

## ⑥ Modèle linéaire généralisé

$Y_1, \dots, Y_n$  indépendants

On a vu que le modèle linéaire  $Y_i = x_i^T \beta + \varepsilon_i$  (on suppose les  $x_i$  déterministes  $x_i$ ) permet de prédire singulièrement un label à valeurs réelles  $Y_i \in \mathbb{R}$  à partir de vecteurs de features  $X_i \in \mathbb{R}^d$  ( $\varepsilon_i$  est un bruit centré et  $\text{Var}(\varepsilon_i) = \sigma^2$ )

Dans le modèle linéaire gaussien on avait supposé  $\varepsilon_i \sim \mathcal{N}(0, \sigma^2)$  de sorte que  $Y_i \sim \mathcal{N}(x_i^T \beta, \sigma^2)$ . Dans un modèle linéaire généralisé on va avoir un label  $Y_i$  distribué différemment selon une loi appartenant à la famille de modèles exponentiels, l'espérance  $\mu = \mathbb{E}Y$  dépend de  $x \in \mathbb{R}^d$  à travers une fonction de lien  $g$  tq

$$g(\mathbb{E}Y_i) = x_i^T \beta \quad \text{ou bien} \quad \mu_i = \mathbb{E}Y_i = g^{-1}(x_i^T \beta)$$

On suppose  $g$  monotone et  $C^1$  en général. Dans ce cadre, la variance s'écrit alors également  $\text{Var} Y_i = V(\mu_i) = V(g^{-1}(x_i^T \beta))$  pour une certaine fonction de variance  $V$ .

Def Un modèle linéaire généralisé pour des données  $(x_1, y_1), \dots, (x_n, y_n)$

consiste en

- 1) Un modèle exponentiel pour  $Y \in \mathbb{R}$
- 2) Un prédicteur linéaire  $y_i = x_i^T \beta$  (ie  $\beta \in \mathbb{R}^d$ )
- 3) Une fonction de lien  $g$   $g(\mathbb{E}Y_i) = x_i^T \beta$

1) On considère donc un modèle exponentiel pour  $Y$  : on considère une densité  $f_{\theta}(y) = \exp(b(\theta)^T T(y) - A(\theta)) h(y)$

où  $b, T, A$  et  $h$  sont des fonctions  $\mathbb{R} \rightarrow \mathbb{R}$



C'est une petite généralisation du modèle exponentiel canonique vu  
avant : si on pose  $\tilde{\theta} = b(\theta)$  (reparamétrisation)

on peut écrire  $f_{\tilde{\theta}}(y) = \exp(\tilde{\theta}^T(y) - \log Z(\tilde{\theta}))$

avec  $Z(\tilde{\theta}) = \exp(A(b^{-1}(\tilde{\theta})))$  (si  $b$  inversible...)

et même dominante  $\tilde{\mu}(dy) = h(y) \mu(dy)$

On a vu que dans le modèle canonique où  $b(\theta) = \theta$  et si  $T(y) = y$  alors

$$\mathbb{E}_{\theta} Y = \nabla A'(\theta) \quad \text{et} \quad \text{Var}_{\theta} Y = \nabla^2 A''(\theta) \quad \text{où} \quad A(\theta) = \log Z(\theta)$$

2) L'idée est de lier ce modèle pour  $Y$  avec la valeur du prédicteur linéaire  
en posant simplement  $\eta_i = x_i^T \beta$  (ce prédicteur linéaire  $\in \mathbb{R}$  en général...)

et en supposant que  $\mu_i = \mathbb{E} Y_i = g^{-1}(\eta_i) = g^{-1}(x_i^T \beta)$  pour une certaine  
fonction de lien  $g$ . Le lien entre le prédicteur linéaire ~~est fait~~ ~~avec~~ ~~avec~~ et  
le modèle exponentiel se fait à travers l'espérance de  $Y$  et la fonction de lien  $g$ .

3) La fonction de lien  $g$  fait ce lien, un grand nombre de choix sont possibles.  
Cette fonction doit cependant aller de l'ensemble des valeurs que peut prendre  $\mathbb{E} Y$   
dans  $\mathbb{R}$  ou  $\mathbb{Z}$ . Il y a cependant un choix particulier de lien : la fonction  
de lien canonique définie lorsque l'on regarde le modèle exponentiel sous  
forme canonique.

Déf. Si  $Y \sim f_{\theta}$  avec  $f_{\theta}(y) = \exp(\theta y - A(\theta))$  alors la fonction

de lien canonique est la fonction  $g$  telle que  $\theta = g(\mu)$  i.e.

$$\theta = g(\mathbb{E} Y)$$



Quand on choisit la fonction de lien canonique on a alors

$$\theta = g(\mu) = x^T \beta \quad \text{ie} \quad \theta_i = g(\mu_i) = x_i^T \beta$$

## 6.1) Modèle Bernoulli, régression logistique

Un exemple très important, modèle de régression logistique et certainement un des modèles les plus utilisés au monde.

Label binaire  $Y_i \in \{0, 1\}$  et features  $x_i \in \mathbb{R}^d$ . Par exemple

clique ou non-clique sur une publicité sur internet... Comme  $Y_i \in \{0, 1\}$

facilement  $Y_i \sim \mathcal{B}(\mu_i)$  (Bernoulli de proba  $\mu_i$ )

Il faut faire un lien entre  $\mu_i$  et  $x_i$ : on commence par

utiliser un prédicteur linéaire  $z_i = x_i^T \beta$  (il faudra estimer  $\beta \in \mathbb{R}^d$ )

et choisir une fonction de lien  $g$  telle que  $g(\mathbb{E} Y_i) = g(\mu_i) = x_i^T \beta$

On a que  $\mu_i \in [0, 1]$  et  $x_i^T \beta \in \mathbb{R}$  ( $\beta \in \mathbb{R}^d$  peut prendre des valeurs arbitraires)

donc il faut choisir une fonction  $g: [0, 1] \rightarrow \mathbb{R}$  qui soit monotone

et continue... De telles fonctions sont en fait des fonctions quantiles

$g = F^{-1}$  où  $F: \mathbb{R} \rightarrow [0, 1]$  est la fonction de répartition d'une certaine

loi. Par exemple si on choisit  $F = \Phi =$  fonction de répartition

de  $\mathcal{N}(0, 1)$  on obtient la régression probit.

Quelle est la fonction de lien canonique pour ce modèle? On a ( $y \in \{0, 1\}$ )

$$b_p(y) = p^y (1-p)^{1-y} = \exp\left(y \log\left(\frac{p}{1-p}\right) + \log(1-p)\right)$$

et ici  $\mathbb{E} Y = p$  et  $\theta = \log\left(\frac{p}{1-p}\right)$  pour obtenir un modèle

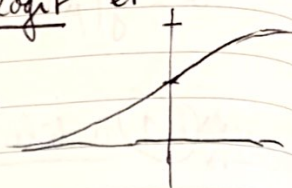


canonique, donc la fonction de lien canonique est ici  $g(\eta) = \log\left(\frac{\eta}{1-\eta}\right)$

car alors  $\theta = g(EY)$ . Cette fonction s'appelle la fonction logit et

$\eta = \frac{1}{1+e^{-\theta}} = g^{-1}(\theta)$  s'appelle la fonction sigmoïde.

(une fonction de répartition bien particulière)



Ce modèle linéaire généralisé s'appelle la régression logistique. On

$$\text{suppose que } P(Y=y|X=x) = \begin{cases} \sigma(\beta^T x) & y=1 \\ 1-\sigma(\beta^T x) & y=0 \end{cases}$$

dans ce modèle. On observe  $(x_1, y_1), \dots, (x_n, y_n)$  et on veut estimer  $\beta \in \mathbb{R}^d$ .

Pour ce faire on utilise à nouveau l'estimateur du maximum de vraisemblance.

On va utiliser plutôt  $Y \in \{-1, 1\}$  à la place de  $Y \in \{0, 1\}$  (sans perte de généralité) car cela nous amènera à des formules plus élégantes.

$$\text{On a } f_{Y|X=x}(y) = \begin{cases} \sigma(\beta^T x) & \text{si } y=1 \\ 1-\sigma(\beta^T x) & \text{si } y=-1 \end{cases}$$

mais  $\sigma$  est la fonction de répartition d'une loi symétrique donc  $\sigma(-z) = 1 - \sigma(z)$

de sorte que  $f_{Y|X=x}(y) = \sigma(y \beta^T x)$  et donc

la vraisemblance de  $(Y_1, \dots, Y_n)$  avec  $Y_i$  indépendantes est donnée par

$$\beta \mapsto \prod_{i=1}^n \sigma(y_i x_i^T \beta) \quad \text{et donc la vraisemblance}$$

$$\text{s'écrit } -\ell(\beta) = -\sum_{i=1}^n \log \sigma(y_i x_i^T \beta)$$

$$= \sum_{i=1}^n \log(1 + e^{-y_i x_i^T \beta})$$



De sorte que  $-\frac{1}{m} \ell(\beta) = \frac{1}{m} \sum_{i=1}^m \log(1 + e^{-y_i x_i^T \beta})$ . Pour estimer  $\beta$  on cherche donc à minimiser la fonction  $-\frac{1}{m} \ell(\beta)$ . Cette fonction est "simple" à minimiser.

Étudions-la un peu: on pose  $f(z) = \log(1 + e^{-z})$

qui s'appelle la fonction de perte logistique. Elle s'interprète bien:

si  $y_i x_i^T \beta \gg 0$  le prédicteur linéaire  $x_i^T \beta$  est très positif et  $y_i = 1$  (resp.  $-1$ )  $\Rightarrow$   $\beta$  permet de bien prédire  $y_i$  et à l'inverse si

$y_i x_i^T \beta \ll 0$  alors le signe de  $x_i^T \beta$  est fortement en désaccord avec celui de  $y_i \Rightarrow \beta$  se trompe fortement sur l'échantillon  $i$ .

En minimisant  $\beta \mapsto -\frac{1}{m} \ell(\beta)$ , on cherche un  $\beta$  qui, en moyenne sur l'échantillon a des petites "pertes" mesurées par la fonction logistique.

C'est un problème d'optimisation convexe: en effet  $f'(z) = -\sigma(-z) (= -\frac{1}{1+e^z})$  et  $f''(z) = \frac{e^z}{(1+e^z)^2} = \sigma(z)\sigma(-z) = \sigma(z)(1-\sigma(z)) \in [0, \frac{1}{4}]$

$$\text{Donc } \nabla_{\beta} \left( -\frac{1}{m} \ell(\beta) \right) = \frac{1}{m} \sum_{i=1}^m \nabla_{\beta} f(y_i x_i^T \beta) = -\frac{1}{m} \sum_{i=1}^m y_i \sigma(-y_i x_i^T \beta) x_i$$

$$\text{et } \nabla_{\beta}^2 \left( -\frac{1}{m} \ell(\beta) \right) = \frac{1}{m} \sum_{i=1}^m \nabla_{\beta}^2 f(y_i x_i^T \beta) = \frac{1}{m} \sum_{i=1}^m \sigma(y_i x_i^T \beta) (1 - \sigma(y_i x_i^T \beta)) x_i x_i^T$$

$$\text{et } 0 \leq \sigma(1-\sigma) \leq \frac{1}{4} \quad \text{et } x_i x_i^T \succ 0 \quad \text{donc } 0 \leq \nabla_{\beta}^2 \left( -\frac{1}{m} \ell(\beta) \right) \left( \leq \frac{1}{4m} \sum x_i x_i^T \right)$$

Hessienne  $\succeq 0 \Rightarrow$  problème convexe.

Rem On peut rendre le problème strictement convexe (1-fortement convexe)

en considérant comme pour la régression linéaire moindres carrés une pénalisation

$$\text{ridge: } \hat{\beta}_{\lambda} \in \underset{\beta \in \mathbb{R}^d}{\text{argmin}} \frac{1}{m} \sum_{i=1}^m \log(1 + e^{-y_i x_i^T \beta}) + \frac{\lambda}{2} \|\beta\|_2^2$$

(solution unique dans ce cas)



D'autres exemples de modèles classiques sont  $\mathcal{P}(\lambda)$  = poisson et  $\mathcal{N}(\mu, \sigma^2)$

loi	Support	lien canonique	Moyenne $\mu$
$\mathcal{N}(\mu, \sigma^2)$	$\mathbb{R}$	$g(\mu) = \mu = x^T \beta$	$\mu = x^T \beta$
$\mathcal{P}(\lambda)$	$\mathbb{N}$	$g(\lambda) = \log \lambda = x^T \beta$	$\mu = e^{x^T \beta}$
$\mathcal{B}(\eta)$	$\{0, 1\}$	$g(\eta) = \log\left(\frac{\eta}{1-\eta}\right)$	$\mu = \frac{1}{1 + e^{-x^T \beta}}$
$\mathcal{E}(\lambda)$	$\mathbb{R}^+$	$g(\lambda) = -\frac{1}{\lambda}$	$\mu = -\frac{1}{x^T \beta}$

## 6.2 Maximum de vraisemblance dans un GLM et TCL

On veut établir un TCL dans le modèle linéaire généralisé avec lien

canonique, i.e.  $Y \sim f_{\beta}$  où  $f_{\beta}(y|x) = \exp(y x^T \beta - A(x^T \beta)) h(y)$

Rappel : pour la régression logistique  $A(z) = -\log(1 + e^z)$

$$A'(z) = -\sigma(-z) \quad A''(z) = -\sigma(z)(1 - \sigma(z))$$

La log vraisemblance s'écrit

$$\ell_n(\beta) = \sum_{i=1}^n \left\{ y_i x_i^T \beta - A(x_i^T \beta) - \log h(y_i) \right\}$$

et le score vaut

$$S_n(\beta) = \nabla \ell_n(\beta) = \sum_{i=1}^n \left\{ y_i - A'(x_i^T \beta) \right\} x_i$$

et l'information de Fisher

$$I_n(\beta) = \nabla^2 \ell_n(\beta) = - \sum_{i=1}^n A''(x_i^T \beta) x_i x_i^T$$

On rappelle que  $\mathbb{E}_{\beta} [S_n(\beta)] = 0$  et  $\text{cov}_{\beta} (S_n(\beta)) = I_n(\beta)$

On va avoir besoin de qq hypothèses qu'on va discuter rapidement:

On commence par poser  $B_n(\delta) = \left\{ \beta \in \mathbb{R}^d : \left\| I_n(\beta_0)^{1/2} (\beta - \beta_0) \right\| \leq \delta \right\}$



Hypothèses On suppose  $Y_1, \dots, Y_n$  i.i.d de densité  $f_{\beta_0}$  pour  $\beta_0 \in \mathbb{R}^d$   
 = le "vrai" paramètre.

Hyp (I)  $\lambda_{\min}(I_n(\beta_0)) \rightarrow +\infty \quad n \rightarrow +\infty$

Hyp (II) On pose  $\tilde{I}_n(\beta) = I_n(\beta_0)^{-1/2} I_n(\beta) I_n(\beta_0)^{-1/2}$

matrice d'information de Fisher "normalisée" on suppose que

$$\forall \delta > 0 \quad \max_{\beta \in B_n(\delta)} \|\tilde{I}_n(\beta) - Id\| \rightarrow 0$$

On a que si  $\beta \in B_n(\delta)$  alors  $\|\beta - \beta_0\| \leq \frac{\delta}{\lambda_{\min}(I_n(\beta_0))} \rightarrow 0$

donc  $B_n(\delta)$  se réduit à  $\{\beta_0\}$  qd  $n \rightarrow +\infty \quad (\forall \delta > 0)$

Donc sous (I)/(II) est une hypothèse asymptotiquement faible. De plus (II) équivaut

à  $\forall \varepsilon > 0 \quad \forall u \in \mathbb{R}^d$  et  $\forall \beta \in B_n(\delta)$  on a

$$|u^T (\tilde{I}_n(\beta) - Id) u| \leq \varepsilon \|u\|^2 \quad \text{pour } n \text{ assez grand}$$

ie en posant  $v = I_n^{-1/2}(\beta_0) u$  on a

$$\forall \varepsilon > 0 \quad \forall v \in \mathbb{R}^d \quad \forall \beta \in B_n(\delta)$$

$$|v^T I_n(\beta) v - v^T I_n(\beta_0) v| \leq \varepsilon v^T I_n(\beta_0) v \quad n \text{ assez grand}$$

ce qui entraîne ( $\varepsilon = 1/2$ )  $I_n(\beta) \not\prec c I(\beta_0) \quad \forall \beta \in B_n(\delta)$  pour  $n$  assez grand

Thm Sous (I)+(II) on a pour  $\hat{\beta}_n$  MLE (existe avec proba  $\rightarrow 1$ )

$$\left[ \text{que} \quad I_n(\beta_0)^{1/2} (\hat{\beta}_n - \beta_0) \rightsquigarrow N(0, Id) \right]$$

Dém Pour l'instant prenons  $\hat{\beta}_n$  tq  $S_n(\hat{\beta}_n) = 0$  (MLE si il existe)  
 et faisons un DL  $\left( \nabla \ell_n(\hat{\beta}_n) \right)$



$$S_n(\beta_0) = S_n(\hat{\beta}_n) + \int_0^1 I_n(\beta_0 + t(\hat{\beta}_n - \beta_0)) dt (\hat{\beta}_n - \beta_0)$$

$$= I_n(\beta_0)^{1/2} \int_0^1 \tilde{I}_n(\beta_0 + t(\hat{\beta}_n - \beta_0)) dt I_n(\beta_0)^{1/2} (\hat{\beta}_n - \beta_0)$$

ie  $I_n(\beta_0)^{-1/2} S_n(\beta_0) = \int_0^1 \tilde{I}_n(\text{---}) dt I_n(\beta_0)^{1/2} (\hat{\beta}_n - \beta_0)$

Donc si (A)  $\int_0^1 \tilde{I}_n(\beta_0 + t(\hat{\beta}_n - \beta_0)) dt \xrightarrow{P} I_d$  et si  
 (B)  $I_n(\beta_0)^{-1/2} S_n(\beta_0) \rightsquigarrow N(0, I_d)$  (C)  $\hat{\beta}_n$  converge avec proba  $\rightarrow 1$

alors on a  $I_n(\beta_0)^{1/2} (\hat{\beta}_n - \beta_0) \rightsquigarrow N(0, I_d)$ .

Donc il faut montrer (A) + (B) + (C).

Preuve de (B)

Soit  $\delta > 0$  et  $u \in \mathbb{R}^d$ ,  $\|u\|=1$ . On pose  $\beta_n = \beta_0 + \delta I_n(\beta_0)^{-1/2} u$  de sorte que  $\|\beta_n - \beta_0\|_{I_n(\beta_0)} = \delta$  ie  $\beta_n \in \partial B_n(\delta) \subset B_n(\delta)$ .

On fait un DL = on a pour un certain  $\tilde{\beta}_n \in [\beta_n, \beta_0] = \{t\beta_n + (1-t)\beta_0 : t \in [0,1]\}$

$$\ln(\beta_n) = \ln(\beta_0) + \langle \beta_n - \beta_0, \nabla \ln(\beta_0) \rangle + \frac{1}{2} (\beta_n - \beta_0)^T \nabla^2 \ln(\tilde{\beta}_n) (\beta_n - \beta_0)$$

$$= \ln(\beta_0) + \delta \langle u, I_n(\beta_0)^{-1/2} S_n(\beta_0) \rangle - \frac{\delta^2}{2} u^T \tilde{I}_n(\tilde{\beta}_n) u$$

et enfin  $e^{\frac{\delta^2}{2} u^T \tilde{I}_n(\tilde{\beta}_n) u} L_n(\beta_n) = e^{\delta u^T I_n(\beta_0)^{-1/2} S_n(\beta_0)} L_n(\beta_0)$

où  $L_n(\beta) = e^{\ln(\beta)}$  = vraisemblance du modèle en  $\beta$ . Donc en intégrant on obtient

$$\mathbb{E}_{\beta_n} \left[ \exp\left(\frac{\delta^2}{2} u^T \tilde{I}_n(\tilde{\beta}_n) u\right) \right] = \mathbb{E}_{\beta_0} \left[ \exp\left(\delta u^T I_n(\beta_0)^{-1/2} S_n(\beta_0)\right) \right]$$



On a  $\beta_n \in B_n(\delta)$  donc (II) entraîne  $\max \|\tilde{I}_n(\tilde{\beta}_n) - I_d\| \rightarrow 0$

donc  $\tilde{\beta}_n \in B_n(\delta)$

et par continuité de l'exponentielle on obtient  $\forall \varepsilon > 0 \left| e^{\frac{\delta^2}{2} u^T \tilde{I}_n(\tilde{\beta}_n) u} - e^{\frac{\delta^2}{2}} \right| \leq \varepsilon$

pour  $n$  assez grand et donc  $E_{\beta_n} \exp\left(\frac{\delta^2}{2} u^T \tilde{I}_n(\tilde{\beta}_n) u\right) \rightarrow e^{\frac{\delta^2}{2}}$

et donc  $E_{\beta_0} \left[ \exp\left(\delta u^T I_n(\beta_0)^{-1/2} S_n(\beta_0)\right) \right] \rightarrow e^{\frac{\delta^2}{2}} \quad n \rightarrow +\infty$

ce qui entraîne que  $I_n(\beta_0)^{-1/2} S_n(\beta_0) \rightsquigarrow N(0, I_d)$

(la transformée de Laplace converge vers la transformée de Laplace de  $N(0, I_d)$ )

Preuve de (A) On veut montrer que

$$\left\| \int_0^1 \tilde{I}_n(\beta_0 + t(\hat{\beta}_n - \beta_0)) dt - I_d \right\| \xrightarrow{P} 0$$

On remarque que  $\beta_0 + t(\hat{\beta}_n - \beta_0) \in B_n(\delta)$  si  $\hat{\beta}_n \in B_n(\delta)$

Il faut que l'on étudie  $\hat{\beta}_n$  comme un maximiseur quand il existe. On rappelle

$$\text{que } I_n(\beta) = -\nabla^2 \ln(\beta) = -\sum_{i=1}^m A''(x_i; \beta) x_i x_i^T < 0 \quad \forall \beta \in \mathbb{R}^d$$

donc  $\ln(\beta)$  est strictement concave sur  $\mathbb{R}^d \Rightarrow$  il y a au plus

un  $\beta \in \mathbb{R}^d$  tq  $S_n(\beta) = 0$  et c'est dans ce cas  $\hat{\beta}_n =$  le MLE.

Si  $\ln(\beta) - \ln(\beta_0) < 0 \quad \forall \beta \in \partial B_n(\delta)$  (pour un certain  $\delta > 0$ )

alors il y a un maximiseur dans  $B_n(\delta)$  (car  $\ln$  bornée au bord de  $B_n(\delta)$ )

et dans ce cas c'est  $\hat{\beta}_n$ . On a montré que

$$\forall \eta > 0 \exists \delta > 0 \text{ tq } P(\ln(\beta) - \ln(\beta_0) < 0 \quad \forall \beta \in \partial B_n(\delta)) \geq 1 - \eta$$

Sur  $\beta \in \partial B_n(\delta)$  de sorte que  $u = \frac{1}{\delta} I_n(\beta_0)^{-1/2} (\beta - \beta_0)$  vérifie  $\|u\| = 1$

et on refait encore le DL

$$\ln(\beta) - \ln(\beta_0) = \delta \langle u, I_n(\beta_0)^{-1/2} S_n(\beta_0) \rangle - \frac{\delta^2}{2} u^T \tilde{I}_n(\tilde{\beta}_n) u \quad \text{pour } \tilde{\beta}_n \in [\beta_0, \beta_n]$$



Rem:  $\max_{\|v\|=1} \langle v, I_n(\beta_0)^{-1/2} S_n(\beta_0) \rangle = \|I_n(\beta_0)^{-1/2} S_n(\beta_0)\|$  et

$\inf_{\|v\|=1} v^T I_n(\tilde{\beta}_n) v = \lambda_{\min}(\tilde{I}_n(\tilde{\beta}_n))$  donc il suffit de m<sub>1</sub>

$$P\left[\|I_n(\beta_0)^{-1/2} S_n(\beta_0)\|^2 < \frac{\delta^2}{4} \lambda_{\min}(\tilde{I}_n(\tilde{\beta}_n))^2\right] \geq 1 - \gamma$$

mais (II)  $\Rightarrow I_n(\beta) \succeq c I_n(\beta_0) \quad \forall \beta \in B_n(\delta)$

donc  $I_n(\beta_0)^{-1/2} I_n(\beta) I_n(\beta_0)^{-1/2} \succeq c I_d$   
 $\tilde{I}_n(\beta)$

donc il suffit de vérifier que  $P\left[\|I_n(\beta_0)^{-1/2} S_n(\beta_0)\|^2 < \frac{(\delta c)^2}{4}\right] \geq 1 - \gamma$

Mais il suffit d'appliquer Markov:  $P\left[\| \cdot \|^2 < \frac{(\delta c)^2}{4}\right] \geq 1 - \frac{4 E \|I_n(\beta_0)^{-1/2} S_n(\beta_0)\|^2}{(\delta c)^2}$

et enfin  $E \|I_n(\beta_0)^{-1/2} S_n(\beta_0)\|^2$

$$= E \left[ S_n(\beta_0)^T I_n(\beta_0)^{-1} S_n(\beta_0) \right] = E \text{tr} (S^T I^{-1} S)$$

$$= E \text{tr} (I^{-1} S S^T) = \text{tr} (I^{-1} E(S S^T))$$

$$= \text{tr} \left( I_n^{-1}(\beta_0) \underbrace{E[S_n(\beta_0) S_n(\beta_0)^T]}_{\parallel} \right)$$

$$\text{Var}_{\beta_0}(S_n(\beta_0)) = I_n(\beta_0)$$

$$= d \quad \text{donc } n \text{ on choisit } \delta = \sqrt{\frac{4d}{c^2 \gamma}} \quad \text{et pour}$$

n assez grand (tq  $I_n(\beta) \succeq I_n(\beta_0) \quad \forall \beta \in B_n(\delta)$ ) on obtient bien

$$P\left[\ln(\beta) - \ln(\beta_0) < 0 \quad \forall \beta \in \partial B_n(\delta)\right] \geq 1 - \gamma$$

et on remarque que sur cet événement  $\hat{\beta}_n$  existe et  $\hat{\beta}_n \in B_n(\delta)$



donc  $P[\hat{\beta}_n \in B_n(\delta)] \geq 1 - \eta$

$$P[\|I_n(\beta_0)^{-1/2}(\hat{\beta}_n - \beta_0)\| \leq \delta] \geq 1 - \eta.$$

On a aussi  $\beta_0 + t(\hat{\beta}_n - \beta_0) \in B_n(\delta)$  si  $\hat{\beta}_n \in B_n(\delta)$  donc en utilisant II on obtient  $\|\tilde{I}_n(\beta_0 + t(\hat{\beta}_n - \beta_0)) - I_d\| \rightarrow 0$

et donc  $\|\int_0^1 \tilde{I}_n(\beta_0 + t(\hat{\beta}_n - \beta_0)) dt - I_d\| \rightarrow 0$  sur cet événement on a donc montré que  $\forall \varepsilon > 0$  et  $\forall \eta > 0$  on peut trouver  $\delta$  et  $n$  assez grand

$$P\left[\left\|\int_0^1 \tilde{I}_n(\text{---}) dt - I_d\right\| > \varepsilon\right] \leq \eta$$

ce qui montre (A) □

Concernant les hypothèses : si  $C_1 \in A''(z) \subseteq C_2$  pour  $z \in \textcircled{4}$   $0 < C_1 < C_2$

et  $(x_i)_{i=1, \dots, n}$  et  $p$  tq  $x_i^T p \in \textcircled{4}$  ( $p \in B$  compact)

alors 
$$I_n(\beta_0) = \sum_{i=1}^n A''(x_i^T \beta_0) x_i x_i^T$$

$$\geq \min_{i=1, \dots, n} A''(x_i^T \beta_0) \sum_{i=1}^n x_i x_i^T$$

si  $\lambda_{\min}\left(\sum_{i=1}^n x_i x_i^T\right) \rightarrow +\infty$  alors  $\lambda_{\min}(I_n(\beta_0)) \rightarrow +\infty$  et (II) est

vérifiée si  $1 - \varepsilon \leq \frac{A''(x_i^T \beta)}{A''(x_i^T \beta_0)} \leq 1 + \varepsilon \quad \forall i$

mais  $A \in C^\infty$  et (convexe) sur  $\textcircled{4}$

$E_x$  log-veg 
$$\frac{\sigma(x_i^T \beta)(1 - \sigma(x_i^T \beta))}{\sigma(x_i^T \beta_0)(1 - \sigma(x_i^T \beta_0))}$$