac{1}{2}}$$

avec un signe positif si $Y_i \geq \hat{\pi}_i$, négatif sinon.

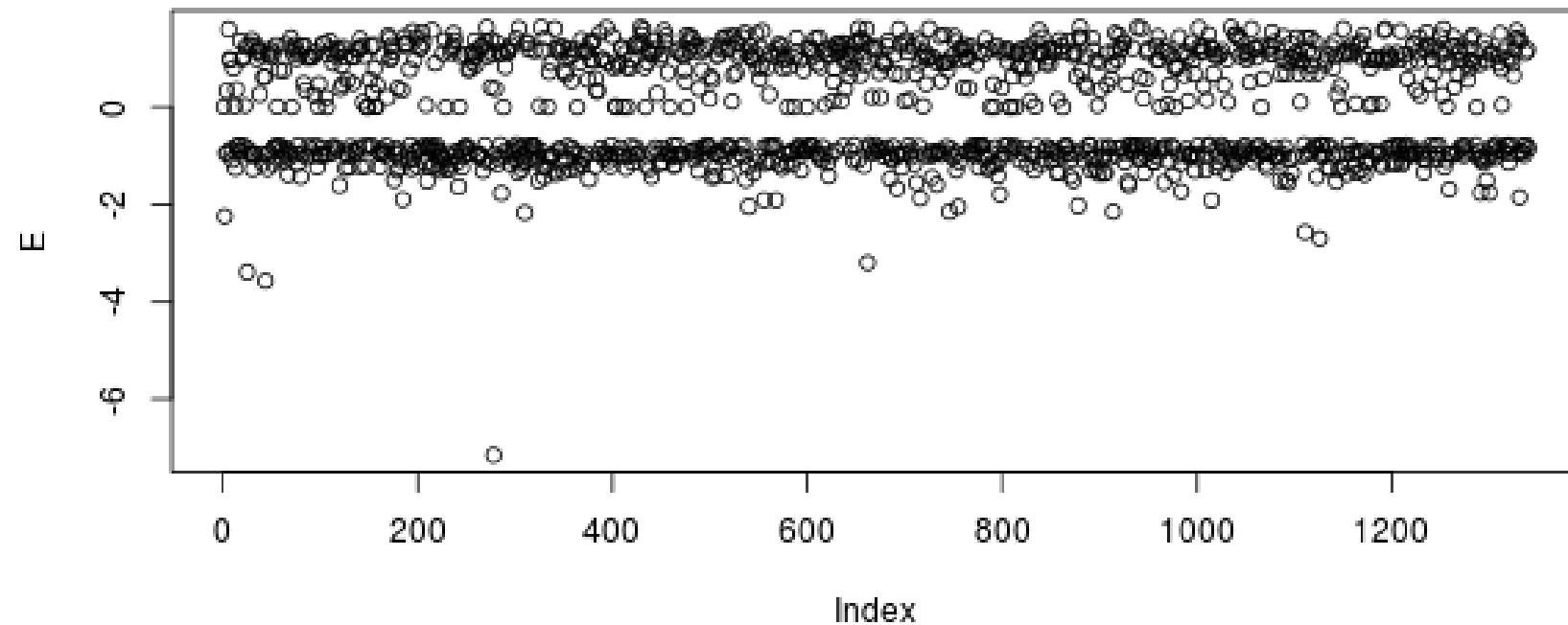
Visualisation des résidus

```
1 > reglogit = glm((ATTORNEY==1) ~ CLMSEX+LOSS, data=avocat, family=binomial)
2 > E = residuals(reglogit, type="response")
```



Visualisation des résidus

```
1 > reglogit = glm((ATTORNEY==1) ~ CLMSEX+LOSS, data=avocat, family=binomial)
2 > E = residuals(reglogit, type="deviance")
```



Prédiction avec un modèle logistique

Sous R, la prédiction sera (par défaut) une prédiction du score, i.e. $\mathbf{x}^\top \hat{\boldsymbol{\beta}}$.

```
1 > predict(reglogit, newdata=data.frame(LOSS=10, CLMSEX="M"))
2
3 1
4 2.530417
```

Il faut alors utiliser la loi logistique pour prévoir π ,

```
1 > exp(predict(reglogit, newdata=data.frame(LOSS=10, CLMSEX="M")))/
2 + (1+exp(predict(reglogit, newdata=data.frame(LOSS=10, CLMSEX="M"))))
3
4 1
5 0.9262468
```

ou utiliser directement

```
1 > predict(reglogit, newdata=data.frame(LOSS=10, CLMSEX="M"),
2 + type="response")
3
4 1
5 0.9262468
```

Modélisation de la survenance de sinistre

On peut aussi modéliser la survenance de sinistres

```

1 > camping=read.table("http://freakonometrics.free.fr/baseN.txt",
2 +                               header=TRUE ,sep=";")
3 > Y <- camping$nombre
4 > table(Y)
5
6      0      1
7 60155 3264

```

i.e. 5.1467% des polices sont touchées par un sinistre.

```

1 > reg=glm(Y~1,family=binomial(link="logit"))
2 > predict(reg,type="response") [1]
3
4 1
5 0.05146723

```

Prise en compte de l'exposition

Parfois, les polices d'assurance ne sont pas observées pendant un an, car les assurés peuvent les résilier. On peut malgré tout utiliser ces informations partielles (censurées, car on cherche une fréquence annuelle).

```

1 > E=camping$exposition
2 > table(E)
3 E
4   0.2   0.25   0.33   0.5      1
5     5       48     309   1804  61253

```

Si on suppose que la survenance d'un évènement est modélisé par un processus de Poisson (absence de mémoire), alors

$$\mathbb{P}(Y = 0) = \mathbb{P}(N = 0)^E$$

Si $N \sim \mathcal{B}(\pi)$,

$$\mathbb{P}(Y_i = y_i) = \pi^{e_i y_i} [1 - \pi^{e_i}]^{1-y_i}$$

et la vraisemblance est alors

$$\log \mathcal{L}(\pi; \mathbf{y}, \mathbf{e}) = \sum_{i=1}^n y_i e_i \log[\pi] + ([1 - y_i]) \log[1 - \pi^{e_i}]$$

de telle sorte que

$$\frac{\partial \log \mathcal{L}(\pi; \mathbf{y}, \mathbf{e})}{\partial \pi} = \sum_{i=1}^n \frac{y_i e_i}{\pi} - \frac{(1 - y_i) e_i \pi^{e_i - 1}}{1 - \pi^{e_i}} = 0.$$

```

1 > logL = function(p){
2 +   -sum(log(dbinom(Y, size=1, prob=1-(1-p)^E)))
3 +
4 > opt=optimize(f=logL, interval=c(1e-5,1e-1))
5 > opt$minimum
6 [1] 0.05225447

```

Plus généralement, on peut alors une régression logistique, en modélisant la non

survenance de sinistre,

$$\mathbb{P}[Y = 0 | \mathbf{X}, E] = \left(\frac{\exp(\mathbf{X}^\top \boldsymbol{\beta})}{1 + \exp(\mathbf{X}^\top \boldsymbol{\beta})} \right)^E$$

Si $\mathbb{P}[Y = 0] \sim 0$,

$$\mathbb{P}[Y = 1] = 1 - \mathbb{P}[Y = 0] = 1 - (1 - \mathbb{P}[N = 1])^E$$

de telle sorte que

$$\mathbb{P}[Y = 1] \sim 1 - 1 + E \cdot \mathbb{P}[N = 1] = E \cdot \mathbb{P}[N = 1]$$

```

1 > reg=glm(Y~1+offset(log(E)),family=binomial(link="log"))
2 > predict(reg,type="response")[1]
3
4 1
5 0.05212615

```

Prise en compte de l'exposition

Remarque ici, l'exposition est clairement endogène

```
1 > with(camping,table(nombre, exposition))
2
3   exposition
4
5   nombre    0.2    0.25    0.33    0.5      1
6
7   0          2       29     215    1230  58679
8
9   1          3       19      94     574   2574
```

Régression logistique et variables catégorielles

Dans la base, on dispose de variables catégorielles, la situation familiale, le sexe, la marque du véhicule, le type de paiement pour la prime, etc.

Formellement, X_j prend les modalités $\{x_{j,1}, \dots, x_{j,d_j}\}$, e.g. X_2 ,

$$\pi_i = \text{logit}^{-1} \left(\beta_0 + \beta_1 X_1 + \sum_{k=2}^{d_2} \beta_{2,k} \mathbf{1}(X_2 = x_{2,k}) \right)$$

On perd généralement la première modalité,

```
1 > levels(camping$situationfamiliale)
2 [1] "Celiba"    "Marie"     "Veuf/Div"
```

Régression logistique et variables catégorielles

```

1 > reg=glm(nombre~ageconducteur+situationfamiliale,data=camping,family
2   =binomial)
3
4 Coefficients:
5
6 (Intercept)           -3.8556504   0.1095669  -35.190  < 2e-16 ***
7 ageconducteur          0.0184951   0.0018124   10.205  < 2e-16 ***
8 situationfamilialeMarie -0.2028149   0.0484682   -4.184  2.86e-05 ***
9 situationfamilialeVeuf /Div -0.0006492   0.0992486   -0.007      0.995

```

Fusion de modalités pour une variable catégorielle

```

1 > levels(camping$situationfamiliale)=c("Celib/Veuf/Div","Marie",
2   Celib/Veuf/Div")
3 > reg=glm(nombre~ageconducteur+situationfamiliale,data=camping,family
4   =binomial)
5 > summary(reg)

5 Coefficients:
6
7             Estimate Std. Error z value Pr(>|z|)
8 (Intercept) -3.855646  0.109566 -35.190 < 2e-16 ***
9 ageconducteur 0.018493  0.001779  10.392 < 2e-16 ***
10 situationfamilialeMarie -0.202681  0.043973 -4.609 4.04e-06 ***

```

Régression et variables continues, *natura non facit saltus*

Les variables continues apparaissent forcément de manière linéaire

$$\pi_i = \text{logit}^{-1}(\beta_0 + \beta_1 x_1) \sim \exp(\beta_0 + \beta_1 x_1)$$

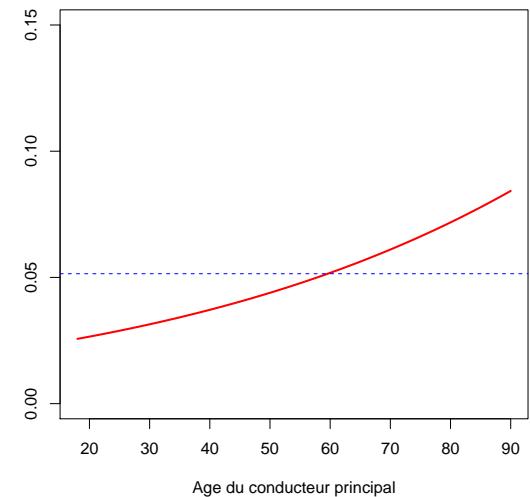
```

1 > reg=glm(nombre~ageconducteur,data=camping,family=binomial)
2 > summary(reg)
3
4 Coefficients:
5             Estimate Std. Error z value Pr(>|z|)
6 (Intercept) -3.949435   0.108492 -36.40  <2e-16 ***
7 ageconducteur 0.017375   0.001771    9.81  <2e-16 ***

```

Régression et variables continues, *natura non facit saltus*

```
1 > reg=glm(nombre~ageconducteur , data=camping ,  
family=binomial)
```



Régression et variables continues

Dans la base, on dispose de variables catégorielles, mais aussi ‘continues’ comme l’âge du conducteur, l’ancienneté du permis, l’âge du véhicule.

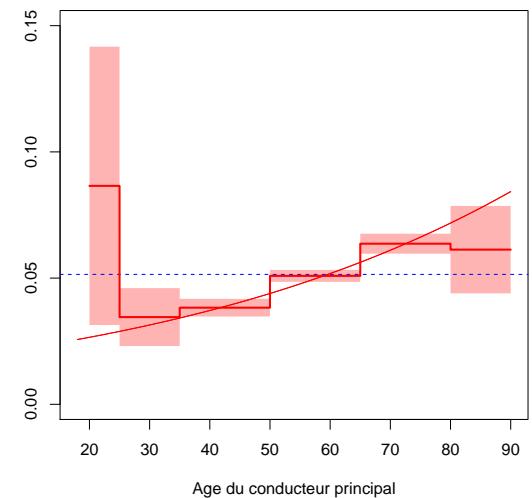
Historiquement, les assureurs aimait constituer des classes de risques, et donc cherchaient à découper les variables continues en classes. On peut considérer un découpage arbitraire en classes d’âges, e.g.

$$[18, 25), [25, 35), [35, 50), [50, 65), [65, 80), [80, 100).$$

Régression et variables continues, *natura non facit saltus*

```

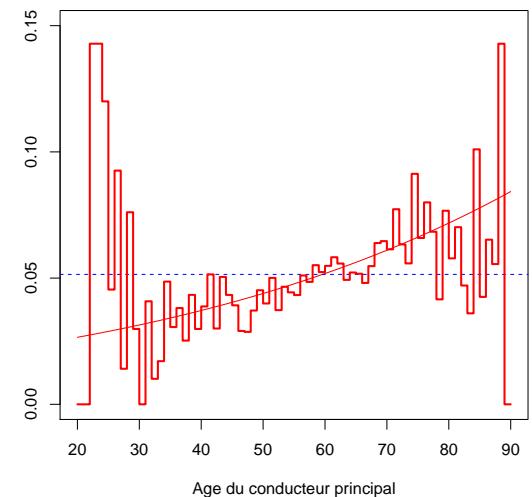
1 > AGE=c(18,25,35,50,65,80,100)
2 > cut_logistic = glm(nombre~cut(ageconducteur,
  breaks=AGE),data=camping,family=binomial)
```



On a l'impression que le modèle linéaire n'est peut-être pas adapté....

Régression et variables continues, *natura non facit saltus*

```
1 > f_logistic = glm(nombre~as.factor(ageconducteur)
  ,data=camping,family=binomial)
```



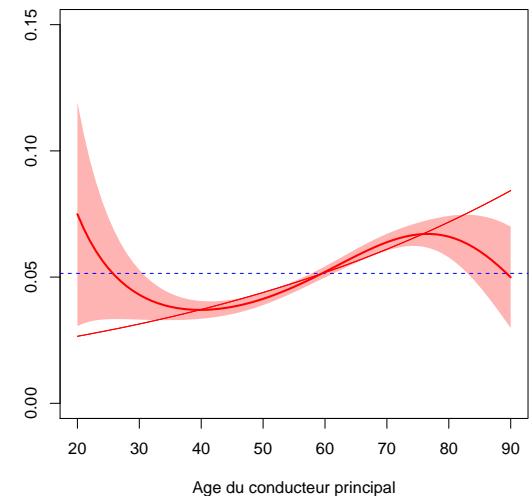
En prenant l'âge comme variable catégorielle, on retrouve l'effet non linéaire

Régression et variables continues, *natura non facit saltus*

```

1 > library(splines)
2 > bs_logistic = glm(nombre~bs(ageconducteur), data=
  camping, family=binomial(link="log"))

```



L'utilisation de splines semble confirmer la non-linéarité.

Lissage et régression locale

Si on veut prédire $\mathbb{E}[Y|X = x]$, on peut construire un modèle au voisinage de x

```
1 > h=1  
2 > w=(abs(db$x-x)<=h)*1  
3 > regl=lm(y~1,data=db,weights=w)
```

Lissage et régression locale

```
1 > h=1  
2 > w=(abs(db$x-x)<=h)*1  
3 > regl=lm(y~x,data=db,weights=w)
```

Lissage et régression locale - à noyau

```
1 > h=1  
2 > w=dnorm(db$x,x, sd=h)  
3 > regl=lm(y~x, data=db, weights=w)
```

Prédicteur linéaire

Si les modèles sont non-linéaires, on construit souvent un **prédicteur linéaire**.

Dans le modèle linéaire

$$\hat{\mathbf{y}} = \mathbf{X} \hat{\boldsymbol{\beta}} = \underbrace{\mathbf{X} [\mathbf{X}^\top \mathbf{X}]^{-1} \mathbf{X}^\top}_{\mathbf{H}} \mathbf{y}$$

Plus généralement, on a un prédicteur linéaire si la prévision \hat{y} au point \mathbf{x} s'écrit

$$\hat{y} = \hat{m}(\mathbf{x}) = \sum_{i=1}^n S_{\mathbf{x},i} y_i = \mathbf{S}_{\mathbf{x}}^\top \mathbf{y}$$

où $\mathbf{S}_{\mathbf{x}}$ est une vecteur de poids (**smoother vector**), que l'on peut associer à une matrice $n \times n$ de lissage

$$\hat{\mathbf{y}} = \mathbf{S}\mathbf{y}$$

où les prévisions sont faites aux points \mathbf{x}_i 's.

Degrés de liberté et complexité de modèle

Pour le modèle linéaire

$$\mathbf{S}_{\mathbf{x}} = \mathbf{H}_{\mathbf{x}} = \mathbf{X}[\mathbf{X}^T \mathbf{X}]^{-1} \mathbf{x}$$

Notons que

$$T = \frac{\|\mathbf{S}\mathbf{Y} - \mathbf{H}\mathbf{Y}\|}{\text{trace}([\mathbf{S} - \mathbf{H}]^T [\mathbf{S} - \mathbf{H}])}$$

est souvent utilisé comme statistique d'un test de linéarité: sous l'hypothèse de modèle linéaire, T suit une loi de Fisher.

Pour les prédicteurs linéaires, $\text{trace}(\mathbf{S})$ est souvent appelé **nombre équivalent de paramètres** et peut être relié à $n -$ nombre de degrés de liberté effectifs (cf [Ruppert et al. \(2003\)](#)).

Lissage, polynômes et *b-splines*

De manière générale, on va décomposer $\mathbb{E}[Y|X = x]$ dans une base de fonctions

$$m(\boldsymbol{x}) = \sum_{m=0}^M \beta_m h_m(\boldsymbol{x}),$$

On va retrouver $h_m(x) = x^m$ pour une régression polynomiale, ou $h_m(x) = (x - s_m)_+$ pour des ensembles de noeuds s_m (avec ici des splines linéaires, mais on peut considérer des splines quadratiques ou cubiques).

Lissage et régression polynômiale

```
1 > regp=lm(y~poly(x,M),data=db)
```

Lissage et régression par splines linéaires

```
1 > bs_reg_lin=lm(y~bs(x,df=M,degree=1),data=db)
```

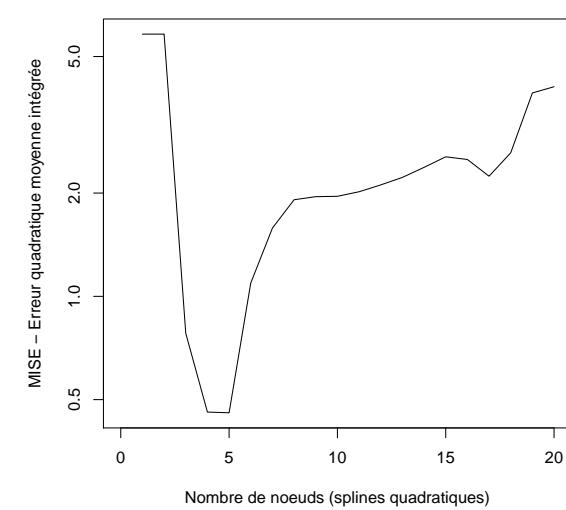
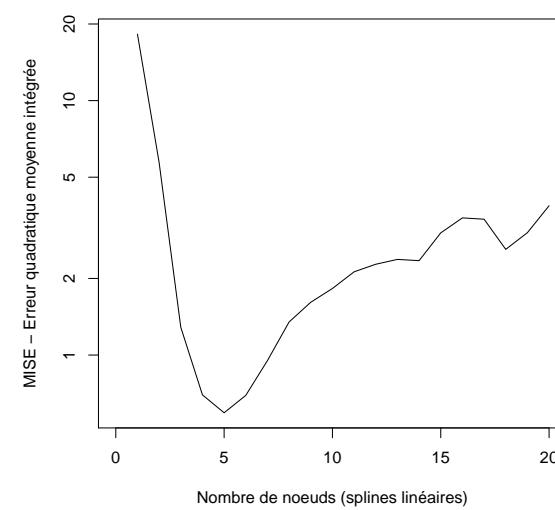
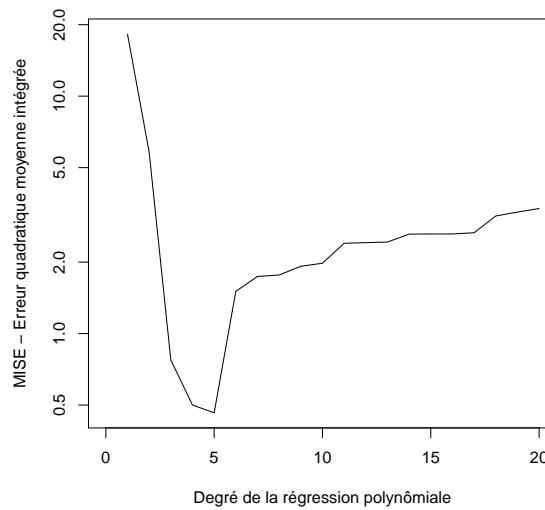
Lissage et régression par splines quadratiques

```
1 > bs_reg_quad=lm(y~bs(x,df=M,degree=2),data=db)
```

Lissage et tradeoff biais-variance

Problème classique en statistique : choix de M (degré du polynôme ou nombre de noeud), **tunning parameter**,

$$\text{mse}_M = \frac{1}{n} \sum_{i=1}^n [y_i - \hat{m}_M(x_i)]^2 = \text{variance} + \text{biais}^2$$



Lissage et surapprentissage

- M trop grand: variance importante, biais faible
- M trop faible: biais important, variance faible

On retrouve ici l'idée de **surapprentissage** avec un modèle trop complexe. Autre soucis,

$$\text{mse}_M = \frac{1}{n} \sum_{i=1}^n [y_i - \hat{m}_M(\mathbf{x}_i)]^2 \rightarrow ?$$

On ne peut invoquer la loi de grands nombres car $\hat{m}_M(\cdot)$ dépend de tous les (\mathbf{x}_i, y_i) .

Surapprentissage et validation croisée

Considérons un échantillon $\mathbf{x} = (x_1, \dots, x_n)$ et une statistique $\widehat{\theta}$. Posons $\widehat{\theta}_n = \widehat{\theta}(\mathbf{x})$

Le **Jackknife** est utilisé pour réduire le biais: posons $\widehat{\theta}_{(-i)} = \widehat{\theta}(\mathbf{x}_{(-i)})$, et

$$\tilde{\theta} = \frac{1}{n} \sum_{i=1}^n \widehat{\theta}_{(-i)}$$

Si $\mathbb{E}(\widehat{\theta}_n) = \theta + O(n^{-1})$ alors $\mathbb{E}(\tilde{\theta}_n) = \theta + O(n^{-2})$.

Cf idé de la validation croisée **leave-one-out** avec $\widehat{m}(\cdot)$

$$\text{mse} = \frac{1}{n} \sum_{i=1}^n [y_i - \widehat{m}_{(-i)}(x_i)]^2$$

L'estimateur par **Bootstrap** est obtenu par rééchantillonnage: posons

$\widehat{\theta}_{(b)} = \widehat{\theta}(\mathbf{x}_{(b)})$, et $\tilde{\theta} = \frac{1}{n} \sum_{i=1}^n \widehat{\theta}_{(b)}$, où $\mathbf{x}_{(b)}$ est un échantillon de taille n , obtenu par tirage avec remise dans $\{x_1, \dots, x_n\}$. Cf [Efron \(1979\)](#).

Modèle Prédicatif et Généralisation

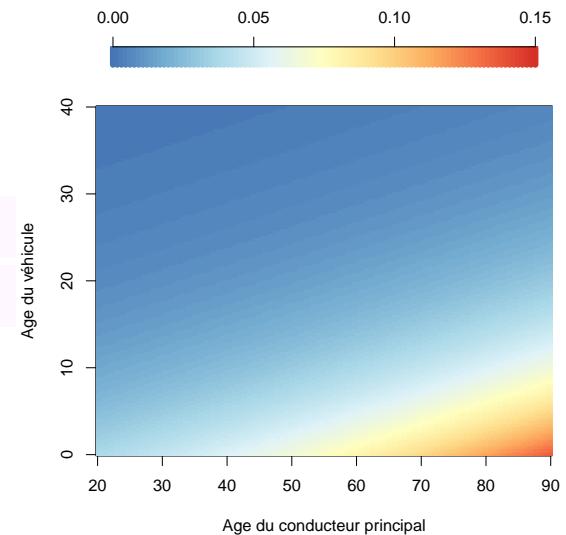
Autre approche classique, utiliser une partie des données pour apprendre, et l'autre pour valider le modèle. Le modèle ne doit pas (seulement) être bon sur la population d'apprentissage, il doit surtout l'être sur la population de test.

Effets croisés

Dans un modèle linéaire, avec deux variables, on considère une combinaison linéaire

$$\pi = \text{logit}^{-1}(\beta_0 + \beta_1 x_1 + \beta_2 x_2)$$

```
1 > reg_glm=glm(nombre~ageconducteur+agevehicule,
  data=camping,family=binomial)
```

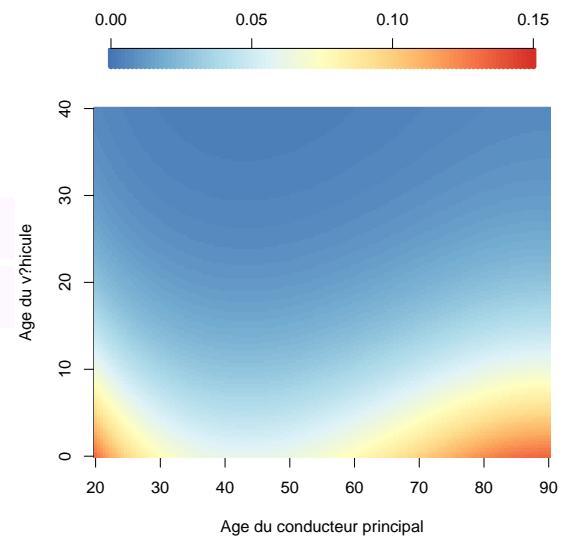


Effets croisés

En lissant les variables continues, le modèle reste additif

$$\pi = \text{logit}^{-1}(\beta_0 + s_1(x_1) + s_2(x_2))$$

```
1 > reg_add=glm(nombre~bs(ageconducteur)+bs(
+ agevehicule), data=camping, family=binomial)
```

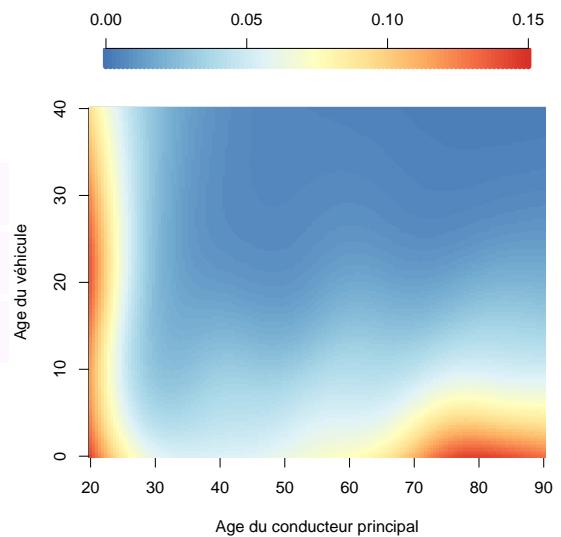


Effets croisés

On peut envisager un lissage bivarié

$$\pi = \text{logit}^{-1}(\beta_0 + s(x_1, x_2))$$

```
1 > library(mgcv)
2 > reg_gam=gam(nombre~s(ageconducteur,agevehicule),
   data=camping,family=binomial)
```

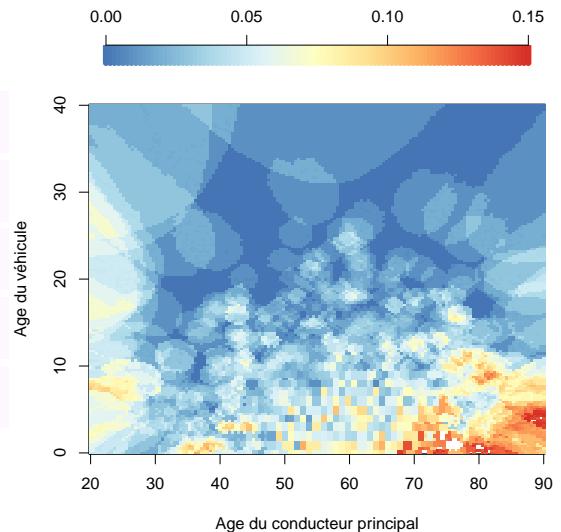


Les k -plus proches voisins

Parmi les modèles nonparamétriques classiques, on peut utiliser les k plus proches voisins

```

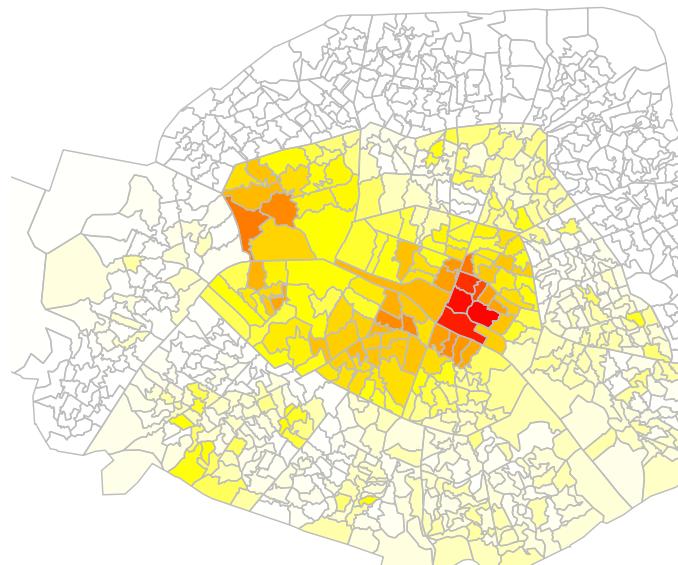
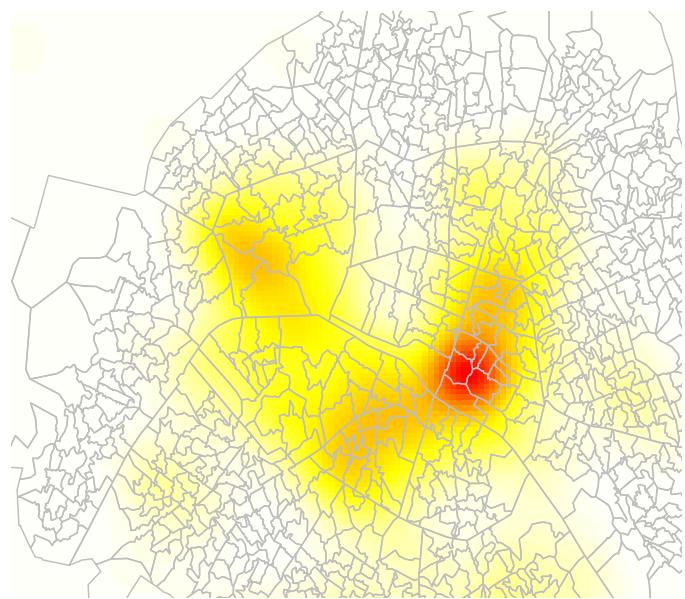
1 > library(caret)
2 > sc=sd(camping$ageconducteur)
3 > sv=sd(camping$agevehicule)
4 > knn=knn3((nombre==1)~I(ageconducteur/sc)+I(
   agevehicule/sv),data=camping,k=100)
```



(attention, la distance ne doit pas être sensible aux unités...)

Lissage spatial

On peut aussi utiliser du lissage spatial, quand les données sont géolocalisées (lieu d'habitation, localisation d'un accident, etc), cf `kde()` de `library(ks)`, ou par *k*-plus proche voisin quand on dispose d'une zone géographique `knn2nb()` de `library(spdep)`



De la loi Binomiale à la loi Multinomiale

Pour une loi de Bernoulli, $y \in \{0, 1\}$,

$$\mathbb{P}(Y = 1) = \frac{e^{\mathbf{X}^\top \boldsymbol{\beta}}}{1 + e^{\mathbf{X}^\top \boldsymbol{\beta}}} = \frac{p_1}{p_0 + p_1} \propto p_1 \text{ et } \mathbb{P}(Y = 0) = \frac{1}{1 + e^{\mathbf{X}^\top \boldsymbol{\beta}}} = \frac{p_0}{p_0 + p_1} \propto p_0$$

Pour une loi multinomiale, $y \in \{A, B, C\}$

$$\mathbb{P}(X = A) = \frac{p_A}{p_A + p_B + p_C} \propto p_A \text{ i.e. } \mathbb{P}(X = A) = \frac{e^{\mathbf{X}^\top \boldsymbol{\beta}_A}}{e^{\mathbf{X}^\top \boldsymbol{\beta}_A} + e^{\mathbf{X}^\top \boldsymbol{\beta}_B} + e^{\mathbf{X}^\top \boldsymbol{\beta}_C} + 1}$$

$$\mathbb{P}(X = B) = \frac{p_B}{p_A + p_B + p_C} \propto p_B \text{ i.e. } \mathbb{P}(X = B) = \frac{e^{\mathbf{X}^\top \boldsymbol{\beta}_B}}{e^{\mathbf{X}^\top \boldsymbol{\beta}_A} + e^{\mathbf{X}^\top \boldsymbol{\beta}_B} + e^{\mathbf{X}^\top \boldsymbol{\beta}_C} + 1}$$

$$\mathbb{P}(X = C) = \frac{p_C}{p_A + p_B + p_C} \propto p_C \text{ i.e. } \mathbb{P}(X = C) = \frac{1}{e^{\mathbf{X}^\top \boldsymbol{\beta}_A} + e^{\mathbf{X}^\top \boldsymbol{\beta}_B} + e^{\mathbf{X}^\top \boldsymbol{\beta}_C} + 1}$$

Courbe ROC, receiver operating characteristic

On dispose d'observations Y_i prenant les valeurs 0 (on dira **négatives**) ou 1 (on dira **positives**).

On a construit un modèle qui prédit π_i . En se fixant un seuil s , si $\hat{\pi}_i \leq s$, on prédit $\hat{Y}_i = 0$, et si $\hat{\pi}_i > s$, on prédit $\hat{Y}_i = 1$.

Le modèle sera bon si les positifs sont prédits positifs, et les négatifs sont prédits négatifs.

Le choix du seuil s permettra de minimiser soit les faux positifs, soit les faux négatifs.

	$Y = 1$	$Y = 0$
$\hat{Y}_i = 1$	TP	FP
$\hat{Y}_i = 0$	FN	TN

où **TP** sont les True Positive, **FP** les False Positive.

Courbe ROC et Qualité d'un Classifieur

	$Y = 0$	$Y = 1$	prevalence
$\hat{Y} = 0$	true negative N_{00}	false negative (type II) N_{01}	negative predictive value $NPV = \frac{N_{00}}{N_0.}$ $FOR = \frac{N_{01}}{N_0.}$
$\hat{Y} = 1$	false positive (type I) N_{10}	true positive N_{11}	false discovery rate $FDR = \frac{N_{10}}{N_1.}$ positive predictive value $PPV = \frac{N_{11}}{N_1.}$ (precision)
negative likelihood ratio $LR^- = FNR/TNR$	true negative rate $TNR = \frac{N_{00}}{N_0.}$ (specificity)	false negative rate $FNR = \frac{N_{01}}{N_1.}$	
positive likelihood ratio $LR^+ = TPR/FPR$	false positive rate $FPR = \frac{N_{10}}{N_0.}$ (fall out)	true positive rate $TPR = \frac{N_{11}}{N_1.}$ (sensitivity)	
diagnostic odds ratio $= LR^+/LR^-$			

Courbe ROC et Qualité d'un Classifieur

La courbe ROC est construite à partir de

```
1 > logistic <- step(glm(PRONO~.,data=myocarde,family=binomial))
2 > Y=myocarde$PRONO
3 > S=as.numeric(predict(logistic,type="response"))
```

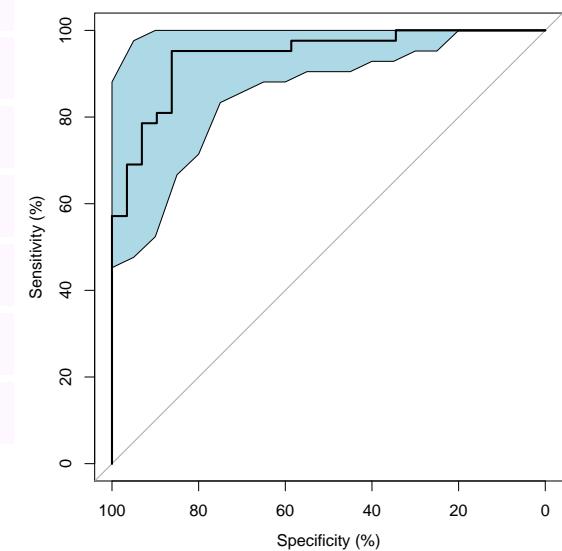
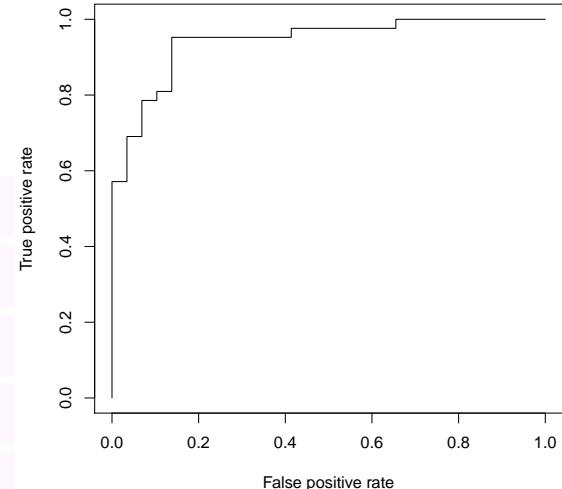
Courbe ROC

avec la régression logistique sur les infarctus

```

1 > logistic <- step(glm(F,data=camping,family=
  binomial))
2 > Y=camping$nombre
3 > S=predict(logistic,type="response")
4 > library(ROCR)
5 > pred=prediction(S,Y)
6 > plot(performance(pred,"tpr","fpr"))
7 > library(pROC)
8 > roc=plot.roc(Y,S,main="",
9 + percent=TRUE,ci=TRUE)
10 > roc.se=ci.se(roc,specificities=seq(0,100,5))
11 > plot(roc.se,type="shape",col="light blue")

```



Mesures construites à partir d'une courbe ROC

L'aire sous la courbe ROC, **Area Under the Curve AUC**, peut être interprétée comme la probabilité qu'une observation soit mieux prédite qu'une prédition purement aléatoire cf [Swets, Dawes & Monahan \(2000\)](#)

Beaucoup d'autre mesures peuvent être construites

```

1 > library(hmeasures)
2 > HMeasure(Y,S)$metrics[,1:5]
3 Class labels have been switched from (Death, Survival) to (0,1)
4          H      Gini      AUC      AUCH        KS
5 scores 0.7323154 0.8834154 0.9417077 0.9568966 0.8144499

```

avec la *H*-mesure (cf [hmeasure](#)), Gini et AUC, ainsi que l'aire sous l'enveloppe convexe de la courbe (**AUCH**).

Mesures construites à partir d'une courbe ROC

La statistique κ compare une **Observed Accuracy** avec une **Expected Accuracy** (random chance), cf [Landis & Koch \(1977\)](#).

	$Y = 0$	$Y = 1$	
$\hat{Y} = 0$	TN	FN	TN+FN
$\hat{Y} = 1$	FP	TP	FP+TP
	TN+FP	FN+TP	n

On va comparer les matrices de confusion observée et aléatoire (cf. chi-deux)

	$Y = 0$	$Y = 1$	
$\hat{Y} = 0$	25	3	28
$\hat{Y} = 1$	4	39	43
	29	42	71

	$Y = 0$	$Y = 1$	
$\hat{Y} = 0$	11.44	16.56	28
$\hat{Y} = 1$	17.56	25.44	43
	29	42	71

$$\text{total accuracy} = \frac{TP + TN}{n} \sim 90.14\%$$

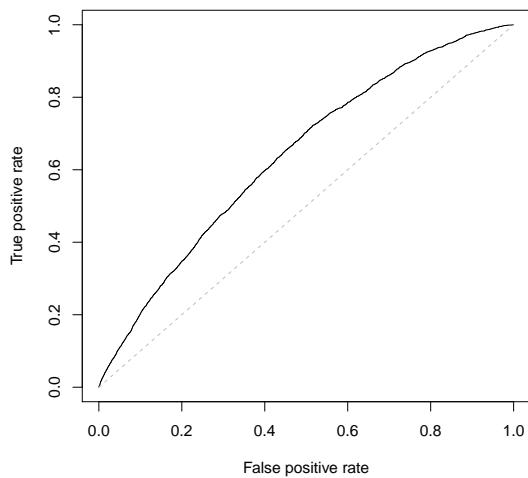
$$\text{random accuracy} = \frac{[TN + FP] \cdot [TP + FN] + [TP + FP] \cdot [TN + FN]}{n^2} \sim 51.93\%$$

$$\kappa = \frac{\text{total accuracy} - \text{random accuracy}}{1 - \text{random accuracy}} \sim 79.48\%$$

```

1 > logistic <- glm(nombre ~ ., data=
   camping, family=binomial)
2 > Y=camping$nombre
3 > S=predict(logistic, type="response")
4 > pred=prediction(S, Y)
5 > plot(performance(pred, "tpr", "fpr"))

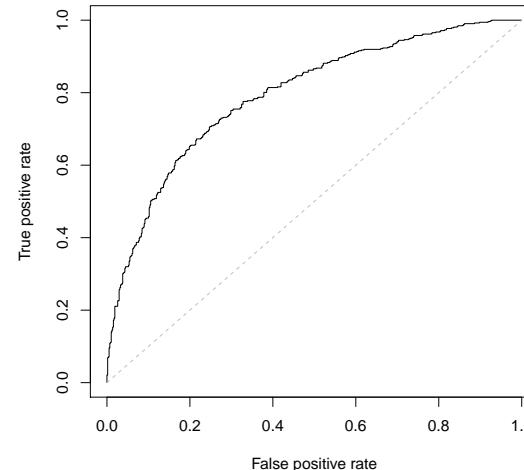
```



```

1 > logistic <- glm((ATTORNEY==2) ~
   ., data=avocat, family=binomial)
2 > Y=avocat$ATTORNEY
3 > S=predict(logistic, type="response")
4 > pred=prediction(S, Y)
5 > plot(performance(pred, "tpr", "fpr"))

```



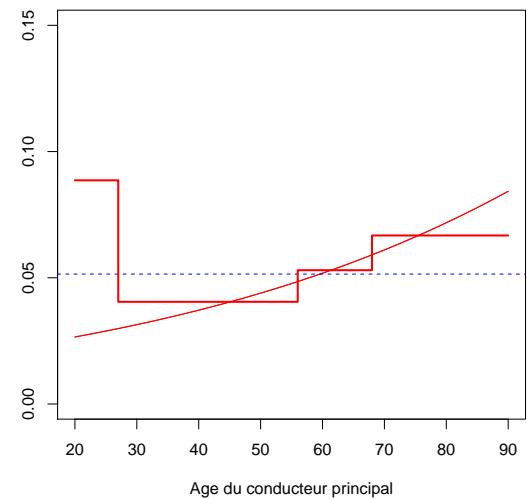
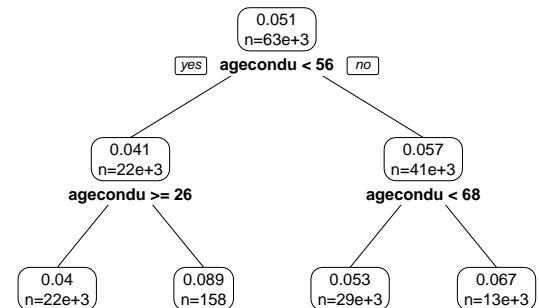
Les arbres de classification

Dans ces structures d'arbre, les feuilles représentent les valeurs de la variable-cible et les embranchements correspondent à des combinaisons de variables d'entrée qui mènent à ces valeurs (Source: wikipedia).

```

1 > library(rpart)
2 > tree=rpart((nombre==1) ~ ageconducteur , data=
   camping , cp=1e-4)
3 > library(rpart.plot)
4 > library(rattle)
5 > prp(tree , type=2 , extra=1)

```



Les arbres de classification

Considérons une feuille $\{N\}$, alors l'indice de Gini sur la feuille est

$$\text{gini}(N) = \frac{n_{x,y}}{n_x} \left(1 - \frac{n_{x,y}}{n_x} \right)$$

et si $\{N\}$ est partagée en $\{N_L, N_R\}$,

$$\text{gini}(N_L, N_R) = - \sum_{x \in \{L, R\}} \frac{n_x}{n} \sum_{y \in \{0, 1\}} \frac{n_{x,y}}{n_x} \left(1 - \frac{n_{x,y}}{n_x} \right)$$

Plus généralement, on dispose d'un indice d'impureté $\mathcal{I}(\cdot)$ et pour partager $\{N\}$ en $\{N_L, N_R\}$, on considère

$$\mathcal{I}(N_L, N_R) = \sum_{x \in \{L, R\}} \cdot \frac{n_x}{n} \mathcal{I}(N_x)$$

e.g. l'indice de Gini (utilisé dans CART, [Breiman et al. \(1984\)](#)), $-p(1-p)$ ou l'indice d'entropie (utilisé dans C4.5 et C5.0) $p \log(p)$.

On considère des feuilles de la forme $N_L = \{x_{i,\textcolor{red}{j}} \leq \textcolor{blue}{s}\}$ et $N_R = \{x_{i,\textcolor{red}{j}} > \textcolor{blue}{s}\}$

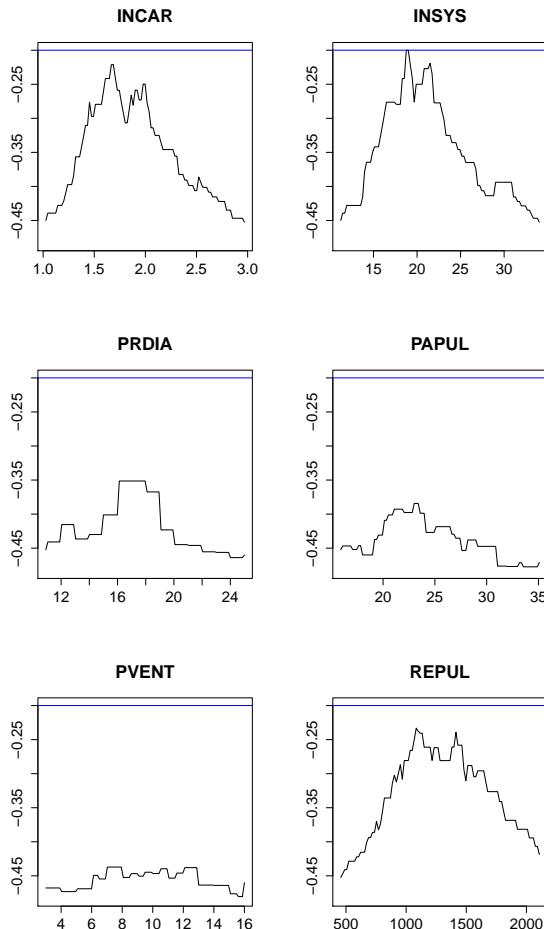
Les arbres de classification

```
1 > library(rpart)
2 > cart<-rpart(PRONO~., data = myocarde)
3 > library(rpart.plot)
4 > library(rattle)
5 > prp(cart, type=2, extra=1)
```

ou

```
1 > fancyRpartPlot(cart, sub="")
```

Les arbres de classification

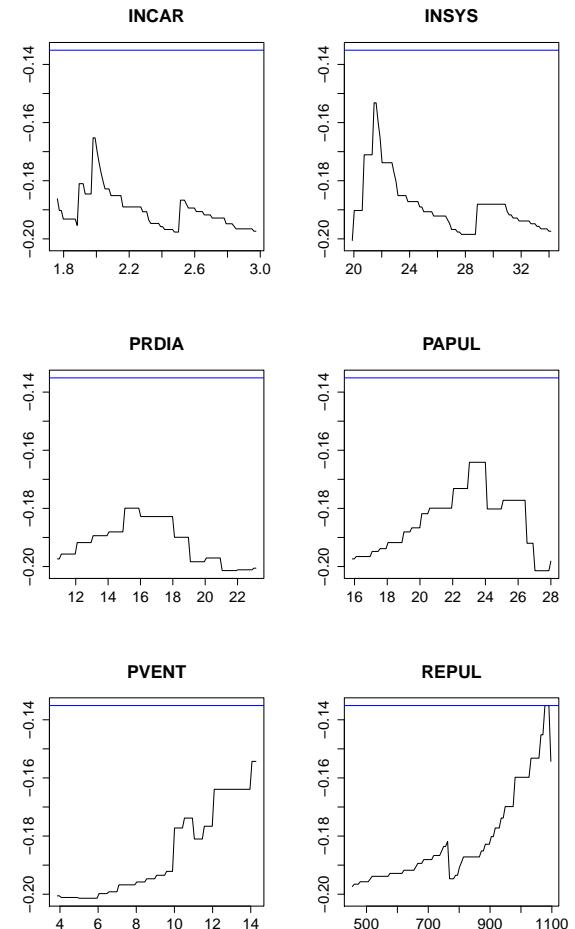


$$N_L: \{x_{i,j} \leq s\} \quad N_R: \{x_{i,j} > s\}$$

résoudre $\max_{j \in \{1, \dots, k\}, s} \{\mathcal{I}(N_L, N_R)\}$

← premier split

second split →



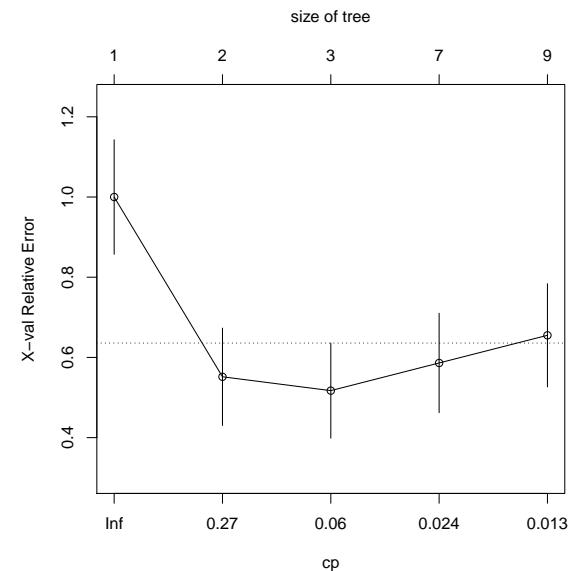
Taille des arbres (élagage, *pruning*)

On ne partage pas à un noeud si $\Delta\mathcal{I}$ est trop faible,

$$\Delta\mathcal{I}(N_L, N_R) = \mathcal{I}(N) - \mathcal{I}(N_L, N_R) = \mathcal{I}(N) - \left(\frac{n_L}{n} \mathcal{I}(N_L) - \frac{n_R}{n} \mathcal{I}(N_R) \right)$$

Remarque on note cp (complexity parameter) le seuil sur $\Delta\mathcal{I}(N_L, N_R)/\mathcal{I}(N)$ (gain relatif d'indice), cf [Breiman et al. \(1984\)](#).

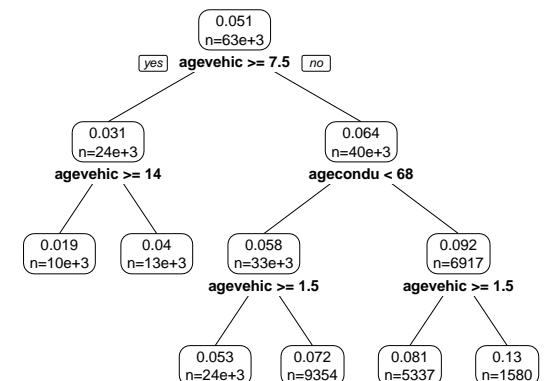
```
1 > cart <- rpart(PRONO ~ ., data = myocarde,
2 + minsplit=3)
3 > plotcp(cart)
```



Les arbres de classification

En diminuant la valeur de `cp`, on peut construire un arbre plus profond (on peut aussi diminuer `mincut`)

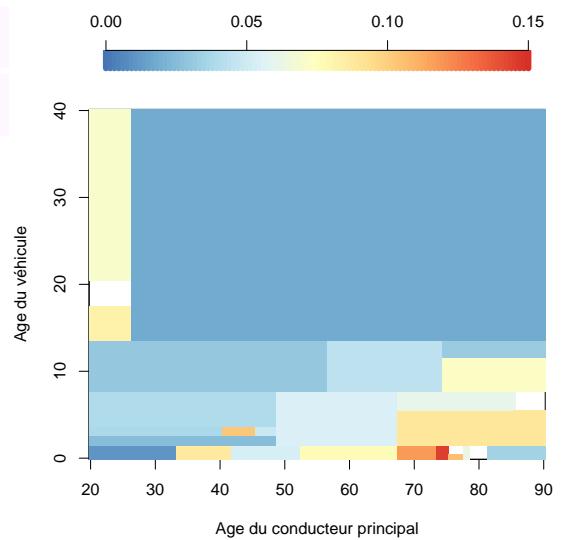
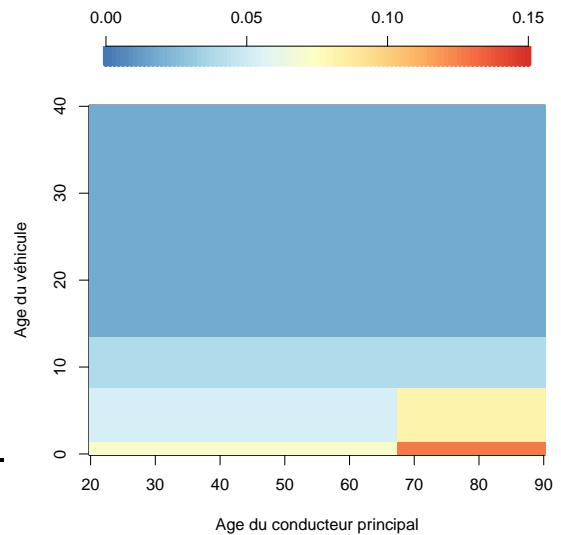
```
1 > tree=rpart((nombre==1)~., data=camping, cp=7e-4)
2 > prp(tree, type=2, extra=1)
```



Effets croisés

On peut utiliser les arbres pour visualiser des possibles effets croisés entre deux variables continues

```
1 > tree=rpart((nombre==1)~ageconducteur+agevehicule  
 , data=camping , cp=7e-4)
```



Méthodes de *Bagging* (ou ensachage)

Les arbres sont (relativement) peu robustes...

On peut construire plusieurs arbres de décision en ré-échantillonnant l'ensemble d'apprentissage [bootstrap], puis en construisant les arbres par une procédure de consensus[aggregation] ou en moyennant les prédictions.

→ marche aussi pour des GLM (cf dernière section du cours) mais est très populaire sur les arbres Breiman (1996)

Le bagging permet de définir la notion de marge,

$$\text{marge}_i = \frac{1}{B} \sum_{b=1}^B \mathbf{1}(\hat{y}_i = y_i) - \frac{1}{B} \sum_{b=1}^B \mathbf{1}(\hat{y}_i \neq y_i)$$

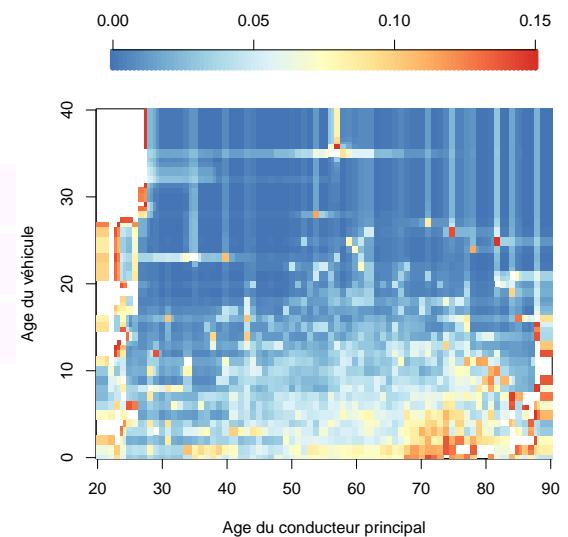
Remarque La probabilité que la ligne i ne soit pas dans un échantillon est $(1 - n^{-1})^n \rightarrow e^{-1} \sim 36.8\%$, cf ratio training / validation (2/3-1/3)

Bagging

```

1 > library(ipred)
2 > bag=bagging((nombre==1)~ageconducteur+
   agevehicule ,data=camping)

```



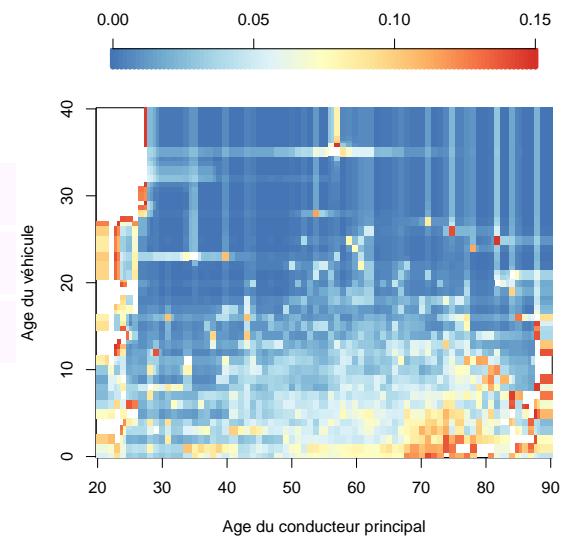
Les forêts aléatoires (*Random Forest*)

Générer des forêts en boostrappant puis en agrégeant pour obtenir une prévision globale est du *bagging*.

In the [random forest](#) algorithm, we combine Breiman's [bagging](#) idea and the random selection of features, introduced independently by [Ho \(1995\)](#) and [Amit & Geman \(1997\)](#)

Par défaut, on sélectionne \sqrt{k} variables explicatives.

```
1 > library(randomForest)
2 > rf=randomForest((nombre==1)~ageconducteur+
  agevehicule,data=camping)
```



Variable Importance

Étant donnée une forêt composée de M arbres, on pose

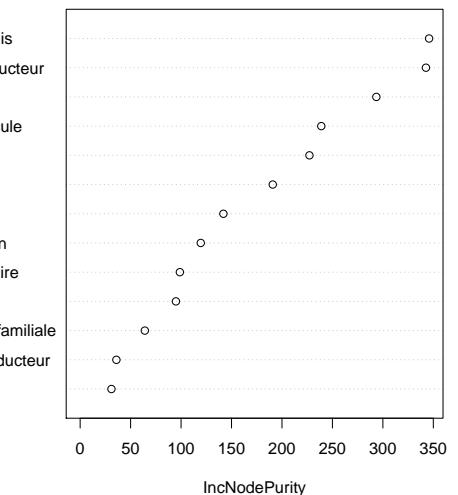
$I(X_k) = \frac{1}{M} \sum_m \sum_{t \in X_k} \frac{N_t}{N} \Delta \mathcal{I}(t)$ où la première somme est sur tous les arbres, et la seconde sur tous les noeuds splités suivant la variable X_k .

```

1 > RF=randomForest((nombre==1) ~ ., data = camping)
2 > varImpPlot(RF, main="")
3 > importance (RF)

4           IncNodePurity
5 ageconducteur      342.68711
6 agepermis          345.79078
7 sexeconducteur    35.90909
8 situationfamiliale 64.09047
9 habitation         119.52306
10 zone              227.22830
11 agevehicule        238.98409
12 proprietaire       98.81689
13 payment            94.92349
14 marque             293.54250
15 poids              141.96059
16 usage              31.00314

```



Partial Response Plots

On peut aussi construire, pour toutes les variables continues des **Partial Response Plots**,

$$x \mapsto \frac{1}{n} \sum_{i=1}^n \hat{\mathbb{E}}[Y|X_k = x, \mathbf{X}_{i,(k)} = \mathbf{x}_{i,(k)}]$$

```

1 > importanceOrder <- order(-RF$importance)
2 > names <- rownames(RF$importance)[importanceOrder
   ]
3 > for (name in names)
4 + partialPlot(RF, camping, eval(name), col="red",
   main="", xlab=name)

```

Gradient Boosting en Classification

Boosting is a machine learning ensemble meta-algorithm for reducing bias primarily and also variance in supervised learning, and a family of machine learning algorithms which convert weak learners to strong ones. (source: [Wikipedia](#))

Avec des notations très générales, on souhaite résoudre

$$m^* = \operatorname{argmin}\{\mathbb{E}[\ell(Y, m(\mathbf{X}))]\}$$

pour une fonction de perte ℓ ([loss function](#)). Heuristiquement, on construit itérativement des modèles sur les erreurs résiduelles

On ajuste un modèle sur \mathbf{y} , $m_1(\cdot)$ à partir de \mathbf{y} et \mathbf{X} , et on définit l'erreur $\varepsilon_1 = \mathbf{y} - m_1(\mathbf{X})$. On ajuste un modèle sur ε_1 , $m_2(\cdot)$ à partir de ε_1 et \mathbf{X} , et on définit l'erreur $\varepsilon_2 = \varepsilon_1 - m_2(\mathbf{X})$, etc. Puis on pose

$$m(\cdot) = \underbrace{m_1(\cdot)}_{\sim y} + \underbrace{m_2(\cdot)}_{\sim \varepsilon_1} + \underbrace{m_3(\cdot)}_{\sim \varepsilon_2} + \cdots + \underbrace{m_k(\cdot)}_{\sim \varepsilon_{k-1}}$$

'Gradient' Boosting ?

Dans un modèle linéaire $\mathbf{y} - m_k(\mathbf{X})$ est un résidu, qui correspond au gradient de $\frac{1}{2}[y - m(\mathbf{x})]^2$

La descente de Gradient vient de l'expansion de Taylor

$$\underbrace{f(\mathbf{x}_k)}_{\langle f, \mathbf{x}_k \rangle} \sim \underbrace{f(\mathbf{x}_{k-1})}_{\langle f, \mathbf{x}_{k-1} \rangle} + \underbrace{(\mathbf{x}_k - \mathbf{x}_{k-1})}_{\alpha} \underbrace{\nabla f(\mathbf{x}_{k-1})}_{\langle \nabla f, \mathbf{x}_{k-1} \rangle}$$

Mais là, ce que l'on considère, c'est

$$\underbrace{f_k(\mathbf{x})}_{\langle f_k, \mathbf{x} \rangle} \sim \underbrace{f_{k-1}(\mathbf{x})}_{\langle f_{k-1}, \mathbf{x} \rangle} + \underbrace{(f_k - f_{k-1})}_{\beta} \underbrace{?}_{\langle f_{k-1}, \nabla \mathbf{x} \rangle}$$

où ? pourrait être interprété comme un 'gradient'.

Boosting avec rétrécissement (*shrinkage*)

L'algorithme est ici

- on considère un modèle (simple) sur y , $\mathbf{y} = h_1(\mathbf{x})$
- on construit les résidus (avec un paramètre ν de shrinkage), $\boldsymbol{\varepsilon}_1 = \mathbf{y} - \nu h_1(\mathbf{x})$

à l'étape j ,

- on considère un modèle (simple) sur ε_j , $\boldsymbol{\varepsilon}_j = h_j(\mathbf{x})$
- on construit les résidus $\boldsymbol{\varepsilon}_{j+1} = \boldsymbol{\varepsilon}_j - \nu h_j(\mathbf{x})$

et on boucle. Finalement

$$\hat{\mathbf{y}} = \sum_{j=1}^M \nu_j h_j(\mathbf{x})$$

Boosting

Apprentissage avec des **stumps** (arbres)

```
1 > df$e <- df$y - mean(df$y)
2 > for(k in 1:1000){
3 + fit <- rpart(e ~ x,data=df)
4 + df$e <- df$y - v*predict(fit)
5 + YP <- cbind(YP,v*predict(fit))
6 +
7 > apply(YP ,1 ,sum)
```

Boosting

Apprentissage avec des **splines** (linéaires)

```

1 > library(freeknotsplines)
2 > df$e <- df$y - mean(df$y)
3 > for(k in 1:1000){
4 + noeuds = freelsgen(df$x,df$e,degree=1,numknot=2)
5 + fit <- lm(e ~ bs(x,degree=1,knots=noeuds@optknot
+ ,data=df)
6 + df$e <- df$y - v*predict(fit)
7 + YP <- cbind(YP,v*predict(fit))
8 +
9 > apply(YP,1,sum)

```

la localisation des noeuds est ici optimisée

Boosting for Classification

On cherche toujours $m^*(\cdot) = \operatorname{argmin}\{\mathbb{E}[\ell(Y, m(\mathbf{X}))]\}$. Classiquement, en machine learning, $y \in \{-1, +1\}$, et on considère comme fonction de perte

$$\ell(y, m(\mathbf{x})) = e^{-y \cdot m(\mathbf{x})}$$

on obtient l'algorithme dit **AdaBoot**, cf [Freund & Schapire \(1996\)](#).

Alors

$$\mathbb{P}[Y = +1 | \mathbf{X} = \mathbf{x}] = \frac{1}{1 + e^{2m^*(\mathbf{x})}}$$

qui rappelle la loi logistique... A l'étape k on résout

$$\operatorname{argmin} \left\{ \sum_{i=1}^n \underbrace{e^{y_i \cdot m_k(\mathbf{x}_i)}}_{\omega_{i,k}} \cdot e^{y_i \cdot h(\mathbf{x}_i)} \right\}$$

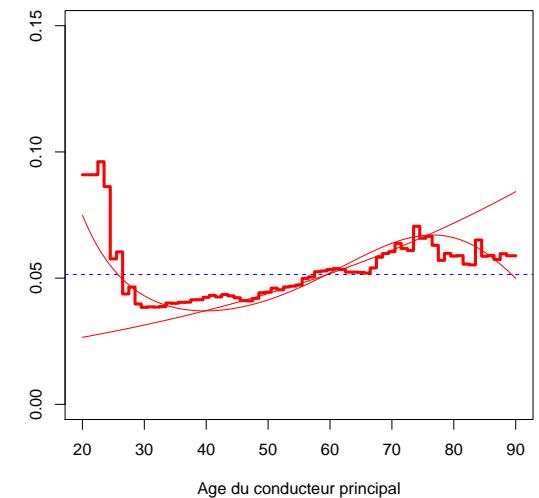
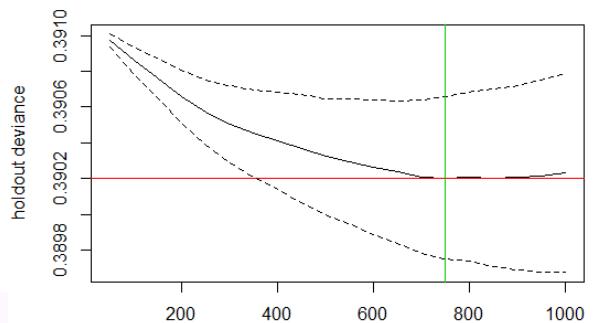
Boosting

Apprentissage avec des stumps (arbres)

```

1 > library(dismo)
2 > library(gbm)
3 > fit <- gbm.step(data=camping, gbm.x=1, gbm.y=13,
   family="bernoulli", tree.complexity=5,
   learning.rate=0.001, bag.fraction=0.5)
4 > predict(fit, type="response", n.trees=700)

```



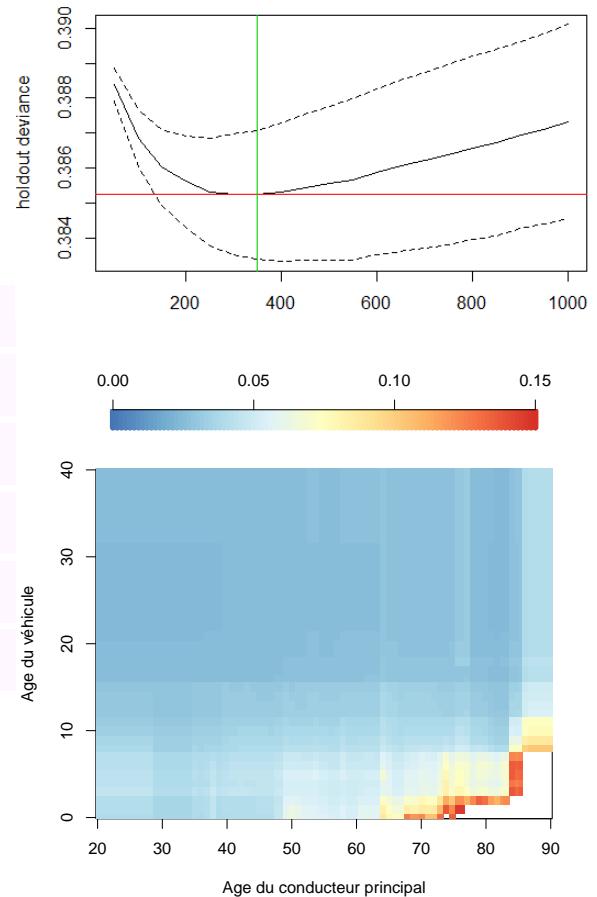
Boosting

Apprentissage avec des stumps (arbres)

```

1 > library(dismo)
2 > library(gbm)
3 > fit <- gbm.step(data=camping, gbm.x=c(1,7), gbm.
   y=13, family="bernoulli", tree.complexity=5,
   learning.rate=0.01, bag.fraction=0.5)
4 > predict(fit, type="response", n.trees=400)

```



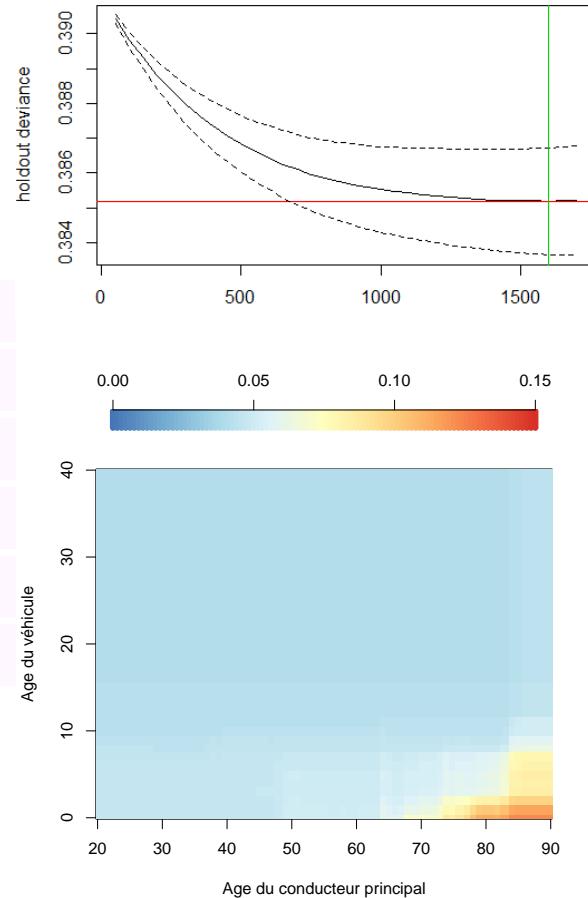
Boosting

Apprentissage avec des stumps (arbres)

```

1 > library(dismo)
2 > library(gbm)
3 > fit <- gbm.step(data=camping, gbm.x=c(1,7), gbm.
   y=13, family="bernoulli", tree.complexity=5,
   learning.rate=0.001, bag.fraction=0.5)
4 > predict(fit, type="response", n.trees=1500)

```



Support Vector Machine

Les SVMs ont été développés dans les années 90 à partir de [Vapnik & Lerner \(1963\)](#), cf [Vailant \(1984\)](#).

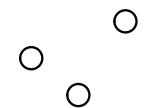
On suppose que les points sont linéairement séparables, i.e. il existe ω et b tels que

$$Y = \begin{cases} +1 & \text{if } \omega^\top \mathbf{x} + b > 0 \\ -1 & \text{if } \omega^\top \mathbf{x} + b < 0 \end{cases}$$

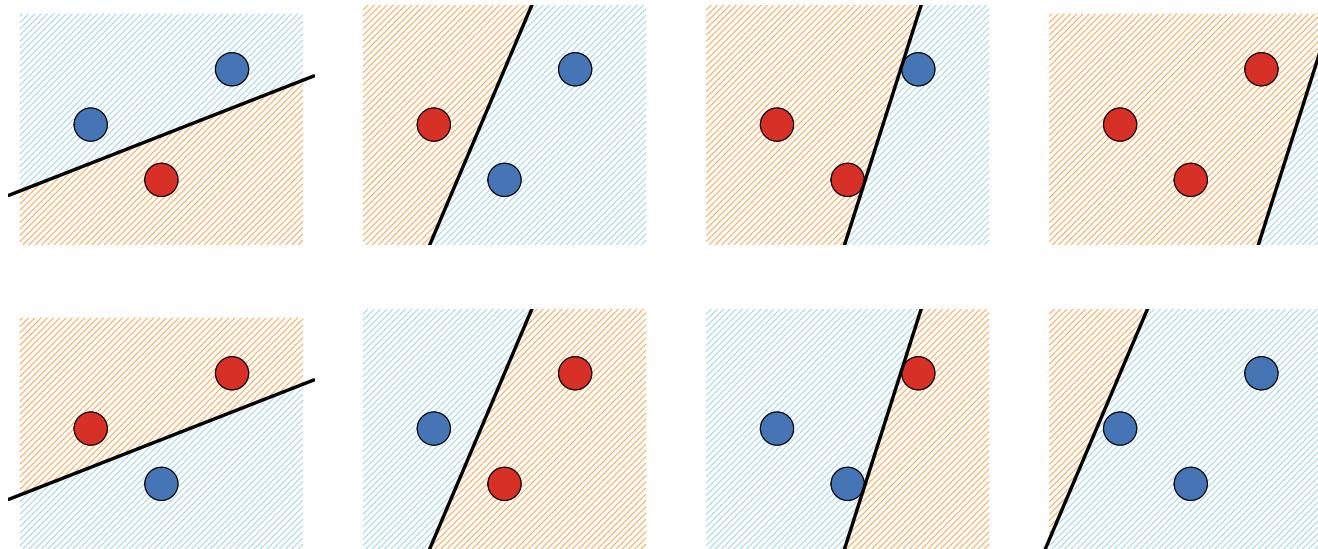
Problème: nombre infini de solutions, il faut en choisir une [bonne](#) qui sépare les deux nuages autant que possible.

Remarque derrière l'idée de séparabilité se cache la [dimension VC](#). Soit $\mathcal{H} : \{h : \mathbb{R}^d \mapsto \{-1, +1\}\}$. Alors \mathcal{H} est un [shatter](#) des points \mathbf{X} si toutes les dichotomies peuvent être obtenues.

E.g. avec trois point, toutes les configurations sont atteignables



Support Vector Machine



Mais avec les quatre point ci-contre, non: plusieurs configuration ne peuvent être obtenues (avec un séparateur linéaire)

Support Vector Machine

La dimension de Vapnik (VC) est la taille du plus grand ensemble séparé (shattered) de \mathbf{X} .

En pratique, où est le plan optimal ?

La distance de \mathbf{x}_0 au plan $\boldsymbol{\omega}^\top \mathbf{x} + b$ est

$$d(\mathbf{x}_0, H_{\boldsymbol{\omega}, b}) = \frac{\boldsymbol{\omega}^\top \mathbf{x}_0 + b}{\|\boldsymbol{\omega}\|}$$

et le **plan optimal** (dans le cas séparable) est

$$\operatorname{argmin} \left\{ \min_{i=1, \dots, n} d(\mathbf{x}_i, H_{\boldsymbol{\omega}, b}) \right\}$$

Support Vector Machine

On définit **vectors supports** comme les observations telles que

$$|\boldsymbol{\omega}^\top \mathbf{x}_i + b| = 1$$

La marge est la distance des supports vecteurs au plan.

La distance des vecteurs supports à $H_{\boldsymbol{\omega}, b}$ est $\|\boldsymbol{\omega}\|^{-1}$, et la marge est $2\|\boldsymbol{\omega}\|^{-1}$.
 → l'algorithme va minimiser l'inverse des marges, sous contrainte que le plan $H_{\boldsymbol{\omega}, b}$ sépare les points (ici ± 1) i.e.

$$\min \left\{ \frac{1}{2} \boldsymbol{\omega}^\top \boldsymbol{\omega} \right\} \text{ s.t. } Y_i(\boldsymbol{\omega}^\top \mathbf{x}_i + b) \geq 1, \quad \forall i.$$

Support Vector Machine

Problème difficile à résoudre (trop de contraintes: n)

→ on va résoudre ce problème d'optimisation dans l'espace dual

Dans l'[espace initial](#), on aurait

$$\boldsymbol{\omega} = \sum \alpha_i Y_i \mathbf{x}_i \text{ avec } \sum_{i=1} \alpha_i Y_i = 0.$$

Dans l'[espace dual](#), le programme s'écrit

$$\max \left\{ \sum_{i=1} \alpha_i - \frac{1}{2} \sum_{i=1} \sum_{j=1} \alpha_i \alpha_j Y_i Y_j \mathbf{x}_i^\top \mathbf{x}_j \right\} \text{ s.t. } \sum_{i=1} \alpha_i Y_i = 0.$$

qui peut s'écrire

$$\min_{\boldsymbol{\alpha}} \left\{ \frac{1}{2} \boldsymbol{\alpha}^\top \mathbf{Q} \boldsymbol{\alpha} - \mathbf{1}^\top \boldsymbol{\alpha} \right\} \text{ s.c. } \begin{cases} 0 \leq \alpha_i \quad \forall i \\ \mathbf{y}^\top \boldsymbol{\alpha} = 0 \end{cases}$$

où $\mathbf{Q} = [\mathbf{Q}_{i,j}]$ et $\mathbf{Q}_{i,j} = y_i y_j \mathbf{x}_i^\top \mathbf{x}_j$.

Support Vector Machine

Quid du cas non-séparable?

Ici on ne **peut pas** avoir $y_i(\boldsymbol{\omega}^\top \mathbf{x}_i + b) \geq 1 \ \forall i$.

On introduit la notion de **slack variables**,

$$\begin{cases} \boldsymbol{\omega}^\top \mathbf{x}_i + b \geq +1 - \xi_i & \text{lorsque } y_i = +1 \\ \boldsymbol{\omega}^\top \mathbf{x}_i + b \leq -1 + \xi_i & \text{lorsque } y_i = -1 \end{cases}$$

où $\xi_i \geq 0 \ \forall i$.

On commet une erreur de classification lorsque $\xi_i > 1$.

L'idée est de résoudre

$$\min \left\{ \frac{1}{2} \boldsymbol{\omega}^\top \boldsymbol{\omega} + \textcolor{red}{C} \mathbf{1}^\top \mathbf{1}_{\xi > 1} \right\}, \text{ au lieu de } \min \left\{ \frac{1}{2} \boldsymbol{\omega}^\top \boldsymbol{\omega} \right\}$$

Support Vector Machines, with a Linear Kernel

Pour l'instant, on avait considéré

$$d(\mathbf{x}_0, H_{\boldsymbol{\omega}, b}) = \min_{\mathbf{x} \in H_{\boldsymbol{\omega}, b}} \{ \|\mathbf{x}_0 - \mathbf{x}\|_{\ell_2} \}$$

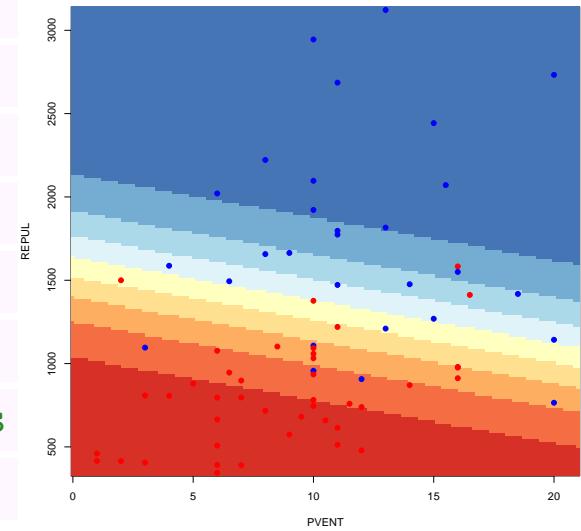
où $\|\cdot\|_{\ell_2}$ est la norme Euclidienne (ℓ_2),

$$\|\mathbf{x}_0 - \mathbf{x}\|_{\ell_2} = \sqrt{(\mathbf{x}_0 - \mathbf{x}) \cdot (\mathbf{x}_0 - \mathbf{x})} = \sqrt{\mathbf{x}_0 \cdot \mathbf{x}_0 - 2\mathbf{x}_0 \cdot \mathbf{x} + \mathbf{x} \cdot \mathbf{x}}$$

```

1 > library(kernlab)
2 > SVM2 <- ksvm(PRONO ~ PVENT + REPUL, data =
  myocarde, prob.model = TRUE, kernel =
  "vanilladot", type="C-svc")
3 > pred_SVM2 = function(p,r){
4 + return(predict(SVM2,newdata=
5 + data.frame(PVENT=p,REPUL=r), type="probabilities
  ")[,2])}

```



Support Vector Machines, with a Non Linear Kernel

Mais plus généralement, on peut considérer

$$d(\mathbf{x}_0, H_{\omega,b}) = \min_{\mathbf{x} \in H_{\omega,b}} \{\|\mathbf{x}_0 - \mathbf{x}\|_k\}$$

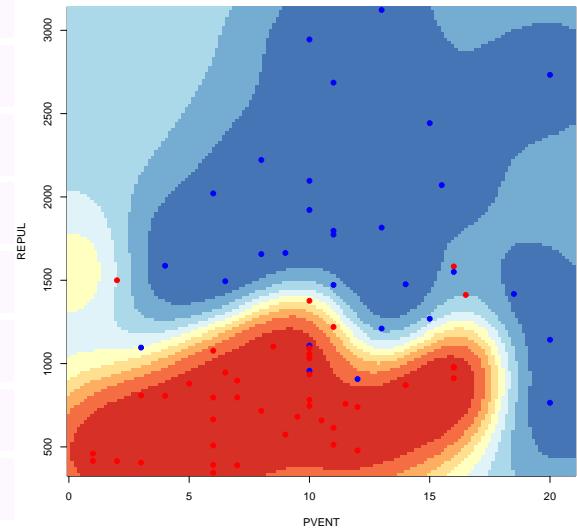
où $\|\cdot\|_k$ est une norme construite avec un noyau

$$\|\mathbf{x}_0 - \mathbf{x}\|_k = \sqrt{k(\mathbf{x}_0, \mathbf{x}_0) - 2k(\mathbf{x}_0, \mathbf{x}) + k(\mathbf{x}, \mathbf{x})}$$

```

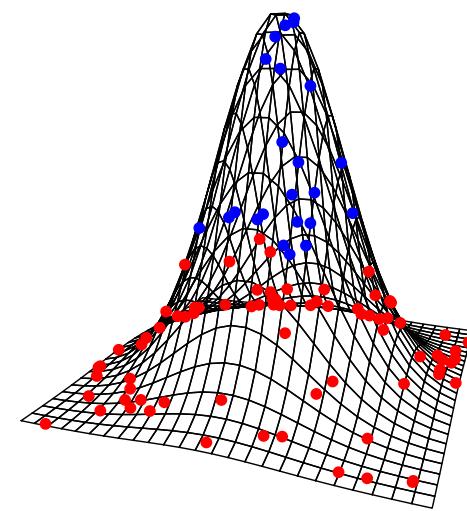
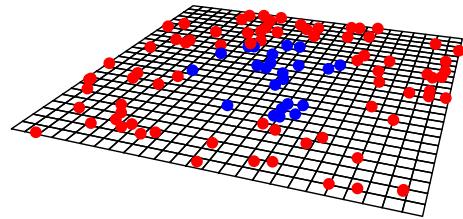
1 > library(kernlab)
2 > SVM2 <- ksvm(PRONO ~ PVENT + REPUL, data =
   myocarde, prob.model = TRUE, kernel = "rbfdot"
   , type="C-svc")
3 > pred_SVM2 = function(p,r){
4 + return(predict(SVM2,newdata=
5 + data.frame(PVENT=p,REPUL=r), type="probabilities
   ")[,2])}

```

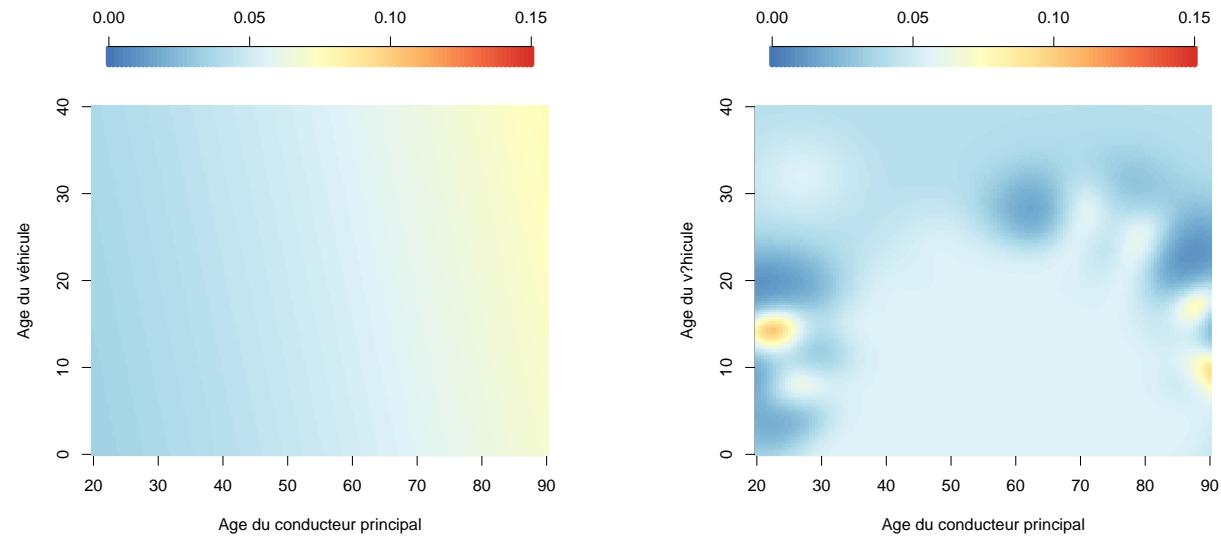


Heuristique sur les SVM à noyau

Une interprétation est que les données ne sont pas linéairement séparables dans l'espace de départ, mais on peut changer d'espace (e.g. noyau Gaussien ici) pour les rendre séparables,



Support Vector Machines



Remarque toutes ces méthodes de type apprentissage statistique sont très sensibles aux tuning parameters.