# Cours n°3 – Régression avancée L. Carraro

Option MAFQ



### Plan

- Bootstrap et bagging
- Retour sur les prédicteurs corrélés
- Régression ridge
- Les degrés de liberté
- LASSO
- Vers les méthodes L<sup>1</sup> avancées



# Le bootstrap

- Méthode générale en statistiques
- Exemple élémentaire : intervalle de confiance
  - Soit  $(x_1,...,x_n)$  un n-échantillon de loi normale N(m,1)
  - Problème : intervalle de confiance pour m basé sur l'échantillon
- Ici, calcul théorique :

$$\overline{X} = \frac{1}{n} \sum_{i=1}^{n} X_i \text{ est de loi } N(m, 1/n)$$



# Cas général

• Un intervalle de confiance de niveau α doit vérifier :

$$P_{\theta_0}(\theta_0 \in [l(X_1,...,X_n),u(X_1,...,X_n)]) = 1 - \alpha$$

- Pour l'exemple qui précède, on trouve un intervalle valable pour tout  $m=\theta$  et non seulement pour le « bon »  $\theta_0$
- Bootstrap:

Remplacer  $P_{\theta_0}$  par une estimation notée  $P_n^*$ 

• choix le plus courant : loi empirique

$$P_n^* = \frac{1}{n} \sum_{i=1}^n \delta_{X_i}$$



# Algorithme

- Pour b=1 à B, faire :
  - simuler un échantillon de taille n parmi  $\{x_1,...,x_n\}$ , noté  $X^{*b}$
  - stocker l'estimation de  $\theta$  obtenue, notée  $\theta^{*b}$
- Fin faire
- Evaluer les quantiles d'ordre  $\alpha/2$  et 1- $\alpha/2$  de l'échantillon des  $\theta^{*b}$ , b=1,...,B

$$IC_{bootstrap} = [q_{\alpha/2}, q_{1-\alpha/2}]$$

Résultats : voir simulations R



# Le bagging

Bagging = Boostraping and Aggregating

### • Idée:

- rééchantillonner parmi les (x<sub>i</sub>, y<sub>i</sub>)
- construire un modèle sur cet échantillon
- moyenner les modèles obtenus



# Algorithme pour le bagging

- Soit Z = (X,Y) ou (X,G) l'ensemble d'apprentissage
  - Pour b=1 à B, faire :
    - simuler un échantillon de taille n parmi  $\{z_1,...,z_n\}$ , noté  $Z^{*b}$
    - estimer un modèle à partir de Z\*b
    - stocker le modèle obtenu, noté  $\hat{Y}_b^*$
  - Fin faire
  - Poser :

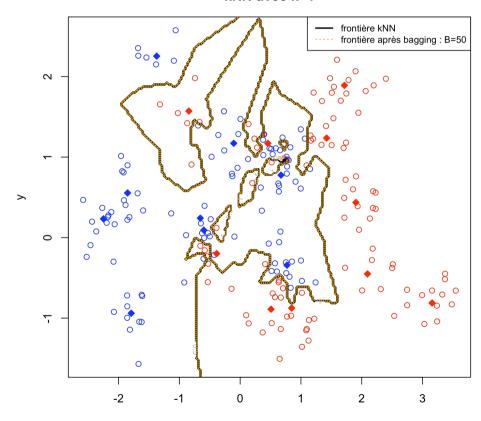
$$\hat{Y}_{bag}(x) = \frac{1}{B} \sum_{b=1}^{B} \hat{Y}_{b}^{*}(x)$$



# Warning

- Inutile si régression linéaire
- Inutile si modèle linéaire en les réponses

#### séparation des deux classes kNN avec k=1





# Retour sur prédicteurs corrélés

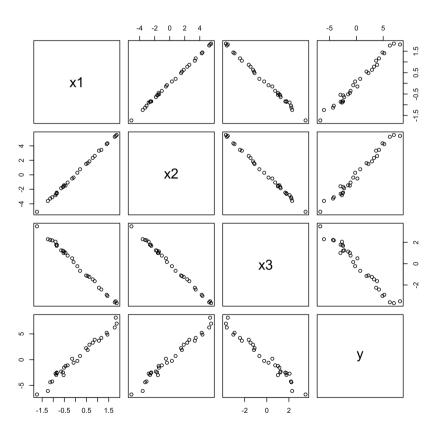
- Exemple synthétique
- On suppose que y dépend linéairement de la moyenne de 3 prédicteurs x<sub>1</sub>, x<sub>2</sub>, x<sub>3</sub> :

$$y = 2 \times (x_1 + x_2 + x_3) + \varepsilon$$

- $-\varepsilon$  de loi N(0,0.1)
- Plan d'expériences à 30 points :
  - $-x_1$  de loi N(0,1),  $x_2/x_1$  de loi N(3 $x_1$ ,0.1)
  - $-x_3/x_1,x_2$  de loi  $N(x_1-x_2,0.1)$



# pairplot





## Modèle linéaire

#### Coefficients:

```
Estimate Std. Error t value Pr(>|t|)
(Intercept) 0.02610 0.07898 0.330 0.743746
x1 8.00830 3.10282 2.581 0.015847 *
x2 1.25169 1.35297 0.925 0.363401
x3 3.89028 0.92633 4.200 0.000277 ***
```



## Analyse de modèles emboités

```
Analysis of Variance Table
> mod1 < -lm(y \sim x1 + x2 + x3)
> mod2 < -lm(v \sim x3 - 1)
> anova(mod2,mod1)
                                     Model 1: y \sim x3-1
                                     Model 2: y \sim x1 + x2 + x3
     Res.Df RSS Df Sum of Sq F Pr(>F)
          29 14,9925
          26 4.6244 3 10.368 19.431 8.084e-07 ***
      Remarque : si conserve le modèle 2 :
    Coefficients:
       Estimate Std. Error t value Pr(>|t|)
    x3 -1.86923 0.06418 -29.12 <2e-16 ***
    Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
    Residual standard error: 0.719 on 29 degrees of freedom
    Multiple R-squared: 0.9669, Adjusted R-squared: 0.9658
    F-statistic: 848.1 on 1 and 29 DF, p-value: < 2.2e-16
```

## Seconde simulation

#### Coefficients:

Rien de significatif avec  $R^2 = 97\%$ !!



## Troisième simulation

#### Coefficients:

	Estimate	Std. Error	t value	Pr(> t )	
(Intercept)	0.11513	0.09604	1.199	0.24141	
x1	0.89016	2.84315	0.313	0.75671	
x2	2.96368	1.25139	2.368	0.02559	*
<b>x</b> 3	2.88135	1.03635	2.780	0.00996	**



# Troisième simulation modèles emboités

```
Analysis of Variance Table
> mod1 < -lm(y \sim x1 + x2 + x3)
> mod2 < -lm(y \sim x2 + x3 - 1)
> anova(mod2,mod1)
                                     Model 1: y \sim x2+x3-1
                                     Model 2: y \sim x1 + x2 + x3
     Res.Df RSS Df Sum of Sq F Pr(>F)
          28 6.8955
         26 6.5343 2 0.3612 0.7186 0.4969
      Ici, on conserve le modèle 2 :
    Coefficients:
       Estimate Std. Error t value Pr(>|t|)
                    0.6725 4.735 5.72e-05 ***
         3.1844
    x2
    x3 2.7640
                   1.0057 2.748 0.0104 *
    Residual standard error: 0.4963 on 28 degrees of freedom
    Multiple R-squared: 0.9852, Adjusted R-squared: 0.9842
    F-statistic: 933.8 on 2 and 28 DF, p-value: < 2.2e-16
```

# Qualité de prédiction dernier modèle

- Plan test factoriel à 1331 points
  - $-x_1$  entre -2 et 2
  - $-x_2$  entre -6 et 6
  - $-x_3$  entre -4 et 4
- Ecart type des écarts prédit simulé :
  - 13.96 (au lieu de 0.5!)
- Pourcentage de réponses dans l'intervalle de prévision à 95%
  - 11,8 %



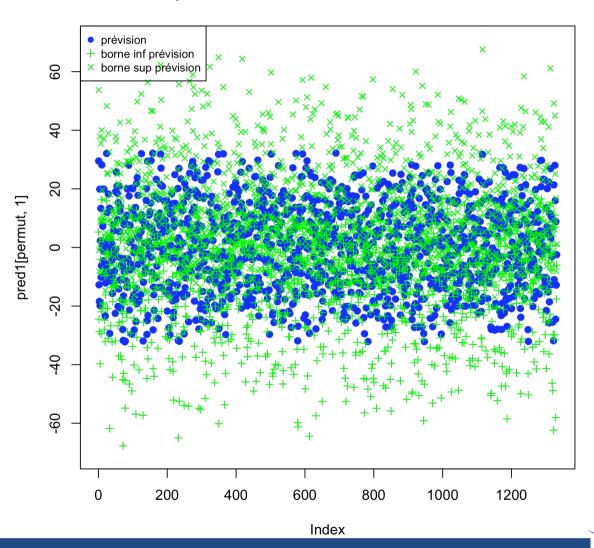
# Qualité de prédiction

- Troisième simulation
- modèle complet
- Ecart type des écarts prédit simulé :
  - 13.96 (au lieu de 0.5!)
- Pourcentage de réponses dans l'intervalle de prévision à 95%
  - 53 %
  - La colinéarité des prédicteurs augmente la taille de l'intervalle de prévision



## Illustration

#### prévisions et intervalles de confiance





## Conclusion

- Corrélation des prédicteurs entraîne :
  - des choix de modèles erronés
  - des modèles souvent structurellement aberrants (cf. signe des coefficients)
  - des modèles instables
  - des modèles strictement valables dans le domaine d'apprentissage



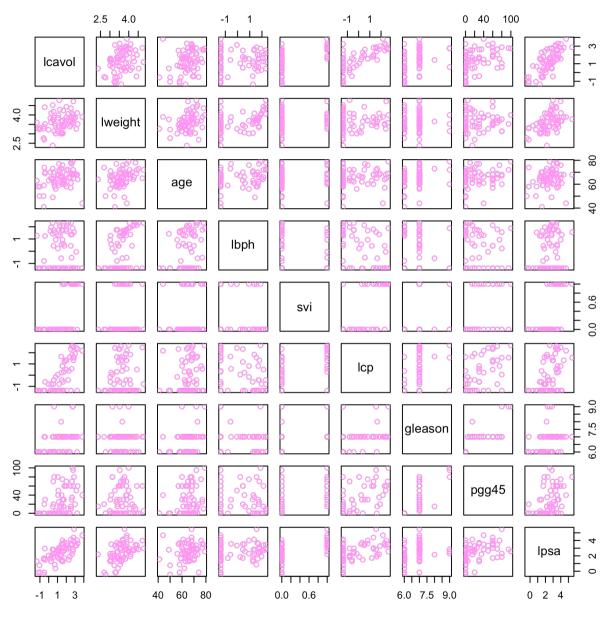
# Exemple « prostate »

• Réponse : taux de PSA (prostate-specific antigen) : lpsa

#### • Prédicteurs :

- lcavol (log volume cancer)
- lweight (log poids prostate)
- age (en années)
- lbph (log benign prostatic hyperplasia)
- svi (invasion vésicule séminale)
- lcp (log pénétration capsulaire)
- gleason (indice : 6, 7, 8, 9)
- pgg45(pourcentage de Gleason 4 ou 5)

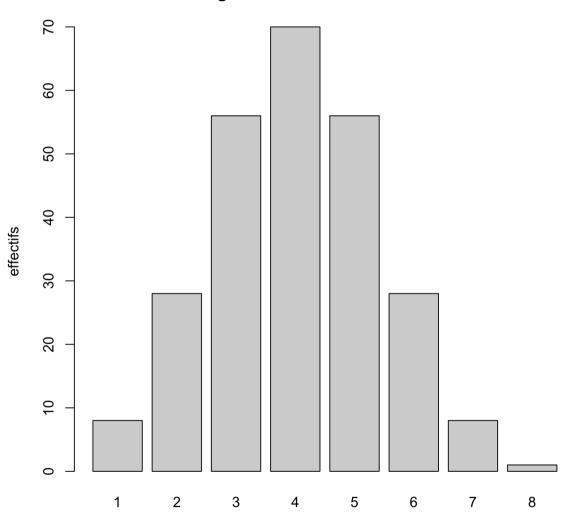




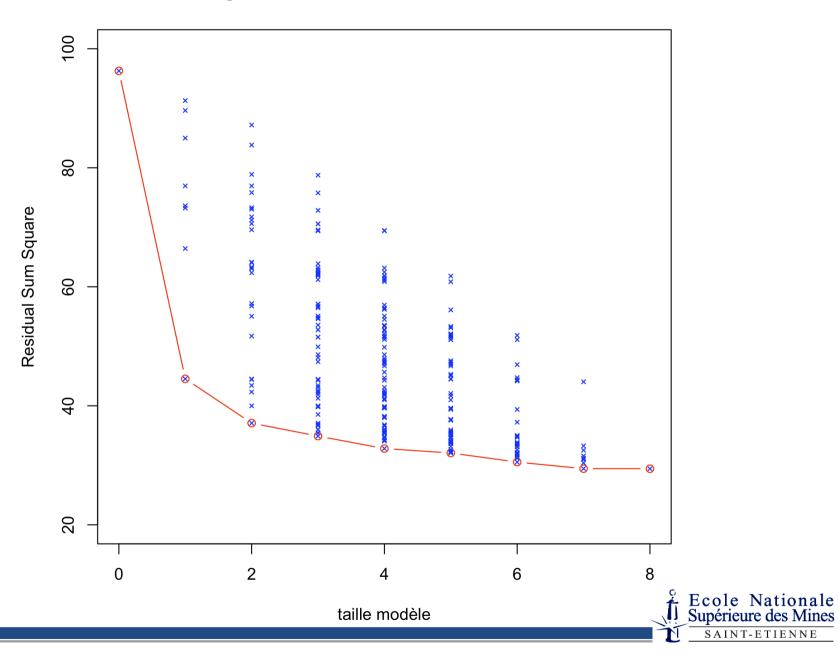


# On essaye tous les sous-modèles

#### histogramme des tailles de modèles



#### régression sur tous les sous-modèles



# Pourquoi un sous-modèle?

### Prévision

- intervalles de prévision corrects
- intervalles de prévision de taille raisonnable
- Interprétation
  - on veut la « big picture »
  - détails sans importance



# Régression robuste « ridge regression »

• On part d'un modèle linéaire :

$$Y = X\beta + \varepsilon$$

on suppose les colonnes de X normées

• Estimation de β

$$Arg\min_{\beta} \left\{ \left\| Y - X\beta \right\|^{2} + \lambda \left\| \beta \right\|^{2} \right\}$$

OU

$$\underset{\|\beta\| \leq c}{Arg\min} \left\{ \left\| Y - X\beta \right\|^{2} \right\}$$

# Ridge: estimation de $\beta$

$$\hat{\beta}_{ridge} = (X'X + \lambda Id)^{-1}X'Y$$

- Augmentation du conditionnement matriciel
- Remarque : si les colonnes de X sont orthogonales

$$\hat{\beta}_{ridge} = \frac{1}{1+\lambda} X'Y = \frac{1}{1+\lambda} \hat{\beta}$$



# Vision bayésienne

- On suppose B de loi a priori  $N(0,\eta^2 \text{ Id})$
- Le modèle de régression donne :
  - Y/B= $\beta$  de loi N(X $\beta$ , $\sigma$ <sup>2</sup> Id)
- Loi (a posteriori) de B:
  - $-f_{B/Y=y}(\beta) = f_{Y/B=\beta}(y) f_B(\beta)/f_Y(y)$
  - Espérance et mode (a posteriori) de Y sont égaux à l'estimateur ridge, avec  $\lambda = \sigma^2/\eta^2$



## Evaluations matricielles

• Décomposition en valeurs singulières de X :

$$-X = UDV'$$

- U n\*p orthogonale : U'U=Id<sub>p</sub>
- D p\*p diagonale
- V p\*p orthogonale : V'V=Id<sub>p</sub>

$$\hat{Y} = X\hat{\beta} = X(X'X)^{-1}X'Y = UU'Y$$

NB: U'Y donne les coordonnées de Y dans la base des colonnes de U



# Cas de l'estimateur ridge

$$\hat{Y}_{ridge} = X\hat{\beta}_{ridge} = X(X'X + \lambda Id)^{-1}X'Y$$

$$= UD(D^2 + \lambda Id)^{-1}DU'Y$$

$$= \sum_{j=1}^{p} u_j \frac{d_j^2}{d_j^2 + \lambda} u_j Y$$

Décomposition de Y sur la base des colonnes de U, puis diminution des coefficients gouvernée par λ



# Les degrés de liberté

• En régression linéaire :

```
    p = dim(ev{colonnes de X})
    = dim(Im(H))
    où H est la matrice chapeau : H=X(X'X)-1X'
    H est un projecteur : dim(Im(H))=Tr(H)
```

- En régression ridge :  $H_{\lambda}=X(X'X+\lambda Id)^{-1}X'$  n'est pas un projecteur
- Définition très générale (ex. splines)



# Degrés de liberté en régression ridge

• On pose:

$$df(\lambda) = Tr(H_{\lambda})$$

• Propriétés :

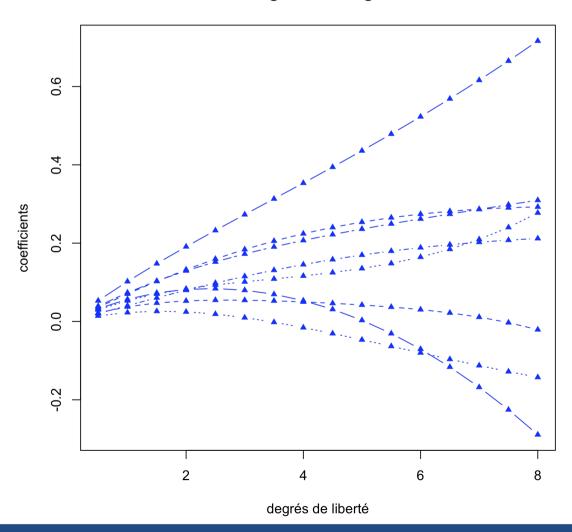
$$Tr(H_{\lambda}) = \sum_{j=1}^{p} \frac{d_{j}^{2}}{d_{j}^{2} + \lambda}$$

- df(0)=p, df(.) tend vers 0 à l'infini
- df(.) est décroissante



## Evolution des coefficients

#### régression ridge





## Le LASSO

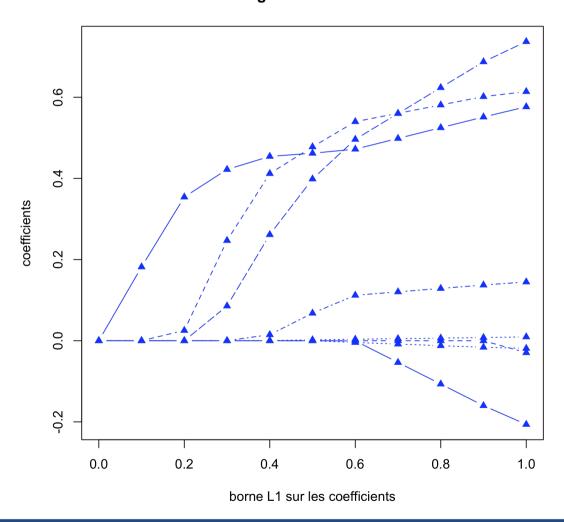
• Même idée que ridge :

$$\underset{\|\beta\|_{1} \leq c}{Arg\min} \left\{ \|Y - X\beta\|^{2} \right\}$$

- Mais:
  - on remplace  $\|\beta\|_2$  par  $\|\beta\|_1$

## Evolution des coefficients

#### régression LASSO





## Différences fondamentales

- Estimation de  $\beta$  par LASSO via un algorithme de programmation quadratique
- Estimation de β par ridge par inversion matricielle
- LASSO donne des coefficients nuls

