

CHOIX DU PARAMETRE DE LISSAGE DANS LES MODELES DE STATISTIQUE NON PARAMETRIQUE EN GRANDE DIMENSION

Soutenance de Stage du Master 1 Statistique et Sciences de Données (SSD)

Kodjo Mawuena AMEKOE

28 août 2020



sous la direction de :

Mme Sana LOUHICHI

M. Didier GIRARD

M. Karim BENHENNI

Plan de l'exposé

- 1 Introduction
 - Contexte et motivation
 - Objectifs
- 2 Approximation de la validation croisée complète (ALO)
- 3 Majoration de l'écart entre ALO et LO
- 4 Simulations et discussions
- 5 Contributions
- 6 Conclusion et perspectives
- 7 Références

Contexte

Considérons le jeu de données $\mathcal{D} = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$, où $x_i \in \mathbb{R}^p$ et $y_i \in \mathbb{R}$ avec les observations (iid), issues de la distribution conjointe $q(y_i | x_i^\top \beta^*) p(x_i)$.

Contexte

Considérons le jeu de données $\mathcal{D} = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$, où $x_i \in \mathbb{R}^p$ et $y_i \in \mathbb{R}$ avec les observations (iid), issues de la distribution conjointe $q(y_i|x_i^\top \beta^*)p(x_i)$.

On cherche à estimer β^* pour faire la prédiction pour une nouvelle observation (x_{new}, y_{new}) issue de la distribution $q(y|x^\top \beta^*)p(x)$, indépendante de \mathcal{D} . Une méthode classique consiste à choisir la fonction de pénalité r , la fonction de perte l , le paramètre de lissage λ et à résoudre le problème d'optimisation :

$$\hat{\beta} = \arg \min_{\beta \in \mathbb{R}^p} \left[\sum_{i=1}^n l(y_i|x_i^\top \beta) + \lambda r(\beta) \right] \quad (1)$$

Contexte

Considérons le jeu de données $\mathcal{D} = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$, où $x_i \in \mathbb{R}^p$ et $y_i \in \mathbb{R}$ avec les observations (iid), issues de la distribution conjointe $q(y_i|x_i^\top \beta^*)p(x_i)$.

On cherche à estimer β^* pour faire la prédiction pour une nouvelle observation (x_{new}, y_{new}) issue de la distribution $q(y|x^\top \beta^*)p(x)$, indépendante de \mathcal{D} . Une méthode classique consiste à choisir la fonction de pénalité r , la fonction de perte l , le paramètre de lissage λ et à résoudre le problème d'optimisation :

$$\hat{\beta} = \arg \min_{\beta \in \mathbb{R}^p} \left[\sum_{i=1}^n l(y_i|x_i^\top \beta) + \lambda r(\beta) \right] \quad (1)$$

De même on s'intéresse à l'estimation à partir de \mathcal{D} de l'erreur de prédiction définie par :

$$Err_{extra} = E[\phi(y_{new}, x_{new}^\top \hat{\beta}) | \mathcal{D}] \quad (2)$$

.

Contexte

Considérons le jeu de données $\mathcal{D} = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$, où $x_i \in \mathbb{R}^p$ et $y_i \in \mathbb{R}$ avec les observations (iid), issues de la distribution conjointe $q(y_i|x_i^\top \beta^*)p(x_i)$.

On cherche à estimer β^* pour faire la prédiction pour une nouvelle observation (x_{new}, y_{new}) issue de la distribution $q(y|x^\top \beta^*)p(x)$, indépendante de \mathcal{D} . Une méthode classique consiste à choisir la fonction de pénalité r , la fonction de perte l , le paramètre de lissage λ et à résoudre le problème d'optimisation :

$$\hat{\beta} = \arg \min_{\beta \in \mathbb{R}^p} \left[\sum_{i=1}^n l(y_i|x_i^\top \beta) + \lambda r(\beta) \right] \quad (1)$$

De même on s'intéresse à l'estimation à partir de \mathcal{D} de l'erreur de prédiction définie par :

$$Err_{extra} = E[\phi(y_{new}, x_{new}^\top \hat{\beta}) | \mathcal{D}] \quad (2)$$

Une technique simple et intuitive est la **validation croisée**.

Motivation

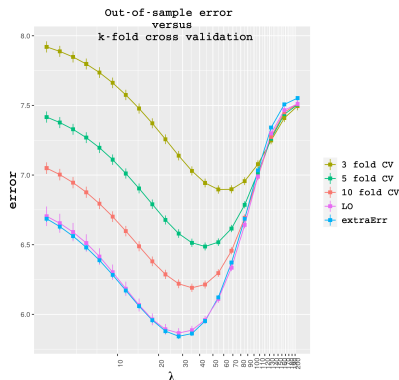


FIGURE: Comparaison de la validation croisée à 3, 5, 10 parties (folds) et du LO avec l'erreur de prédiction.

Motivation

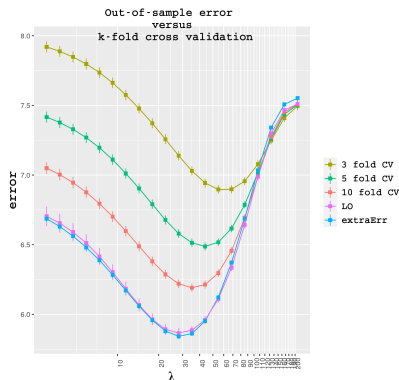


FIGURE: Comparaison de la validation croisée à 3, 5, 10 parties (folds) et du LO avec l'erreur de prédiction.

20 différentes valeurs λ pour une régression linéaire avec la pénalité LASSO.
(p, n, k) = (1000, 250, 50). $y \sim \mathcal{N}(X\beta^*, \sigma^2 I)$ et $Err_{extra} = \sigma^2 + \|\beta^* - \hat{\beta}\|_2^2$.

- La technique de validation croisée souffre d'un biais large sauf si le nombre de folds (parties) est grand.

Motivation

- La technique de validation croisée souffre d'un biais large sauf si le nombre de folds (parties) est grand.
- Le Leave-One-Out (LO) est très coûteux en terme de ressources de calcul en grande dimension.

Challenge

Trouver une technique aussi précise que le LO mais moins coûteuse en ressource

Alternative

Approximate Leave-One-Out (ALO) pour une large classe des estimateurs comme LASSO [Tib96], Bridge [IEFF93] Elastic-net [ZH05] et méthodes de Régression linéaire, logistique, Poisson, robuste proposée par [RM18].

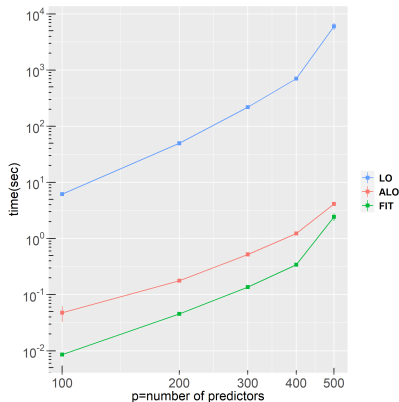


FIGURE: Régression linéaire avec pénalisation
Elastic-net pour $\frac{n}{p} = 5$.

- Comprendre le contexte du travail de [RM18] sur l'approximation du leave-one-out.

- Comprendre le contexte du travail de [RM18] sur l'approximation du leave-one-out.
- Discuter les résultats.

- Comprendre le contexte du travail de [RM18] sur l'approximation du leave-one-out.
- Discuter les résultats.
- Reproduire et faire d'autres simulations pour les configurations $n < p$, $n > p$ et $n = p$.

- Comprendre le contexte du travail de [RM18] sur l'approximation du leave-one-out.
- Discuter les résultats.
- Reproduire et faire d'autres simulations pour les configurations $n < p$, $n > p$ et $n = p$.
- Prendre connaissance des problèmes de statistique en grande dimension.

- Comprendre le contexte du travail de [RM18] sur l'approximation du leave-one-out.
- Discuter les résultats.
- Reproduire et faire d'autres simulations pour les configurations $n < p$, $n > p$ et $n = p$.
- Prendre connaissance des problèmes de statistique en grande dimension.
- Développer les compétences en programmation.

Plan de l'exposé

- 1 Introduction
- 2 Approximation de la validation croisée complète (ALO)
 - Fonctions de perte et pénalité lisses
 - Fonction de pénalité non lisse
- 3 Majoration de l'écart entre ALO et LO
- 4 Simulations et discussions
- 5 Contributions
- 6 Conclusion et perspectives
- 7 Références

La formule de l'estimation Leave-One-Out (LO) est donnée par :

$$LO = \frac{1}{n} \sum_{i=1}^n \phi(y_i, x_i^\top \hat{\beta}_{-/i}) \quad (3)$$

où

La formule de l'estimation Leave-One-Out (LO) est donnée par :

$$\text{LO} = \frac{1}{n} \sum_{i=1}^n \phi(y_i, \mathbf{x}_i^\top \hat{\beta}_{/i}) \quad (3)$$

où

$$\hat{\beta}_{/i} = \arg \min_{\beta \in \mathbb{R}^p} \left[\sum_{j \neq i} l(y_j | \mathbf{x}_j^\top \beta) + \lambda r(\beta) \right] \quad (4)$$

La formule de l'estimation Leave-One-Out (LO) est donnée par :

$$\text{LO} = \frac{1}{n} \sum_{i=1}^n \phi(y_i, x_i^\top \hat{\beta}_{/i}) \quad (3)$$

où

$$\hat{\beta}_{/i} = \arg \min_{\beta \in \mathbb{R}^p} \left[\sum_{j \neq i} l(y_j | x_j^\top \beta) + \lambda r(\beta) \right] \quad (4)$$

Challenge

LO est très coûteux en grande dimension : n calculs de $\hat{\beta}_{/i}$.

Solution

En supposant que $\hat{\beta}_{/i}$ est proche de $\hat{\beta}$, on utilise la méthode de Newton qui permet d'obtenir une approximation de $\hat{\beta}_{/i}$.

Solution

Avec quelques manipulations algébriques, on obtient ALO :

$$\text{ALO} = \frac{1}{n} \sum_{i=1}^n \phi(y_i, x_i^\top \tilde{\beta}_{/i}) = \frac{1}{n} \sum_{i=1}^n \phi \left(y_i, x_i^\top \hat{\beta} + \left(\frac{\dot{l}_i(\hat{\beta})}{\ddot{l}_i(\hat{\beta})} \right) \left(\frac{H_{ii}}{1 - H_{ii}} \right) \right) \quad (5)$$

avec

$$H = X(\lambda \text{diag}[\ddot{r}(\hat{\beta})] + X^\top \text{diag}[\ddot{l}(\hat{\beta})]X)^{-1} X^\top \text{diag}[\ddot{l}(\hat{\beta})]. \quad (6)$$

Challenge

L'approche utilisée ci-dessus nécessite des fonctions de perte et pénalité lisses, donc ne peut pas s'appliquer directement à :

- LASSO :

$$\hat{\beta} = \arg \min_{\beta \in \mathbb{R}^p} \left[\sum_{i=1}^n l(y_i | x_i^\top \beta) + \lambda \|\beta\|_1 \right] \quad (7)$$

Challenge

L'approche utilisée ci-dessus nécessite des fonctions de perte et pénalité lisses, donc ne peut pas s'appliquer directement à :

- LASSO :

$$\hat{\beta} = \arg \min_{\beta \in \mathbb{R}^p} \left[\sum_{i=1}^n l(y_i | x_i^\top \beta) + \lambda \|\beta\|_1 \right] \quad (7)$$

- Elastic-net :

$$\hat{\beta} = \arg \min_{\beta \in \mathbb{R}^p} \left[\sum_{i=1}^n l(y_i | x_i^\top \beta) + \lambda_1 \|\beta\|_2^2 + \lambda_2 \|\beta\|_1 \right]. \quad (8)$$

Solution

Faire une approximation lisse de la fonction de pénalité r : approximation de Schmidt [MGR07] pour LASSO.

Plan de l'exposé

- 1 Introduction
- 2 Approximation de la validation croisée complète (ALO)
- 3 Majoration de l'écart entre ALO et LO**
- 4 Simulations et discussions
- 5 Contributions
- 6 Conclusion et perspectives
- 7 Références

Théorème important

[RM18] a montré avec une probabilité qui tend vers 1 avec n et p que

$$|ALO - LO| \leq O_p\left(\frac{\text{PolyLog}(n)}{\sqrt{n}}\right) \text{ pour } r(\beta) = \gamma\beta^2 + (1 - \gamma)r^\alpha(\beta) \text{ avec } 0 < \gamma < 1$$

(Elastic-net).

En supposant que :

- $\frac{n}{p} = \delta_0$, avec δ_0 un nombre fini ,borné loin de zéro(0) ;

Théorème important

[RM18] a montré avec une probabilité qui tend vers 1 avec n et p que

$|ALO - LO| \leq O_p\left(\frac{\text{PolyLog}(n)}{\sqrt{n}}\right)$ pour $r(\beta) = \gamma\beta^2 + (1 - \gamma)r^\alpha(\beta)$ avec $0 < \gamma < 1$ (Elastic-net).

En supposant que :

- $\frac{n}{p} = \delta_0$, avec δ_0 un nombre fini, borné loin de zéro(0) ;
- les lignes de $X \in \mathbb{R}^{n \times p}$ sont indépendantes et suivent une loi Gaussienne de moyenne 0, de matrice de covariance Σ dont la valeur propre maximale est $\rho_{\max} = \frac{c}{n}$, c une constante ;

Théorème important

[RM18] a montré avec une probabilité qui tend vers 1 avec n et p que

$$|ALO - LO| \leq O_p\left(\frac{\text{PolyLog}(n)}{\sqrt{n}}\right) \text{ pour } r(\beta) = \gamma\beta^2 + (1 - \gamma)r^\alpha(\beta) \text{ avec } 0 < \gamma < 1$$

(Elastic-net).

En supposant que :

- $\frac{n}{p} = \delta_0$, avec δ_0 un nombre fini ,borné loin de zéro(0) ;
- les lignes de $X \in \mathbb{R}^{n \times p}$ sont indépendantes et suivent une loi Gaussienne de moyenne 0 , de matrice de covariance Σ dont la valeur propre maximale est $\rho_{\max} = \frac{c}{n}$, c une constante ;
- la dérivée seconde des fonctions de perte et pénalité est régulière ;

Théorème important

[RM18] a montré avec une probabilité qui tend vers 1 avec n et p que

$|ALO - LO| \leq O_p\left(\frac{\text{PolyLog}(n)}{\sqrt{n}}\right)$ pour $r(\beta) = \gamma\beta^2 + (1 - \gamma)r^\alpha(\beta)$ avec $0 < \gamma < 1$ (Elastic-net).

En supposant que :

- $\frac{n}{p} = \delta_0$, avec δ_0 un nombre fini ,borné loin de zéro(0) ;
- les lignes de $X \in \mathbb{R}^{n \times p}$ sont indépendantes et suivent une loi Gaussienne de moyenne 0 , de matrice de covariance Σ dont la valeur propre maximale est $\rho_{\max} = \frac{c}{n}$, c une constante ;
- la dérivée seconde des fonctions de perte et pénalité est régulière ;
- $\phi(.,.) = l(.,.)$.

Plan de l'exposé

- 1 Introduction
- 2 Approximation de la validation croisée complète (ALO)
- 3 Majoration de l'écart entre ALO et LO
- 4 Simulations et discussions**
- 5 Contributions
- 6 Conclusion et perspectives
- 7 Références

Résultats expérimentaux de l'approximation du LO par ALO pour :

- Régression linéaire avec pénalisation Elastic-net ;

Résultats expérimentaux de l'approximation du LO par ALO pour :

- Régression linéaire avec pénalisation Elastic-net ;
- Régression logistique avec pénalisation LASSO ;

Résultats expérimentaux de l'approximation du LO par ALO pour :

- Régression linéaire avec pénalisation Elastic-net ;
- Régression logistique avec pénalisation LASSO ;
- Régression Poisson Elastic-net ;

Régression linéaire avec pénalisation Elastic-net

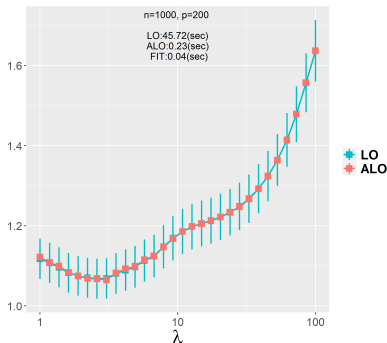


FIGURE: $n > p$

Régression linéaire avec pénalisation Elastic-net

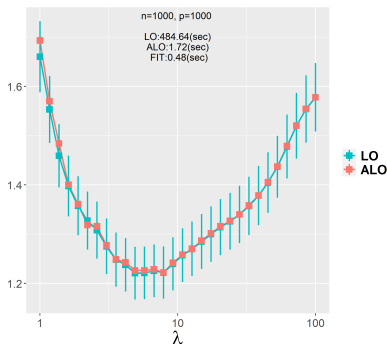


FIGURE: $n = p$

Régression linéaire avec pénalisation Elastic-net

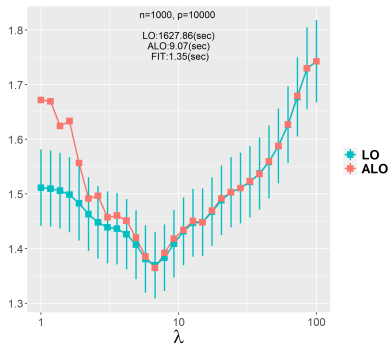


FIGURE: $n < p$

Plan de l'exposé

- 1 Introduction
- 2 Approximation de la validation croisée complète (ALO)
- 3 Majoration de l'écart entre ALO et LO
- 4 Simulations et discussions
- 5 Contributions**
- 6 Conclusion et perspectives
- 7 Références

Impact du paramètre de pondération de l'elastic-net mis à zéro : $\gamma = 0$.

Approximation de Schmidt :

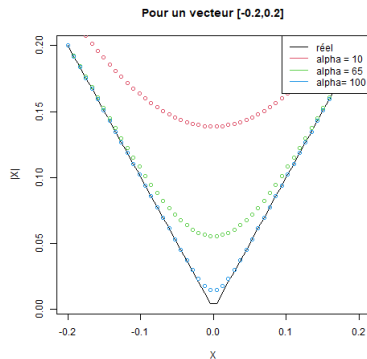


FIGURE: Séquence de 50 points dans l'intervalle $[-0.2, 0.2]$

Pénalisation de Schmidt :

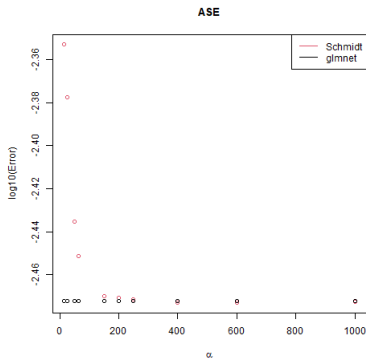


FIGURE: Comparaison ASE Estimation de Schmidt Vs LASSO glmnet

Programmation de ALO^α

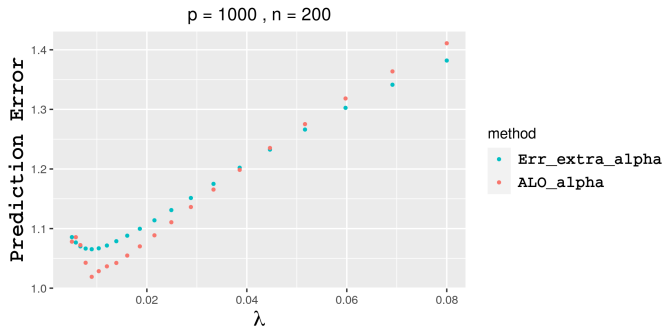


FIGURE: Comparaison ASE Estimation de Schmidt Vs LASSO glmnet

ALO^α VS ALO

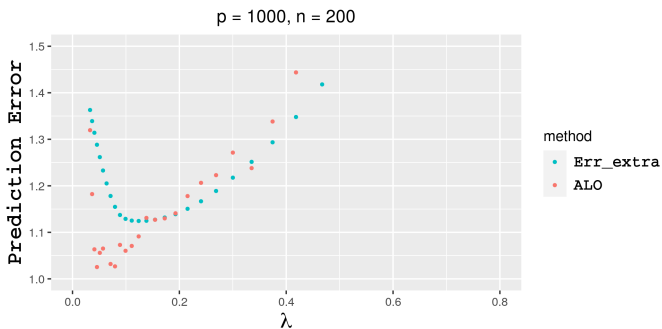


FIGURE: Comparaison de l'erreur de prédiction et ALO pour une régression linéaire LASSO

Influence de fortes et faibles corrélations : X avec la structure de corrélation Toeplitz : $\text{corr}(X_{ij}, X_{ij'}) = \rho^{|j-j'|}$ avec $i = 1, 2, \dots, n$, $j, j' = 1, 2, \dots, p$ et $\rho \in \{0.02, 0.7, 0.99\}$.

Plan de l'exposé

- 1 Introduction
- 2 Approximation de la validation croisée complète (ALO)
- 3 Majoration de l'écart entre ALO et LO
- 4 Simulations et discussions
- 5 Contributions
- 6 Conclusion et perspectives**
- 7 Références

Conclusion

- Missions remplies.

Conclusion

- Missions remplies.
- Importance d'avoir les observations(iid) et $\frac{n}{p}$ fixe.

Conclusion

- Missions remplies.
- Importance d'avoir les observations(iid) et $\frac{n}{p}$ fixe.
- Visualisations avec ggplot2, optimisation des programmes, parallélisation sous R.

- Ressources de calcul, temps limité.

Difficultés et Perspectives

- Ressources de calcul, temps limité.
- Traitement des données réelles.






- Ressources de calcul, temps limité.
- Traitement des données réelles.
- Optimisation du code de l'approximation de Schmidt et calcul de ALO^α .

- Ressources de calcul, temps limité.
- Traitement des données réelles.
- Optimisation du code de l'approximation de Schmidt et calcul de ALO^α .
- Utilisation d'approximations de Monte-Carlo.

Plan de l'exposé

- 1 Introduction
- 2 Approximation de la validation croisée complète (ALO)
- 3 Majoration de l'écart entre ALO et LO
- 4 Simulations et discussions
- 5 Contributions
- 6 Conclusion et perspectives
- 7 Références**

References I

-  Ildiko E. Frank and Jerome H. Friedman, *A statistical view of some chemometrics regression tools*, Technometrics **35** (1993), no. 2, 109–135.
-  Schmidt M., Fung G., and Rosales R, *Fast optimization methods for l1 regularization : A comparative study and two new approaches*, ECML **4701** (2007), 286–297.
-  Kamiar Rahnama Rad and Arian Maleki, *A scalable estimate of the extra-sample prediction error via approximate leave-one-out*, arXiv : Methodology (2018).
-  Robert Tibshirani, *Regression shrinkage and selection via the lasso*, Journal of the Royal Statistical Society. Series B (Methodological) **58** (1996), no. 1, 267–288.
-  Hui Zou and Trevor Hastie, *Regularization and variable selection via the elastic net*, Journal of the Royal Statistical Society. Series B (Statistical Methodology) **67** (2005), no. 2, 301–320.

Merci pour votre attention.