

Modèles Linéaires Appliqués

Arthur Charpentier

Automne 2020

OLS #16 (multicolinéarité, Ridge)

Régression Linéaire sans \mathcal{H}_1

- ▶ \mathbf{X} est mal conditionnée, i.e. certaines valeurs propres de $\mathbf{X}^\top \mathbf{X}$ sont proches de 0 $\Rightarrow \det(\mathbf{X}^\top \mathbf{X}) \simeq 0$, phénomène de colinéarité.
- ▶ (ou/et) $p > n$ le nombre de covariables est nettement supérieur au nombre d'observations.

L'hypothèse initiale \mathcal{H}_1 est donc largement remise en cause!

Régression Ridge

- ▶ Idée: perturber la matrice \mathbf{X} pour éloigner ses valeurs propres de 0 et ainsi stabiliser l'inversion de $\mathbf{X}^\top \mathbf{X}$.
- ▶ Notation: $0 \leq \mu_1 \leq \dots \mu_p$ valeurs propres ordonnées de $\mathbf{X}^\top \mathbf{X}$ et soit \mathbf{P} la matrice orthogonale telle que $\mathbf{X}^\top \mathbf{X} = \mathbf{P} \mathbf{D} \mathbf{P}^\top$ avec $\mathbf{D} = \text{diag}(\mu_1, \dots, \mu_p)$.
- ▶ Soit $\lambda \geq 0$, $\mathbf{X}^\top \mathbf{X} + \lambda \mathbb{I}_p$ a les mêmes vecteurs propres que $\mathbf{X}^\top \mathbf{X}$ et pour valeurs propres $\mu_j + \lambda$ pour $j = 1, \dots, p$.
- ▶ Hoerl and Kennard (1970): remplacer $\mathbf{X}^\top \mathbf{X}$ par $\mathbf{X}^\top \mathbf{X} + \lambda \mathbb{I}_p$ dans la définition des MCO:

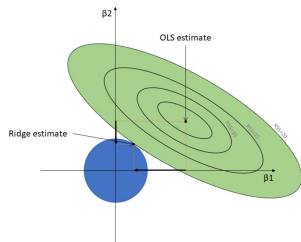
$$\hat{\beta}_{\text{ridge}}(\lambda) = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbb{I}_p)^{-1} \mathbf{X}^\top \mathbf{Y}$$

- ▶ Si $\lambda = 0$, alors $\hat{\beta}_{\text{ridge}} = \hat{\beta}^{\text{mco}}$; si $\lambda \rightarrow \infty$, alors $\hat{\beta}_{\text{ridge}} \rightarrow 0$.
- ▶ Fixer λ est important. Quelles sont les conséquences sur le biais et variance des estimateurs?

Régression Ridge

- Soit $\lambda \geq 0$, on peut montrer que $\hat{\beta}_{\text{ridge}}(\lambda)$ minimise

$$\begin{aligned} & \|\mathbf{Y} - \mathbf{X}\beta\|^2 + \lambda\|\beta\|^2 \\ &= \sum_{i=1}^n (Y_i - \mathbf{x}_i\beta)^2 + \lambda \sum_{j=1}^p \beta_j^2 \end{aligned}$$



<https://rstatisticsblog.com/>

- Soit $\tilde{\beta}$ l'estimateur minimisant

$$\|\mathbf{Y} - \mathbf{X}\beta\|^2 \text{ sous la contrainte que } \|\beta\|^2 \leq \delta.$$

$\forall \lambda > 0, \exists \delta$ tel que les deux solutions coïncident.

Régression Ridge

- ▶ $\hat{\beta}_{\text{ridge}} = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbb{I}_p)^{-1} \mathbf{X}^\top \mathbf{X} \hat{\beta}.$
- ▶ $\mathbb{E}(\hat{\beta}_{\text{ridge}}) = \beta - \lambda (\mathbf{X}^\top \mathbf{X} + \lambda \mathbb{I}_p)^{-1} \beta.$
- ▶ $\text{Var}(\hat{\beta}_{\text{ridge}}) = \sigma^2 (\mathbf{X}^\top \mathbf{X} + \lambda \mathbb{I}_p)^{-1} \mathbf{X}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X} + \lambda \mathbb{I}_p)^{-1}.$
- ▶ $\text{EQM}(\hat{\beta}_{\text{ridge}}) (\mathbf{X}^\top \mathbf{X} + \lambda \mathbb{I}_p)^{-1} (\sigma^2 (\mathbf{X}^\top \mathbf{X}) + \lambda^2 \beta \beta^\top) (\mathbf{X}^\top \mathbf{X} + \lambda \mathbb{I}_p)^{-1}$

alors que, pour rappel, $\text{EQM}(\hat{\beta}) = \text{Var}(\hat{\beta}) = \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1}.$

Régression Ridge

Si μ_j valeur propre de $\mathbf{X}^\top \mathbf{X}$

- ▶ Sous \mathcal{H}_1 (nécessairement $r = p$) $\text{tr}(\text{EQM}(\hat{\boldsymbol{\beta}})) = \sum_{j=1}^p \frac{\sigma^2}{\mu_j}$;
- ▶ (potentiellement $r \leq p$):

$$\text{tr}(\text{EQM}(\hat{\boldsymbol{\beta}}_{\text{ridge}})) = \sum_{j=1}^r \frac{\sigma^2 \mu_j + \lambda^2 (\mathbf{P}^\top \boldsymbol{\beta})_j^2}{(\mu_j + \lambda)^2};$$

- ▶ $\text{tr}(\text{EQM}(\hat{\boldsymbol{\beta}}_{\text{ridge}})) \leq \text{tr}(\text{EQM}(\hat{\boldsymbol{\beta}})) \Leftrightarrow \lambda \leq \frac{2\sigma^2}{\boldsymbol{\beta}^\top \boldsymbol{\beta}}.$

Rétrécissement / *shrinkage*

- ▶ Soit $\hat{\theta}$ un estimateur de θ sans biais et de variance σ^2 , donc $\text{EQM}(\hat{\theta}) = \sigma^2$.
- ▶ Soit $\lambda > 0$ et $\tilde{\lambda} = \frac{\hat{\theta}}{1 + \lambda}$ alors

$$\mathbb{E}(\tilde{\theta}) = \frac{\theta}{1 + \lambda}, \quad \text{Var}(\tilde{\theta}) = \frac{\sigma^2}{(1 + \lambda)^2} \quad \text{et} \quad \text{EQM}(\tilde{\theta}) = \frac{\lambda^2 \theta^2 + \sigma^2}{(1 + \lambda)^2}.$$

- ▶ $\tilde{\theta}$ est donc biaisé mais de variance plus faible que $\hat{\theta}$.

On peut trouver λ^* telle que $\lambda \mapsto \text{tr}(\text{EQM}(\hat{\beta}_{\text{ridge}}(\lambda)))$ soit minimale

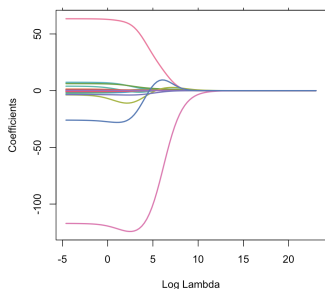
Régression Ridge

On peut utiliser

```
1 > library(MASS)
2 > ?lm.ridge
```

ou

```
1 > library(ISLR)
2 > library(glmnet)
3 > Hitters = na.omit(Hitters)
4 > x = model.matrix(Salary~.,
  Hitters)[,-1]
5 > y = Hitters$Salary
6 > ridge_mod = glmnet(x, y, alpha =
  0)
7 > plot(ridge_mod, var="lambda")
```



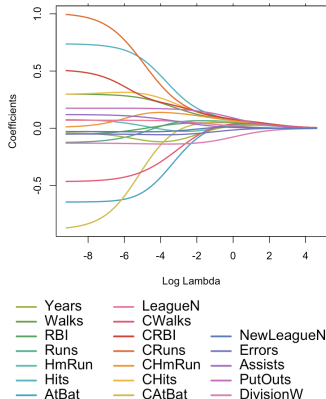
Years	LeagueN	
Walks	CWalks	
RBI	CRBI	NewLeagueN
Runs	CRuns	Errors
HmRun	CHmRun	Assists
Hits	CHits	PutOuts
AtBat	CAtBat	DivisionW

Régression Ridge

au lieu de $\|\beta\|_2$ (norme Euclidienne), on peut être tenté d'utiliser la norme de **Mahalanobis** i.e. on centre et on réduit les variables explicatives

$$\mathbf{x}_j \mapsto \frac{\mathbf{x}_j - \bar{\mathbf{x}}_j}{s_{x_j}}$$

```
1 > ys = (y-mean(y))/sd(y)
2 > xs = x
3 > for(i in 1:ncol(x)) xs[,i] = (x[,i]-mean(x[,i]))/sd(x[,i])
4 > ridge_mod_s = glmnet(xs, ys, alpha = 0)
5 > plot(ridge_mod_s, var="lambda")
```



Régression LASSO

- ▶ Ici, on ne contraint pas la norme euclidienne $\|\beta\| = \|\beta\|_2$ (i.e. la norme ℓ^2) mais la norme ℓ^1 des coefficients.
- ▶ La méthode Lasso consiste à minimiser

$$\|\mathbf{Y} - \mathbf{X}\beta\|^2 \text{ sous la contrainte } \|\beta\|_1 \leq \delta.$$

- ▶ Il n'existe pas de solution exacte (plusieurs algorithmes ont été proposés - LARS, coordinate descent algorithm).
- ▶ On peut montrer que le problème est équivalent à minimiser le problème

$$\hat{\beta}_{\text{lasso}}(\lambda) = \underset{\beta \in \mathbb{R}^p}{\operatorname{argmin}} \frac{1}{2} \|\mathbf{Y} - \mathbf{X}\beta\|^2 + \lambda \|\beta\|_1$$

($\forall \lambda > 0$, $\exists \delta > 0$ tel que les solutions de ces deux problèmes coïncident)...