

Validation croisée (*cross-validation*)

Marie-Anne.Poursat@math.u-psud.fr

Master 1 Mathématiques Appliquées
Université Paris-Saclay

4 avril 2019

L'objectif est de construire à partir des données

$$\mathcal{D} = \{(\mathbf{X}_1, Y_1), \dots, (\mathbf{X}_n, Y_n)\}, \quad \mathbf{X}_i = (X_{i1}, \dots, X_{ip}),$$

un prédicteur $\hat{f}(\cdot) = \hat{f}(\cdot; \mathcal{D})$ (fonction de régression ou de classification).

- 1 Avant d'utiliser les prévisions que $\hat{f}(\cdot)$ fournit, il faut évaluer son risque.
- 2 Bien souvent on dispose de plusieurs règles d'apprentissage (différents choix d'hyperparamètres, différentes méthodes), comment choisir la meilleure méthode et ses paramètres ?

Évaluation de la qualité de prédiction

Risque en régression (*Mean Squared Error*) :

$$MSE = \mathbb{E} \left[(Y - \hat{f}(X))^2 \right]$$

ou en Classification (*Probabilité de mauvais classement*) :

$$P \left(Y \neq \hat{f}(X) \right)$$

Comment l'estimer ?

Approche naïve : par l'erreur d'apprentissage.

$$\frac{1}{n} \sum_{i=1}^n \left(Y_i - \hat{f}(X_i) \right)^2 \text{ (RSS) ou } \frac{1}{n} \sum_{i=1}^n \mathbb{I} \left\{ Y_i \neq \hat{f}(X_i) \right\} \text{ (taux de mal classés)}$$

↔ **sur-apprentissage**

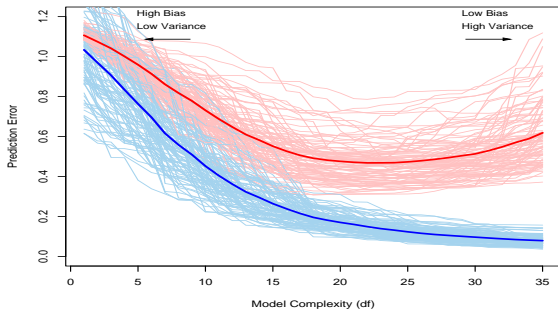


FIGURE 7.1. Behavior of test sample and training sample error as the model complexity is varied. The light blue curves show the training error $\overline{\text{err}}$, while the light red curves show the conditional test error Err_T for 100 training sets of size 50 each, as the model complexity is increased. The solid curves show the expected test error Err and the expected training error $\text{E}[\overline{\text{err}}]$.

Plus un modèle est complexe (plus de paramètres) et flexible (s'adapte mieux aux données), plus son erreur d'apprentissage est faible.

MAIS ce modèle pourra s'avérer mauvais dans un but de **prédiction** ou de **généralisation** (s'adapter à de nouvelles données).

↔ Minimiser l'erreur d'apprentissage conduit au **sur-apprentissage**.

Comment obtenir un bon estimateur de l'erreur de prédiction ?

↔ si l'on dispose de beaucoup de données, on partage le jeu de données en un ensemble d'apprentissage et un ensemble test :

- le modèle est **ajusté** sur les **données d'apprentissage** : calcul de \hat{f} .
- le modèle est **évalué** sur les **données test** : calcul de l'erreur test

$$n = n_{app} + n_{test}, \{(x_i, y_i), i = 1, \dots, n_{app}\} \mapsto \hat{f}$$

$$\text{Erreur test} = \frac{1}{n_{test}} \sum_{t=1}^{n_{test}} \left(y_t - \hat{f}(x_t) \right)^2$$

Avantages :

- simple,
- facile et rapide à implémenter

Inconvénients :

- l'erreur test peut être **variable**, dépendre du découpage des données
- l'erreur test peut **sur-estimer** l'erreur de prédiction en n'utilisant qu'une partie des données pour estimer f .

- Utiliser un jeu de données test de grande taille, rarement disponible.
- Corriger mathématiquement l'erreur d'apprentissage par pénalisation : AIC, BIC, Cp de Mallows,...
- Technique de simulation très utilisée : **validation croisée** (*cross-validation*)

$$Y = f(x) + \varepsilon, E(\varepsilon) = 0, V(\varepsilon) = \sigma^2.$$

$$E \left[(Y - \hat{f}(x))^2 \right] = \sigma^2 + \text{Biais}^2(\hat{f}(x), f(x)) + \text{Var}(\hat{f}(x)).$$

→ compromis biais-variance qui minimise l'erreur de prédiction ?
(*cost-complexity*)

$$E \left[(Y - \hat{f}(x))^2 \right] \simeq E[RSS/n] + \frac{2}{n} \sum \text{Cov}(y_i, \hat{y}_i)$$

et $\sum \text{Cov}(y_i, \hat{y}_i) = (p + 1)\sigma^2$ pour une régression linéaire gaussienne à p covariables.

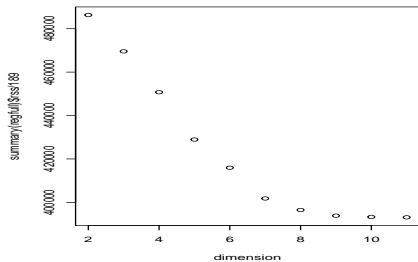
D'où

$$C_p = \frac{RSS}{n} + \frac{2}{n}(p + 1)\sigma^2$$

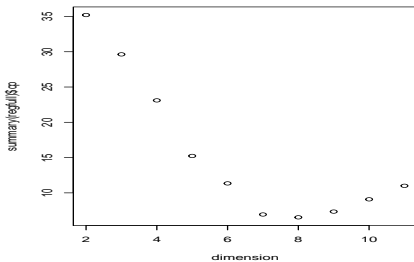
$$\text{et } AIC = -\frac{2}{n} \log \text{Lik} + \frac{2}{n}(p + 1).$$

Exemple Birthwt

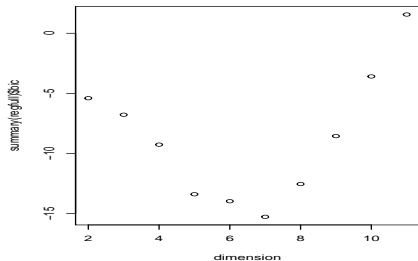
Erreur d'apprentissage



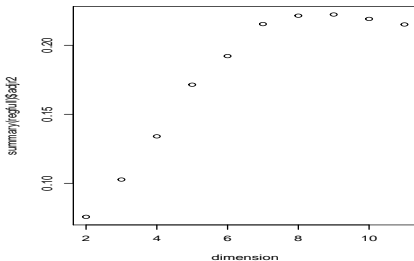
Cp



BIC



Adjusted R2



- **Méthode universelle** : mise en oeuvre dans un cadre statistique général et pour la plupart des procédures d'estimation
- **Principe** : séparer les données d'*apprentissage* et de *test*; construire l'estimateur sur l'échantillon d'apprentissage et utiliser l'échantillon test pour calculer un risque de prédiction (*model assessment*).
Répéter plusieurs fois et moyenner les risques de prédiction obtenus
- Technique très utilisée pour choisir les **tuning parameters** (*model selection*)
En particulier : paramètre de complexité d'un arbre CART (*cp*)

Schémas de validation croisée

- *K-fold CV* : partition des données en K sous-ensembles. Chaque sous-ensemble sert successivement d'échantillon test, le reste d'échantillon d'apprentissage.
En pratique : K entre 5 et 10.
- *Leave-one-out* : n -fold CV
- *Leave- q -out* : chaque sous-ensemble de cardinal q est retiré comme échantillon test, le reste servant d'apprentissage

La procédure complète partitionne en 3 les données :

apprentissage + validation (2/3 données) + test (1/3 données)

- 1 sur les données d'apprentissage : calcul de l'estimateur (algorithme d'apprentissage)
- 2 sur les données de validation : *model selection*, choix d'un meilleur modèle dans la classe de modèles considérée, par sélection de variables et/ou optimisation des hyper-paramètres de l'algorithme (*tuning parameters*).
- 3 sur les données test : calcul de l'erreur de prédiction finale

Les 2 dernières étapes sont souvent réalisées par *validation-croisée*.