


Data Science for Actuaries (ACT6100)

Arthur Charpentier

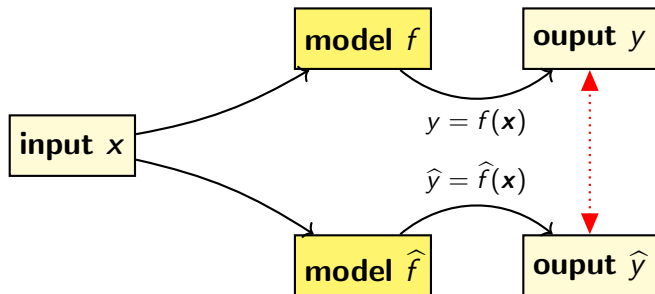
Supervisé # 1.1 (Concepts Fondamentaux)

automne 2Q20

 <https://github.com/freakonometrics/ACT6100/>

Modèle

On suppose qu'il existe une fonction $f : \mathcal{X} \rightarrow \mathcal{Y}$, telle que $y = f(\mathbf{x})$.



On dispose de données $\mathcal{D}_n = \{(y_i, \mathbf{x}_i)\}$ vues comme réalisation de n variables i.i.d. (Y_i, \mathbf{X}_i) .

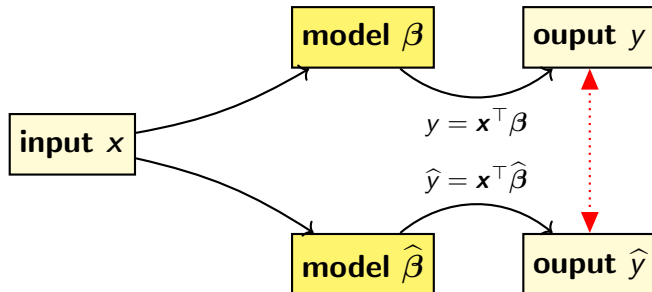
\mathcal{D}_n est une réalisation de $D_n = \{(Y_i, \mathbf{X}_i)\}$

On a associé un **modèle** $\hat{f} = f(\cdot | \mathcal{D}_n) \in \mathcal{M}$ à partir des données \mathcal{D}_n

- ▶ $\hat{y} = f(\mathbf{x} | \mathcal{D}_n)$ est la prévision associée à \mathbf{x}
- ▶ $\hat{Y} = f(\mathbf{x} | \mathcal{D}_n)$ est la prévision vue comme une variable aléatoire

Modèle

On suppose qu'il existe un paramètre β , tel que $y = \mathbf{x}^\top \beta + \epsilon$.



On dispose de données $\mathcal{D}_n = \{(y_i, \mathbf{x}_i)\}$ vues comme réalisation de n variables i.i.d. (Y_i, \mathbf{X}_i) .

\mathcal{D}_n est une réalisation de $D_n = \{(Y_i, \mathbf{X}_i)\}$

On a associé un **modèle** (estimateur) $\hat{\beta}_n$ à partir des données \mathcal{D}_n

- ▶ $\hat{y} = \mathbf{x}^\top \hat{\beta}_n$ est la prévision associée à \mathbf{x}
- ▶ $\hat{Y} = \mathbf{X}^\top \hat{\beta}_n$ est la prévision vue comme une variable aléatoire

Interpréter un modèle

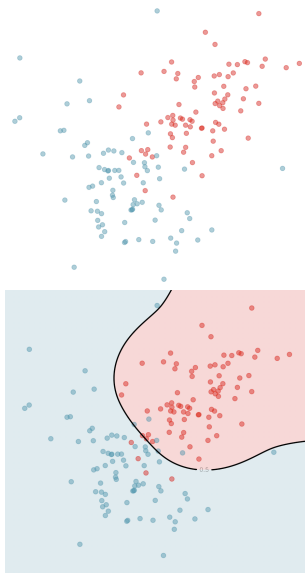
Considérons un problème de classification,
($x_{1,i}, x_{2,i}, y_i$) avec $y_i \in \{0, 1\}$,
et un classifieur $\hat{f} : \mathcal{X} \rightarrow \{0, 1\}$
decision boundary $\mathcal{D}(\mathbf{x})$ —
(région ambiguë)

- ▶ **interprétable par nature**

arbre de classification + modèle linéaire

- ▶ **surgate local** : au voisinage d'un point, on essaye de trouver un modèle linéaire qui approche bien

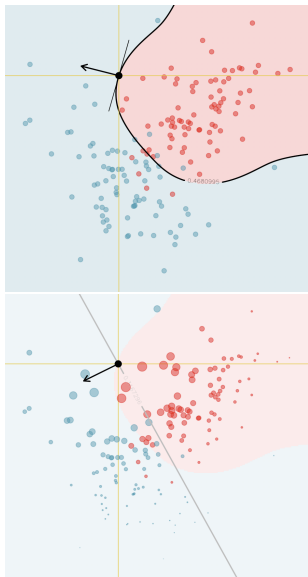
see also Local interpretable model-agnostic explanations (LIME)



Interpréter un modèle

- chercher le plan tangent à $\mathcal{D}(\mathbf{x})$, de vecteur normal $\vec{u}_{(x_1, x_2)}$, interprété comme $\hat{\beta}_{(x_1, x_2)}$
- faire une régression linéaire locale et utiliser $\hat{\beta}_{(x_1, x_2)}$ pour interpréter localement, au voisinage de $\mathbf{x} = (x_1, x_2)$.

- **surgate global** : trouver un modèle linéaire proche / fidèle



Fonction de perte (coût) et risque

- ▶ On définit une **fonction de perte** (ou de **coût**) comme étant une fonction définie sur $\mathcal{Y} \times \mathcal{Y}$, à valeurs réelles, telle que $\ell(y, y') \geq 0$ et $\ell(y, y) = 0$.
- ▶ Ainsi, le **risque d'un prédicteur** \hat{f} pour f (inconnue) est

$$\mathcal{R}(\hat{f}) = \mathbb{E} \left[\ell \left(Y, \hat{f}(\mathbf{X}) \right) \right],$$

- ▶ en régression, on travaille généralement avec une fonction de coût quadratique,

$$\ell(y, y') = (y - y')^2.$$

On obtient alors la fonction de risque

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

où $\hat{f}(\mathbf{X})$ représente une prédiction.

- ▶ On obtient un **risque quadratique**

Risque empirique

- ▶ La distribution exacte de Y n'est pas connue, il n'est donc pas possible de calculer la fonction de risque. On construit la **fonction de risque empirique**

$$\widehat{\mathcal{R}}_n(\widehat{f}) = \frac{1}{n} \sum_{i=1}^n \ell \left(\widehat{f}(\mathbf{x}_i), y_i \right),$$

où $\{\mathbf{x}_i, y_i\}_{i=1, \dots, n}$ est un échantillon aléatoire

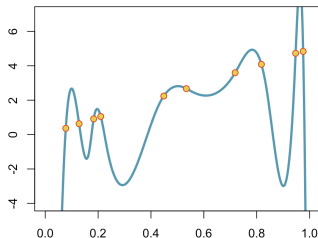
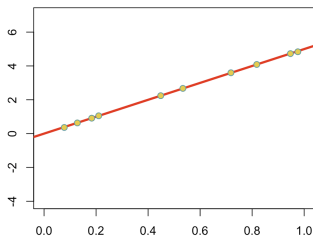
- ▶ Dans un contexte de régression, on travaille généralement avec l'**erreur quadratique moyenne**, ou *Mean Squared Error* (MSE), donnée par

$$\text{EQM}_n = \text{MSE}_n = \frac{1}{n} \sum_{i=1}^n \left(y_i - \widehat{f}(\mathbf{x}_i) \right)^2.$$

Risque empirique

- sur les données suivantes $\{(x_1, y_1), \dots, (x_n, y_n)\}$

$$\hat{\mathcal{R}}_n(\hat{f}) = \frac{1}{n} \sum_{i=1}^n \ell(\hat{f}(\mathbf{x}_i), y_i) = 0 \quad \text{mais aussi } \hat{f}$$



Notion de généralisation...

Base d'entraînement

- Dans la formule

$$\widehat{\mathcal{R}}^{\text{IS}}_n(\widehat{f}_n) = \frac{1}{n} \sum_{i=1}^n \ell(\widehat{f}_n(\mathbf{x}_i), y_i) = 0$$

l'échantillon utilisé pour calculer l'erreur quadratique moyenne est le même que celui utilisé pour ajuster le modèle: erreur d'entraînement, **in-sample risk**

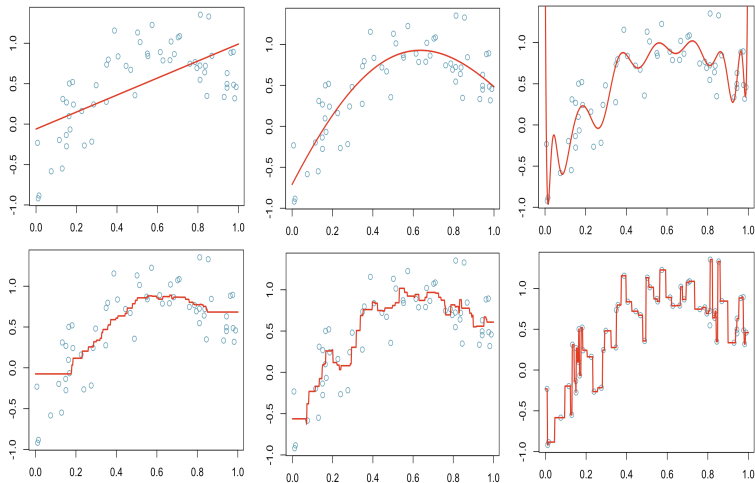
- si on construit

$$\widehat{f}_n = \operatorname{argmin}_{f \in \mathcal{M}} \left\{ \widehat{\mathcal{R}}^{\text{IS}}_n(f) \right\}$$

on aura tendance à capturer beaucoup de bruit et à sur-ajuster les données : sur-apprentissage

Base d'entraînement

\mathcal{M}_k : régression polynomiale de degré k ou k -nearest neighbors



Base d'entraînement et base de validation

- ▶ Pour éviter ce problème de surapprentissage, on va diviser aléatoirement la base de données initiale en une base d'entraînement et une base de validation.
- ▶ La base d'entraînement, de taille $n_T < n$, sera utilisée pour estimer les paramètres du modèle,

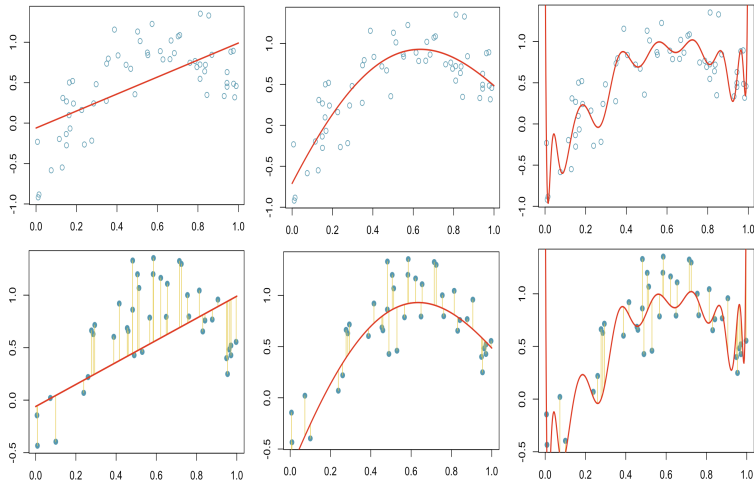
$$\hat{f}_{n_T} = \operatorname{argmin}_{f \in \mathcal{M}} \left\{ \hat{\mathcal{R}}_{n_T}^{IS}(f) \right\}$$

- ▶ La base de validation, de taille $n_V = n - n_T$, sera utilisée pour sélectionner le modèle en minimisant une erreur quadratique moyenne de validation, out-of-sample risk

$$\hat{\mathcal{R}}_{n_V}^{OS}(\hat{f}_{n_T}) = \frac{1}{n} \sum_{i=1}^{n_V} \ell \left(\hat{f}_{n_T}(\mathbf{x}_i), y_i \right) = 0$$

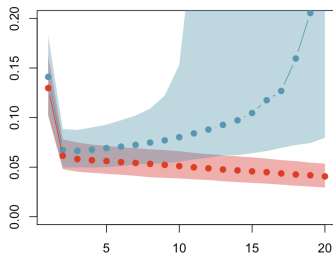
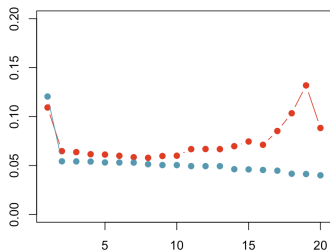
Base d'entraînement et de validation

\mathcal{M}_k : régression polynomiale de degré k



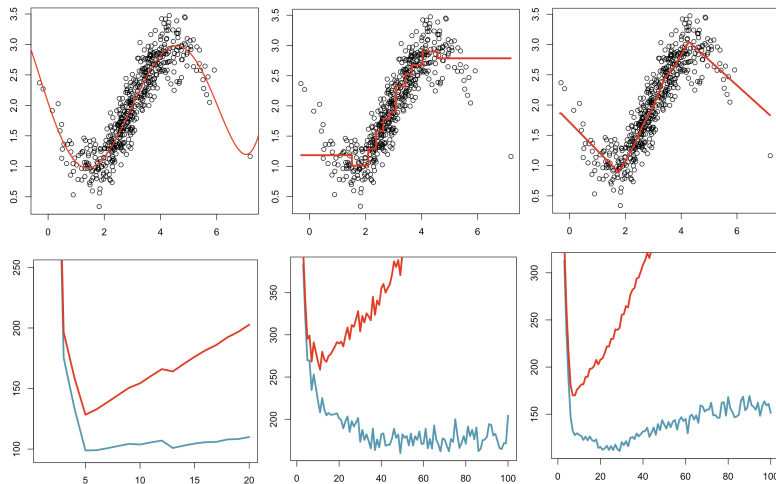
Base d'entrainement et base de validation

```
1 set.seed(1)
2 n=100
3 x=runif(n)
4 y=sin((x-1/6)*pi)+rnorm(n)/4
5 base=data.frame(x=x,y=y)
6 train=1:60
7 test=61:n
8 EQM = function(k){
9   reg=lm(y~poly(x,k),data=base[
10     train,])
11   base$yp=predict(reg,newdata=
12     base)
13   eqm_v=sum((base[test,"y"]-base
14     [test,"yp"])^2)
15   eqm_t=sum((base[train,"y"]-
16     base[train,"yp"])^2)
17   c(eqm_t/60,eqm_v/40)}
18 VE=Vectorize(EQM)(1:20)
```



Base d'entraînement et base de validation

Régression polynomiale, constante par morceaux, et linéaire par morceaux (splines)



Décomposition de l'erreur quadratique moyenne

On a

$$\begin{aligned}\mathbb{E}[(Y - \hat{f}(\mathbf{X}))^2] &= \mathbb{E}[(f(\mathbf{X}) + \epsilon - \hat{f}(\mathbf{X}))^2] \\&= \mathbb{E}[(f(\mathbf{X}) - \hat{f}(\mathbf{X}))^2] + \mathbb{E}[\epsilon^2] \\&\quad + 2\mathbb{E}[\epsilon (f(\mathbf{X}) - \hat{f}(\mathbf{X}))] \\&= \mathbb{E}[(f(\mathbf{X}) - \hat{f}(\mathbf{X}))^2] + \mathbb{E}[\epsilon^2] \\&= \mathbb{E}[(f(\mathbf{X}) - \hat{f}(\mathbf{X}))^2] + \text{Var}[\epsilon] \\&= \text{Var}[(f(\mathbf{X}) - \hat{f}(\mathbf{X}))] + \mathbb{E}[(f(\mathbf{X}) - \hat{f}(\mathbf{X}))]^2 + \text{Var}[\epsilon] \\&= \underbrace{\text{Var}[\hat{f}(\mathbf{X})]}_{(1)} + \underbrace{\mathbb{E}[(f(\mathbf{X}) - \hat{f}(\mathbf{X}))]^2}_{(2)} + \underbrace{\text{Var}[\epsilon]}_{(3)}.\end{aligned}$$

Décomposition

- ▶ On obtient ainsi une erreur quadratique moyenne qui comprend trois termes: la variance de l'estimateur, le biais au carré et l'erreur stochastique (irréductible).
- ▶ Un *bon* modèle doit avoir simultanément une variance faible et un biais faible.
- ▶ Un modèle trop flexible n'est pas approprié car sa variance est généralement trop élevée: il est trop sensible au bruit présent dans les données.
- ▶ Un modèle peu flexible n'est pas approprié car son biais est généralement trop élevé: il n'arrive pas à capturer le comportement de la fonction f inconnue.
- ▶ compromis entre variance et biais

Un peu de formalisme

De manière générale, étant donnée une fonction de perte $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}_+$, le risque est

$$\mathbb{E}_{\mathbf{X}} \left[\mathbb{E}_{D_n} \left(\mathbb{E}_{Y|\mathbf{X}} \left[\ell(Y, \hat{f}(\mathbf{X}) | D_n) \right] \right) \right]$$

qu'on ne peut pas calculer sans connaître la loi de (Y, \mathbf{X}) .
Si ℓ est la perte quadratique $\ell(y, \hat{y}) = (y - \hat{y})^2$,

$$\begin{aligned} \mathcal{R}(\hat{f}) &= \mathbb{E}_{D_n} \left(\mathbb{E}_{Y|\mathbf{X}} \left[\ell(Y, \hat{f}(\mathbf{X}) | D_n) \right] \right) \\ &= \left(\mathbb{E}_{Y|\mathbf{X}}(Y) - \mathbb{E}_{D_n}[\hat{f}(\mathbf{X}) | D_n] \right)^2 \\ &\quad + \mathbb{E}_{Y|\mathbf{X}} \left[(Y - \mathbb{E}_{Y|\mathbf{X}}(Y))^2 \right] \\ &\quad + \mathbb{E}_{D_n} \left[(\hat{f}(\mathbf{X}) | D_n - \mathbb{E}_{D_n}[\hat{f}(\mathbf{X}) | D_n])^2 \right] \end{aligned}$$

On reconnaît le bias^2 , l'erreur stochastique, et la variance de l'estimateur