

# La REGRESSION RIDGE

La régression Ridge ordinaire ou bornée ordinaire a été proposée par E. Hoerl et Kennard dans " Ridge regression : biased estimation for nonorthogonal problems" Technometrics, Vol. 1 Février 1970.

On part de la constatation suivante : lorsque l'on se trouve en face d'un problème de forte colinéarité les valeurs propres de  ${}^tXX$ , les  $d_i^2 = \lambda_i$  correspondant à la colinéarité sont très petits. Hoerl et Kennard proposent de les augmenter faiblement en ajoutant à tous les  $d_i^2$  une même constante  $K > 0$ , ceci dans le but de rendre stables les coefficients. Regardons ce que donne cette proposition.

Remarque: Ils ont également proposé une autre méthode, la Ridge Généralisée, pour laquelle la constante ajoutée varie suivant les variables; cette méthode est nettement plus complexe.

## 1 Définition de la Ridge ordinaire (R.O.)

### 1.1 Sur le modèle canonique

Nous redéfinissons le modèle canonique ( voir cours sur la colinéarité) dans un modèle cette fois **centré et réduit**.

Soit le modèle  $\vec{Y} = X\vec{a} + \vec{\epsilon}$ ,  $V$  la matrice orthogonale des vecteurs propres ( $V^tV = I$ ) de  ${}^tXX$  et  $\Lambda$  sa matrice diagonale des valeurs propres  $\lambda_i$ .

On  ${}^tXX = V\Lambda^tV$  et  $\vec{Y} = X\vec{a} + \vec{\epsilon} = XV^tV\vec{a} + \vec{\epsilon}$

En notant le vecteur  $XV = Z$  et  ${}^tV\vec{a} = \vec{\gamma}$  le modèle peut s'écrire  $\vec{Y} = Z\vec{\gamma} + \vec{\epsilon}$  c'est le modèle sous sa forme canonique

- On constate que  ${}^tZZ = {}^tV^tXXV = {}^tVV\Lambda^tVV = \Lambda$  matrice diagonale des vecteurs propres ou des carrés des valeurs singulières

$${}^tZZ = \begin{pmatrix} \lambda_1 & 0 & 0 \\ 0 & \lambda_i & 0 \\ 0 & 0 & \lambda_k \end{pmatrix} = \begin{pmatrix} d_1^2 & 0 & 0 \\ 0 & d_i^2 & 0 \\ 0 & 0 & d_k^2 \end{pmatrix}$$

donc  $V_{\vec{\gamma}} = \sigma^2 ({}^tZZ)^{-1} = \sigma^2\Lambda^{-1}$  matrice diagonale.

- Dans ces conditions la régression Ridge Ordinaire donne le résultat suivant

Hoerl et Kennard proposent donc d'augmenter "un peu" les  $d_i^2$  en ajoutant à tous (même à ceux qui sont grands) la même constante  $K > 0$ . Cela conduit à augmenter tous les termes de la valeur  $K$ , comme  ${}^tZZ$  est une matrice diagonale, on obtient

$${}^tzz + KI = \begin{pmatrix} d_1^2 + K & .. & 0 \\ .. & .. & .. \\ 0 & .. & d_k^2 + K \end{pmatrix}$$

L'estimateur canonique de la Ridge s'écrit donc

$$\vec{\widehat{\gamma r}} = ({}^tZZ + KI)^{-1} {}^t_z \vec{y}$$

Nous discuterons par la suite de la valeur à donner à K. Il devra être petit, mais que signifie "petit"? On le cernerá mieux en prenant comme on l'a indiqué au début un modèle dans lequel les variables sont centrées et réduites.

## 1.2 Estimateur Ridge dans le modèle de base

$$\vec{\widehat{ar}} = V\vec{\widehat{\gamma r}} = V({}^tV{}^tXXV + K{}^tVV)^{-1} {}^tV{}^tX\vec{y} = V[{}^tV({}^tXX + KI)V]^{-1} {}^tV{}^tX\vec{y}$$

$$\vec{\widehat{ar}} = V{}^tV({}^tXX + KI)^{-1} {}^tV{}^tX\vec{y}$$

$$\vec{\widehat{ar}} = ({}^tXX + KI)^{-1} X\vec{y}$$

Car la matrice des vecteurs propres V est orthogonale. La ridge consiste donc à ajouter la constante K à tous les éléments diagonaux de  ${}^tXX$ . C'est la seconde définition de la ridge.

## 2 Propriétés de l'estimateur R.O.

### 2.1 Espérance de l'estimateur Ridge

$$E(\vec{\widehat{ar}}) = E(({}^tXX + KI)^{-1} X\vec{y}) = E(({}^tXX + KI)^{-1} X(X\vec{a} + \vec{\epsilon}))$$

Si on fait les hypothèses classiques des erreurs d'espérance nulle et des variables explicatives non aléatoires, on obtient,

$$E(\vec{\widehat{ar}}) = E(({}^tXX + KI)^{-1} X(X\vec{a} + \vec{\epsilon})) = ({}^tXX + KI)^{-1} X\vec{a}$$

$$E(\vec{\widehat{ar}}) = ({}^tXX + KI)^{-1} ({}^tXX + KI - KI)\vec{a} = \vec{a} - ({}^tXX + KI)^{-1} \vec{a}$$

De même

$$E(\vec{\widehat{\gamma r}}) = \vec{\gamma} - ({}^tZZ + KI)^{-1} \vec{\gamma}$$

L'estimateur Ridge est donc biaisé.

### 2.2 Variance de l'estimateur Ridge

$$V_{\vec{\widehat{ar}}} = ({}^tXX + KI)^{-1} X V_{\vec{y}} X ({}^tXX + KI)^{-1}$$

Si on fait l'hypothèse de non autocorrélation et homoscedasticité alors  $V_{\vec{y}} = \sigma^2 I$

$$V_{\vec{\widehat{ar}}} = \sigma^2 ({}^tXX + KI)^{-1} X X ({}^tXX + KI)^{-1}$$

de même

$$V_{\vec{\widehat{\gamma r}}} = \sigma^2 ({}^tZZ + KI)^{-1} Z Z ({}^tZZ + KI)^{-1}$$

## 2.3 Conclusion

Nous nous trouvons face à un estimateur biaisé. Est-il meilleur que les MCO ? Mystère car pour l'instant nous ne savons comparer que des estimateurs sans biais : on prend celui de variance Minimum. Mais pour comparer des estimateurs sans biais avec des estimateurs biaisés il faut introduire un nouveau critère celui du M.S.E.

## 3 Le critère du Mean Square Erreur (M.S.E)

### 3.1 Définition du MSE

La colinéarité entraîne comme nous l'avons vu des estimateurs qui bien que sans biais ont des variances trop grandes donc des intervalles de confiance trop grands. Il est parfois plus judicieux de prendre des estimateurs un peu biaisés mais avec des variances beaucoup plus faibles, ce qui conduira à des intervalles de confiance plus petits

(-----[----- $\vec{a} - \vec{b}$ -----]-----)

Avec  $\vec{b} = E(\vec{\hat{a}})$  où  $\vec{\hat{a}}$  est un estimateur biaisé de  $\vec{a}$

Ainsi sur ce graphique (très laid) l'estimateur biaisé est bien meilleur (intervalle de confiance entre les deux crochets) que l'estimateur sans biais mais qui a un intervalle de confiance (entre les deux parenthèses) beaucoup plus grand. Mais comment comparer ces deux estimateurs? Quand des estimateurs sont sans biais le meilleur est celui qui a la plus petite variance. Mais quand l'un est biaisé et l'autre sans biais comment choisir un critère sachant que le biais ne doit pas être grand.

Il existe un critère, celui du MSE (mean square estimator), erreur quadratique moyenne.

Définition : Si le vecteur  $\vec{a}$  est estimé par  $\vec{\hat{a}}$  estimateur biaisé de  $\vec{a}$  donc d'espérance  $E(\vec{\hat{a}}) \neq \vec{a}$ , on va définir le MSE par la matrice

$$MSE(\vec{\hat{a}}) = E[(\vec{\hat{a}} - \vec{a})^t (\vec{\hat{a}} - \vec{a})]$$

Le travail se fera plutôt avec la trace de cette matrice en utilisant les propriétés de la trace qui est la somme des éléments de la diagonale du MSE

$$trMSE(\vec{\hat{a}}) = E[t(\vec{\hat{a}} - \vec{a}) (\vec{\hat{a}} - \vec{a})]$$

$$\begin{aligned} trMSE(\vec{\hat{a}}) &= E[t(\vec{\hat{a}} - E(\vec{\hat{a}}) + E(\vec{\hat{a}}) - \vec{a}) (\vec{\hat{a}} - E(\vec{\hat{a}}) + E(\vec{\hat{a}}) - \vec{a})] \\ &= E[t(\vec{\hat{a}} - E(\vec{\hat{a}})) (\vec{\hat{a}} - E(\vec{\hat{a}}))] + E[t(E(\vec{\hat{a}}) - \vec{a}) (\vec{\hat{a}} - E(\vec{\hat{a}}))] \\ &\quad + E[t(\vec{\hat{a}} - E(\vec{\hat{a}})) (E(\vec{\hat{a}}) - \vec{a})] + E[t(E(\vec{\hat{a}}) - \vec{a}) (E(\vec{\hat{a}}) - \vec{a})] \end{aligned}$$

or  $E[t(\vec{\hat{a}} - E(\vec{\hat{a}})) (E(\vec{\hat{a}}) - \vec{a})] = [t(E(\vec{\hat{a}}) - E(\vec{\hat{a}})) (E(\vec{\hat{a}}) - \vec{a})] = \vec{0} = E[t(E(\vec{\hat{a}}) - \vec{a}) (\vec{\hat{a}} - E(\vec{\hat{a}}))]$

donc

$$trMSE(\vec{\hat{a}}) = E[t(\vec{\hat{a}} - E(\vec{\hat{a}})) (\vec{\hat{a}} - E(\vec{\hat{a}}))] + E[t(\vec{\hat{a}} - E(\vec{\hat{a}})) (E(\vec{\hat{a}}) - \vec{a})]$$

$$trMSE(\vec{\hat{a}}) = \sum_1^k var(\hat{a}_i) + \sum_1^k (biais)^2$$

Dans le cas d'un estimateur sans biais la trace du MSE =  $\sum_1^k var(\hat{a}_i)$

### 3.2 Propriétés du MSE

Le MSE permet de comparer deux estimateurs qu'ils soient sans biais ou biaisés. On prendra l'estimateur qui a la plus petite trMSE.

Théorème : la trMSE d'un modèle est identique à celui de son modèle canonique qui est beaucoup plus facile à calculer.

#### 3.2.1 Rappel du modèle canonique

Soit le modèle  $\vec{Y} = X\vec{a} + \vec{\epsilon}$ ,  $V$  la matrice orthogonale des vecteurs propres ( $V^tV = I$ ) de  ${}^tXX$  et  $\Lambda$  sa matrice diagonale des valeurs propres  $\lambda_i$ .

On  ${}^tXX = V\Lambda^tV$  et  $\vec{Y} = X\vec{a} + \vec{\epsilon} = XV^tV\vec{a} + \vec{\epsilon}$

En notant le vecteur  $XV = Z$  et  ${}^tV\vec{a} = \vec{\gamma}$  le modèle peut s'écrire  $\vec{Y} = Z\vec{\gamma} + \vec{\epsilon}$  c'est le modèle sous sa forme canonique

avec  $V_{\vec{\gamma}} = \sigma^2 ({}^tZZ)^{-1} = \sigma^2\Lambda^{-1}$  matrice diagonale.

- le MSE du modèle canonique s'écrit dans un modèle sans biais:

$$trMSE(\vec{\gamma}) = \sigma^2 \sum_1^k 1/\lambda_i = \sigma^2 \sum_1^k 1/d_i^2 \text{ ou les } d_i \text{ sont les valeurs singulières de la matrice } X$$

#### 3.2.2 Théorème

Le modèle de base et son modèle canonique ont même trMSE.  $trMSE(\vec{\hat{a}}) = trMSE(\vec{\hat{\gamma}})$

Démonstration:

$$trMSE(\vec{\hat{a}}) = E[{}^t(\vec{\hat{a}} - \vec{a})(\vec{\hat{a}} - \vec{a})] = E[{}^t(V\vec{\hat{\gamma}} - V\vec{\gamma})(V\vec{\hat{\gamma}} - V\vec{\gamma})]$$

$$trMSE(\vec{\hat{a}}) = E[{}^t(\vec{\hat{\gamma}} - \vec{\gamma}){}^tVV(\vec{\hat{\gamma}} - \vec{\gamma})] = E[{}^t(\vec{\hat{\gamma}} - \vec{\gamma})(\vec{\hat{\gamma}} - \vec{\gamma})] = MSE(\vec{\hat{\gamma}})$$

## 4 Comparaison des MCO et de la R.O.

On va effectuer les calculs sur les modèles canoniques et toujours sur des données centrées réduites.

### 4.1 MSE de la R.O.

$$\begin{aligned} trMSE(\vec{\hat{ar}}) &= trMSE(\vec{\hat{\gamma r}}) = \sum var(\hat{\gamma r}_i) + \sum biais^2 \\ &= \sum var(\hat{\gamma r}_i) + \sum (\gamma r_i - E(\hat{\gamma r}_i))^2 \end{aligned}$$

or

$$E(\vec{\widehat{\gamma r}}) = \vec{\gamma} - \begin{pmatrix} K/(d_1^2 + K) & 0 & \dots & 0 \\ 0 & K/(d_2^2 + K) & \dots & 0 \\ 0 & \dots & \dots & 0 \\ 0 & 0 & \dots & K/(d_k^2 + K) \end{pmatrix} \begin{pmatrix} \gamma_1 \\ \gamma_2 \\ \dots \\ \gamma_k \end{pmatrix}$$

$$\sum (\gamma_i - E(\widehat{\gamma r}_i))^2 = \sum_{i=1}^k K^2 \gamma_i^2 / (d_i^2 + K)^2$$

La matrice de Variance-covariance de  $\vec{\widehat{\gamma}}$  s'écrit:

$$\begin{aligned} Var_{\vec{\widehat{\gamma r}}} &= \sigma^2 \begin{pmatrix} 1/(d_1^2 + K) & \dots & 0 \\ 0 & \dots & 0 \\ 0 & \dots & 1/(d_k^2 + K) \end{pmatrix} \begin{pmatrix} d_1^2 & 0 \\ 0 & \dots & 0 \\ 0 & \dots & d_k^2 \end{pmatrix} \begin{pmatrix} 1/(d_1^2 + K) & \dots & 0 \\ 0 & \dots & 0 \\ 0 & \dots & 1/(d_k^2 + K) \end{pmatrix} \\ &= \sigma^2 \begin{pmatrix} d_1^2 / (d_1^2 + K)^2 & \dots & 0 \\ 0 & \dots & 0 \\ 0 & \dots & d_k^2 / (d_k^2 + K)^2 \end{pmatrix} \end{aligned}$$

donc  $\sum var(\widehat{\gamma r}_i) = \sigma^2 \sum d_i^2 / (d_i^2 + K)^2$

$$trMSE(\vec{\widehat{\gamma r}}) = \sigma^2 \sum_{i=1}^k d_i^2 / (d_i^2 + K)^2 + \sum_{i=1}^k K^2 \gamma_i^2 / (d_i^2 + K)^2$$

La trMSE de l'estimateur Ridge s'écrit

$$trMSE(\vec{\widehat{\gamma r}}) = \sum \frac{\sigma^2 d_i^2 + K^2 \gamma_i^2}{(d_i^2 + K)^2}$$

## 4.2 trMSE des MCO

La trMSE de l'estimateur des MCO  $\vec{\widehat{\gamma}}$  est ( pour K=0)

$$trMSE(\vec{\widehat{\gamma}}) = \sum \frac{\sigma^2 d_i^2}{(d_i^2)^2} = \sum \frac{\sigma^2}{d_i^2}$$

## 4.3 Comparaison

Nous avons vu que la trMSE de l'estimateur Ridge s'écrit

$$trMSE(\vec{\widehat{\gamma r}}) = \sum \frac{\sigma^2 d_i^2 + K^2 \gamma_i^2}{(d_i^2 + K)^2}$$

La trMSE de l'estimateur des MCO  $\vec{\widehat{\gamma}}$  est ( pour K=0)

$$trMSE(\vec{\widehat{\gamma}}) = \sum \frac{\sigma^2 d_i^2}{(d_i^2)^2} = \sum \frac{\sigma^2}{d_i^2}$$

Pour que la Ridge soit meilleur estimateur que les MCO il faut suivant le critère du MSE que

$$trMSE(\vec{\gamma r}) < trMSE(\vec{\gamma})$$

$$\sum \frac{\sigma^2 d_i^2 + K^2 \gamma_i^2}{(d_i^2 + K)^2} < \sum \frac{\sigma^2}{d_i^2}$$

Pour cela il suffit de vérifier l'inégalité pour chaque  $i$ . Nous allons donc essayer de trouver une condition suffisante pour que la Ridge soit meilleure ( au sens du MSE) que les MCO. Il suffit donc d'avoir pour chaque  $i$

$$\frac{\sigma^2 d_i^2 + K^2 \gamma_i^2}{(d_i^2 + K)^2} < \frac{\sigma^2}{d_i^2}$$

$$\sigma^2 d_i^4 + K^2 \gamma_i^2 d_i^2 < \sigma^2 (d_i^4 + K^2 + 2K d_i^2)$$

$$K^2 (\gamma_i^2 d_i^2 - \sigma^2) - 2K \sigma^2 d_i^2 < 0$$

C'est une inégalité du second degré en  $K$

- Si  $\gamma_i^2 d_i^2 - \sigma^2 > 0$  l'expression est négative entre les racines de l'équation du second degré soit entre 0 et une  $K_i^* = \frac{2\sigma^2 d_i^2}{\gamma_i^2 d_i^2 - \sigma^2} > 0$
- Si  $\gamma_i^2 d_i^2 - \sigma^2 < 0$  les deux expressions sont négatives et la condition est toujours vérifiée

On en déduit donc que le MSE de la Ridge est inférieur au MSE des MCO si  $0 < K < \min K_i^*$ . Ce minimum ne peut être calculé en pratique car il dépend des valeurs théorique des coefficients  $\gamma_i^2$ .

On prendra des valeurs de  $K$  très proches de 0 pour être dans le bon cas.

On vient donc de démontrer qu'il y a des valeurs de  $K$  pour lesquelles la régression Ridge est meilleure que les MCO.

## 5 Application de la ridge

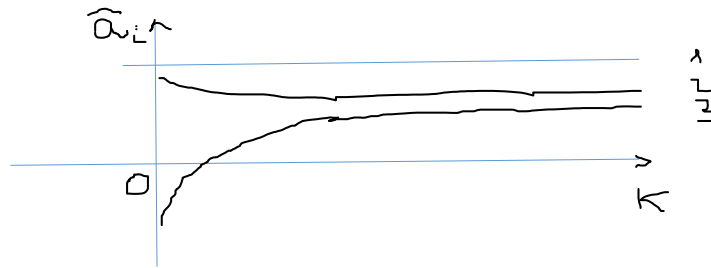
Sur un modèle à 3 variables. On utilise des valeurs de  $K$  très petites et on constate que certains coefficients restent fixes quand  $K$  augmente, d'autres bougent au début pour se stabiliser.

On remarque en théorie que pour  $K=0$  on a les MCO et que si  $K \rightarrow \infty$  alors les coefficients de la Ridge tendent vers 0, ce qui n'est pas le but recherché car  $K$  est choisi très petit d'après la théorie vue ci-dessus.

Interprétation du graphe (toujours aussi laid):

prenons l'exemple d'un modèle à 3 variables 1 2 et 3

- Variable 1 : elle ne joue pas de rôle dans la colinéarité donc si on ajoute  $K$  très petit, son coefficient ne change pas



- Variable 2 : elle joue un rôle dans la colinéarité avec la variable 3: son coefficient positif en  $K=0$  (MCO) baisse un peu puis se stabilise.
- Variable 3 : Son coefficient qui était négatif devient positif puis se stabilise.

Dans la plupart des modèles, on obtient ce genre de résultat. Si les variables ne sont pas colinéaires alors leur coefficient ne bouge pas en appliquant la ridge. Dans le cas des variables colinéaires certains coefficients bougent un peu et d'autres beaucoup avant de se stabiliser, toujours avec  $K$  très petit.

## 6 Choix de $K$

Nous rappelons que pour des raisons d'élimination des unités, nous travaillons sur des variables CENTREES REDUITES comme le suggèrent Hoerl et Kennard.

Le résultat de la Ridge  $\vec{\hat{a}}_r = ({}^tXX + KI)^{-1} {}^tX\vec{y}$

La matrice à inverser  ${}^tXX + KI$  correspond à la matrice

$$X_r = \begin{pmatrix} X \\ \sqrt{KI} \end{pmatrix}$$

telle que  ${}^tXX + KI = {}^tX_r X_r$ . Les valeurs propres de  ${}^tXX + KI$  donc sont les valeurs propres de  ${}^tXX$  plus  $K$ , donc les  $d_i^2 + K$

Le conditionnement de la matrice  $X$  est  $\text{Cond } X = d_{\max}/d_{\min}$

Le conditionnement de la matrice  $X_r$  est  $\text{Cond } X_r = \sqrt{d_{\max}^2 + K} / \sqrt{d_{\min}^2 + K}$

S'il y a colinéarité le conditionnement de la matrice  $X$  est très grand. Pour faire baisser ce conditionnement, on va passer à  $X_r$  en lui imposant un conditionnement beaucoup plus faible ( par exemple 50 ou 30) . Cela entraîne une valeur de  $K$

Si par exemple on impose 30

$$\text{Cond}X_r = \sqrt{d_{\max}^2 + K} / \sqrt{d_{\min}^2 + K} = 30$$

$$(d_{\max}^2 + K) = 30^2(d_{\min}^2 + K)$$

$$K = \frac{d_{\max}^2 - 30^2 d_{\min}^2}{30^2 - 1}$$

On prend cette valeur de K pour l'estimation de la régression Ridge.

On peut aussi (vivement recommandé) faire le graphe de l'évolution des coefficients et regarder quand les coefficients se stabilisent.