

# Apprentissage automatique / Statistique

## Qualité de Prévision et Risque

PHILIPPE BESSE

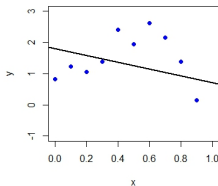
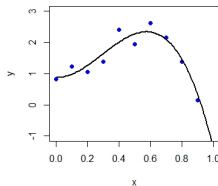
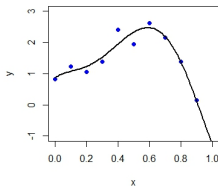
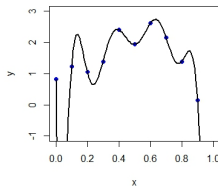
INSA de Toulouse  
Institut de Mathématiques

## Objectif

- Mesurer la **performance** d'un modèle, sa capacité de **prévision** ou de **généralisation**
  - Optimiser la **sélection** au sein d'une famille de modèles
  - **choix de la méthode** en comparant chacun des modèles
  - estimer la **confiance** accordée à une prévision
  - **Enjeu** : estimation sans biais de l'erreur de prévision
  - **Capacité** de généralisation du modèle ou algorithme

## Sans modèle probabiliste, trois stratégies :

- 1 **Pénalisation** de l'erreur d'ajustement ou risque empirique
- 2 **Partition** de l'échantillon : apprentissage, (validation), test
- 3 **Simulation** : validation croisée, bootstrap

Modèle linéaire;  $R^2=0.11$ Modèle cubique;  $R^2=0.95$ Degré 5;  $R^2=0.96$ Degré 10;  $R^2=1$ 

## Régression polynomiale de degré 1, 2, 5 et 10

## Notations

- $D_n$  observations d'un  $n$ -échantillon  
 $D_n = \{(X_1, Y_1), \dots, (X_n, Y_n)\}$  de loi conjointe inconnue  $P$  sur  $\mathcal{X} \times Y$
- $x$  observation de la variable  $X$  multidimensionnelle
- $D_n$  est appelé **échantillon d'apprentissage**
- $D_n$  est supposé indépendant de  $(X, Y)$
- Une **règle de prévision** (ou prédicteur) est une fonction (mesurable)  $f : \mathcal{X} \rightarrow \mathcal{Y}$ ,  $x \rightarrow f(x)$
- Une fonction  $l : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}_+$  est une **fonction de perte** si  $l(y, y) = 0$  et  $l(y, y') > 0$  pour  $y \neq y'$
- Si  $f$  est une règle de prévision,  $l(y, f(x))$  mesure la perte de  $f$  en  $x$

## Définitions

- **Régression réelle** : pertes  $\mathbb{L}^p$  ( $p \geq 1$ ) :  $l(y, y') = |y - y'|^p$   
perte absolue si  $p = 1$ , perte quadratique si  $p = 2$
- **Discrimination binaire** :  $l(y, y') = \mathbb{1}_{y \neq y'} = \frac{|y - y'|}{2} = \frac{(y - y')^2}{4}$
- **Risque d'une règle**  $f : R_P(f) = \mathbb{E}_{(X, Y) \sim P}[l(Y, f(X))]$
- **Risque empirique** associé à  $\mathbf{D}_n$  :  
 $\widehat{R}_n(f, \mathbf{D}_n) = \frac{1}{n} \sum_{i=1}^n l(Y_i, f(\mathbf{X}_i))$
- **Minimisation** du risque empirique sur un sous-ensemble  $F$   
(un modèle) de  $\mathcal{F} : \hat{f}_F(\mathbf{D}_n) \in \operatorname{argmin}_{f \in F} \widehat{R}_n(f, \mathbf{D}_n)$
- **Problème** : choix de  $F$  !
- La règle **oracle** est telle que :  $R_P(f^*) = \inf_f R_P(f)$

## Décomposition du risque empirique

$$R_P(\hat{f}_F(\mathbf{D}_n)) - R_P(f^*) = \underbrace{\left\{ R_P(\hat{f}_F(\mathbf{D}_n)) - \inf_{f \in F} R_P(f) \right\}}_{\substack{\text{Erreur d'estimation} \\ \text{(Variance)}}} + \underbrace{\left\{ \inf_{f \in F} R_P(f) - R_P(f^*) \right\}}_{\substack{\text{d'approximation} \\ \text{(Biais)}}}$$

↗
↘ (taille de  $F$ )

- **Plus** le modèle  $F$  est complexe ou flexible,
- plus le biais est réduit mais
- plus la partie variance risque d'augmenter
- **Enjeu** : meilleur compromis biais / variance

## Risque empirique ou qualité d'ajustement

$$\widehat{R}_n(\widehat{f}(\mathbf{D}_n), \mathbf{D}_n) = \frac{1}{n} \sum_{i=1}^n l(y_i, \widehat{f}(\mathbf{D}_n)(\mathbf{x}_i))$$

- Minimum des moindres carrés dans le cas quantitatif
- Taux de mal classés dans le cas qualitatif
- Estimation biaisée, par optimisme

## Estimation sans biais sur un échantillon indépendant

- **Partition** :  $\mathbf{D}_n = \mathbf{D}_{n_1}^{\text{Appr}} \cup \mathbf{D}_{n_2}^{\text{Valid}} \cup \mathbf{D}_{n_3}^{\text{Test}}$
- $\widehat{R}_n(\widehat{f}(\mathbf{D}_{n_1}^{\text{Appr}}), \mathbf{D}_{n_1}^{\text{Appr}})$  pour **estimer** un modèle choisi  $\widehat{f}(\mathbf{D}_{n_1}^{\text{Appr}})$
- $\widehat{R}_n(\widehat{f}(\mathbf{D}_{n_1}^{\text{Appr}}), \mathbf{D}_{n_2}^{\text{Valid}})$  pour **optimiser** un modèle
- $\widehat{R}_n(\widehat{f}, \mathbf{D}_{n_3}^{\text{Test}})$  pour **comparer** les meilleurs modèles



## $C_p$ de Mallows

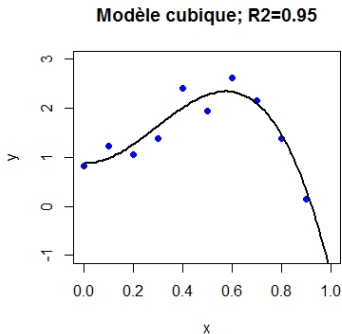
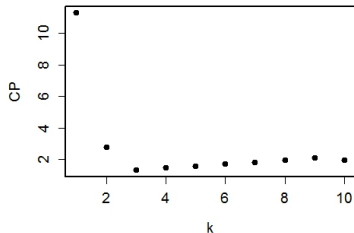
- Décomposition de l'**erreur** de prévision ou risque quadratique :

$$\widehat{R}_P(\widehat{f}(\mathbf{d}^n)) = \widehat{R}_n(\widehat{f}(\mathbf{d}^n), \mathbf{d}^n) + \text{Optim}$$

- Estimation normalisée :

$$C_p = \widehat{R}_n(\widehat{f}(\mathbf{d}^n), \mathbf{d}^n) + 2\frac{d}{n}\widehat{\sigma}^2$$

- $d$  : nombre de paramètres du modèle
- $n$  : nombre d'observations
- $s^2$  : estimation de la variance de l'erreur par modèle de **faible biais**



*Régression polynomiale :  $C_p$  de Mallows fonction du degré et modèle sélectionné.*

## Critère d'Akaïke

- Basé sur la **dissemblance** de Kullback
- compare la loi de  $Y$  et celle de  $\hat{Y}$ 
  - Suppose que la famille de lois du modèle contient la “vraie” loi de  $Y$
  - Pour tout modèle estimé par minimisation d'une **log-vraisemblance**  $\mathcal{L}$

$$AIC = -2\mathcal{L} + 2\frac{d}{n}$$

- Cas gaussien et variance connue : **AIC** et  $C_p$  **équivalents**
- $AIC_c$  adapté aux petits échantillons gaussiens

$$AIC_c = -2\mathcal{L} + \frac{n + d}{n - d - 2}$$

## Critère BIC de Schwarz

- BIC (**B**ayesian **i**nformation **c**riterion)
  - modèle de plus grande probabilité *a posteriori*

$$\text{BIC} = -2\mathcal{L} + \log(n)\frac{d}{n}$$

- Cas gaussien et variance connue : BIC proportionnel à AIC
- $n > e^2 \approx 7,4$ , BIC pénalise plus les modèles complexes
- Asymptotiquement, la probabilité pour BIC de choisir le bon modèle tend vers 1
- différent d'AIC qui tend à choisir des modèles trop complexes
- À Taille fini, BIC risque de se limiter à des modèles trop simples

## Algorithme de *V-fold cross validation*

- *Découper* aléatoirement l'échantillon en  $V$  segments ( $V$ -fold) de tailles similaires selon une loi uniforme ;
- **for**  $k=1$  à  $V$ 
  - Mettre de côté le segment  $k$ ,
  - Estimer le modèle sur les  $V - 1$  segments restants,
  - Calculer la moyenne des erreurs sur le segment  $k$
- **end for**
- Moyenner toutes les erreurs

## Utilisation

- Soit  $\tau : \{1, \dots, n\} \mapsto \{1, \dots, V\}$  la fonction d'**indexation**
- $\hat{f}^{(-v)}$  estimation de  $f$  sans le  $v$ ème segment de l'échantillon
- Estimation par **validation croisée** de l'erreur de prévision :

$$\widehat{R_{CV}} = \frac{1}{n} \sum_{i=1}^n l(y_i, \hat{f}^{(-\tau(i))}(x_i))$$

- Choix de  $V : n$  (variance), **petit** (biais), **10** par défaut
- **Utilisation fréquente** en choix de modèle :

$$\hat{\theta} = \arg \min_{\theta} \widehat{R_{CV}}(\theta)$$

## Introduction au Bootstrap

- Simulation (**Monte Carlo**) de la distribution d'un **estimateur**
- Principe : substituer  $P_n$ , à la distribution inconnue  $P$
- Tirage avec remise d'un **échantillon bootstrap** de même taille
- Itération et convergence

## Estimateur bootstrap naïf

- Échantillon bootstrap :  $\mathbf{z}^*$
- Estimateur **plug-in** (remplacer  $F$  par  $\hat{F}$  de  $R_P(\hat{f}(\mathbf{d}_n))$  :  
$$\hat{R}_n(\hat{f}_{\mathbf{z}^*}, \mathbf{d}_n) = \frac{1}{n} \sum_{i=1}^n l(y_i, \hat{f}_{\mathbf{z}^*}(\mathbf{x}_i))$$
- $\hat{f}_{\mathbf{z}^*}$  désigne l'estimation de  $f$  à partir de  $\mathbf{z}^*$
- Estimation bootstrap de l'erreur moyenne de prévision  
$$\mathbb{E}_{\mathbf{D}_n \sim P^{\otimes n}} [R_P(\hat{f}(\mathbf{D}_n))] :$$
- $R_{\text{Boot}} = E_{\mathbf{Z}^* \sim \hat{F}} [\hat{R}_n(\hat{f}_{\mathbf{Z}^*}, \mathbf{D}_n)] = E_{\mathbf{Z}^* \sim \hat{F}} \left[ \frac{1}{n} \sum_{i=1}^n l(y_i, \hat{f}_{\mathbf{Z}^*}(\mathbf{x}_i)) \right]$



## Estimation bootstrap par simulation

$$\widehat{R}_{\text{Boot}} = \frac{1}{B} \sum_{b=1}^B \frac{1}{n} \sum_{i=1}^n l(y_i, f_{\mathbf{z}^{*b}}(\mathbf{x}_i))$$

Estimation **biaisée** par optimisme

## Estimateur bootstrap out-of-bag

Distinguer les observations de l'échantillon **bootstrap** et les autres

$$\widehat{R}_{\text{oob}} = \frac{1}{n} \sum_{i=1}^n \frac{1}{B_i} \sum_{b \in K_i} l(y_i, f_{z^{*b}}(\mathbf{x}_i))$$

- $K_i$  est l'ensemble des indices  $b$  des échantillons **bootstrap** ne contenant pas la  $i$ ème observation à l'issue des  $B$  simulations
- $B_i = |K_i|$  est le nombre de ces échantillons
- $\widehat{R}_{\text{oob}}$  résout le problème d'un **biais optimiste** de  $\widehat{R}_{\text{Boot}}$  mais biais pessimiste comme en validation croisée ( $\widehat{R}_{\text{CV}}$ )

## Estimateur .632-bootstrap

- **Correctif** basé sur la
- **probabilité** qu'une observation soit **tirée** dans un échantillon **bootstrap** :

$$P[\mathbf{x}_i \in \mathbf{x}^{*b}] = 1 - \left(1 - \frac{1}{n}\right)^n \approx 1 - \frac{1}{e} \approx 0,632$$

- **Sur-évaluation** de l'erreur analogue à celle de la validation croisée avec  $K = 2$
- **Compensation** :

$$\widehat{R}_{.632} = 0,368 \times \widehat{R}_n(\widehat{f}(\mathbf{D}_n), \mathbf{D}_n) + 0,632 \times \widehat{R}_{\text{oob}}$$

## Remarques

- Estimations de l'erreur asymptotiquement équivalentes
- Pas de choix *a priori*
- **Bootstrap** plus compliqué et encore peu utilisé mais
- Central dans les algorithmes de combinaison de modèles
- Problèmes du .632-bootstrap en sur-ajustement
- Rectificatif complémentaire : le **.632+bootstrap**
- Utiliser le même estimateur pour comparer deux méthodes

## Matrice de confusion - Notations

**Prévision** : Si  $\hat{\pi}_i > s$ ,  $\hat{y}_i = 1$  sinon  $\hat{y}_i = 0$

Prévision	Observation		Total
	$Y = 1$	$Y = 0$	
$\hat{y}_i = 1$	$n_{11}(s)$	$n_{10}(s)$	$n_{1+}(s)$
$\hat{y}_i = 0$	$n_{01}(s)$	$n_{00}(s)$	$n_{0+}(s)$
Total	$n_{+1}$	$n_{+0}$	$n$

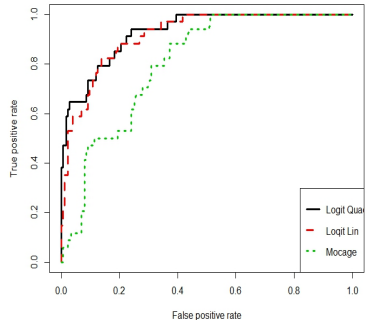
- Nombre de conditions positives  $P = n_{+1}$
- Nombre de conditions négatives  $N = n_{0+}$
- Vrais positifs  $TP = n_{11}(s)$  bien classés ( $\hat{y}_i = 1$  et  $Y = 1$ ),
- Vrais négatifs  $TN = n_{00}(s)$  bien classés ( $\hat{y}_i = 0$  et  $Y = 0$ ),
- Faux négatifs  $FN = n_{01}(s)$  mal classés ( $\hat{y}_i = 0$  et  $Y = 1$ ),
- Faux positifs  $FP = n_{10}(s)$  mal classés ( $\hat{y}_i = 1$  et  $Y = 0$ ),

## Notations

- **Accuracy** et Taux d'erreur :  $ACC = \frac{TN+TP}{N+P} = 1 - \frac{FN+FP}{N+P}$
- **Taux de vrais positifs**, *sensitivity*, *recall*  $TPR = \frac{TP}{P}$
- Taux de vrais négatifs *specificity*, *selectivity*  $TNR = \frac{TN}{N}$
- Précision ou **positive predictive value**  
 $PPV = \frac{TP}{TP+FP} = 1 - FDR$
- **Taux de faux positifs**  $FPR = \frac{FP}{N} = 1 - TNR$
- Taux de faux négatifs  $FNR = \frac{FN}{P} = 1 - TPR$
- Taux de fausses découvertes  $FDR = \frac{FP}{FN+TN}$
- $F_1$  score ou moyenne harmonique de la précision et de la sensibilité,  $F_1 = 2 \times \frac{PPV \times TPR}{PPV + TPR} = \frac{2 \times TP}{2 \times TP + FP + FN}$ ,
- $F_\beta (\beta \in \mathbb{R}^+)$  **score**,  $F_\beta = (1 + \beta^2) \frac{PPV \times TPR}{\beta^2 PPV + TPR}$ .

- Taux de vrais positifs : *TPR* ou *sensibilité*
- Taux de faux positifs :  
=  $FPR = 1 - \text{Spécificité}$
- **AUC** : aire sous la courbe
- **Score de Pierce** :  
 $H - F = TPR - FPN$
- **Log loss**

$$LI = -\frac{1}{n} \sum_{i=1}^n \sum_{k=1}^L y_{ik} \log(\widehat{\pi}_{ik})$$



*Banque : Courbes ROC et aire sous la courbe*