

Data Science for Actuaries (ACT6100)

Arthur Charpentier

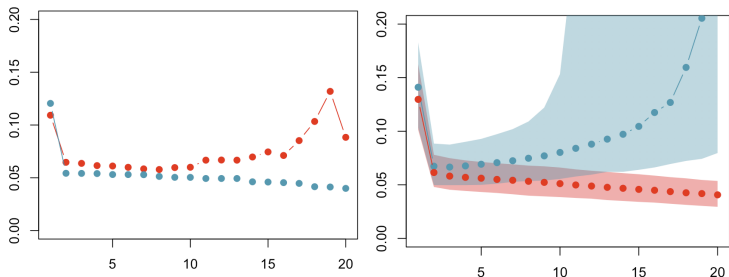
Supervisé # 1.2 (Concepts Fondamentaux - 2)

automne 2020

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

Base d'entrainement et base de validation

- Tel que mentionné précédemment, il est possible de diviser la base de données initiale en une base d'entrainement ($\sim 70\%$) et une base de validation ($\sim 30\%$).



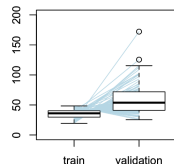
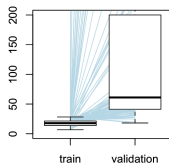
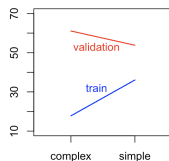
Base d'entraînement et base de validation

Classical Approach : split the sample \mathcal{D}_n in two parts

Hold-Out Cross Validation

1. Split $\{1, 2, \dots, n\}$ in T (training) and V (validation)
2. Estimate \hat{m} on sample (y_i, \mathbf{x}_i) , $i \in T$: \hat{m}_T
3. Compute $\frac{1}{|V|} \sum_{i \in V} \ell(y_i, \hat{m}_T(x_{1,i}, \dots, x_{p,i}))$

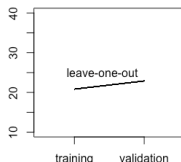
```
1 chicago = read.table("http://
    freakonometrics.free.fr/chicago.txt",
    header=TRUE, sep=";")
2 idx = sample(1:nrow(chicago), nrow(chicago)
    *.7)
3 train = chicago[idx,]
4 valid = chicago[-idx,]
```



Validation croisée, Leave-One-Out

Leave-one-Out Cross Validation

1. Estimate n models : \hat{m}_{-j} on sample (y_i, \mathbf{x}_i) , $i \in \{1, \dots, n\} \setminus \{j\}$
2. Compute $\frac{1}{n} \sum_{i=1}^n \ell(y_i, \hat{m}_{-i}(\mathbf{x}_i))$



Can be computationally intensive...

In the case of a linear regression, there is a simple formula to compute $\hat{\beta}_{-j}$ when observation j is removed. Let

$$\mathbf{H} = \mathcal{P}_{V(\mathbf{X})} = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$$

$$(\mathbf{X}_{(j)}^\top \mathbf{X}_{(j)})^{-1} \mathbf{X}_{(j)}^\top = \frac{1}{1 - H_{j,j}} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}_{(j)}^\top$$



Validation croisée, k -fold

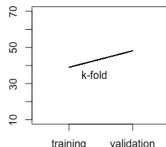
Instead of fitting n models, fit only K

K -Fold Cross Validation

1. Split $\{1, 2, \dots, n\}$ in K groups V_1, \dots, V_K
2. Estimate K models : \hat{m}_k on sample (y_i, \mathbf{x}_i) , $i \in \{1, \dots, n\} \setminus V_k$
3. Compute
$$\frac{1}{K} \sum_{k=1}^K \frac{1}{|V_k|} \sum_{i \in V_k} \ell(y_i, \hat{m}_k(x_{1,i}, \dots, x_{p,i}))$$

If $K = 10$ we fit on 90% of the observation, and validate on the remaining 10%

Here the K groups should be created randomly...



k-validation croisée

- ▶ En pratique, on pose généralement $k = 10$ (*ten-fold cross validation*) ou $k = n$ (*leave-one-out cross validation*).
- ▶ Lorsque le modèle optimal est sélectionné, on cherche à estimer l'erreur quadratique de prédiction, c'est-à-dire

$$\mathbb{E}\left(Y^* - \hat{f}(\mathbf{X}^*)\right)^2,$$

où $(Y^*; \mathbf{X}^*)$ est une nouvelle observation.

- ▶ Il faut alors faire attention car le tMSE généralement sous-estime cette valeur.
- ▶ Pour contourner ce problème, si la base de données est de taille suffisante, on peut garder une portion des données pour constituer une base de test.
- ▶ On estime alors l'erreur quadratique de prédiction comme étant

$$\widehat{\mathbb{E}\left(Y^* - \hat{f}(\mathbf{X}^*)\right)^2} = \frac{1}{n_{\text{test}}} \sum_{i=1}^{n_{\text{test}}} \left(y_i^* - \hat{f}(\mathbf{x}_i^*)\right)^2.$$

Validation croisée

Bootstrap Cross Validation

1. Generate B bootstrap samples from $\{1, 2, \dots, n\}$, l_1, \dots, l_B
2. Estimate B models : \hat{m}_b on sample (y_i, \mathbf{x}_i) , $i \in l_b$
3. Compute
$$\frac{1}{B} \sum_{b=1}^B \frac{1}{n - |l_b|} \sum_{i \notin l_b} \ell(y_i, \hat{m}_b(x_{1,i}, \dots, x_{p,i}))$$

The probability that $i \notin l_b$ is

$$\left(1 - \frac{1}{n}\right)^n \sim e^{-1} (= 36.78\%)$$

At stage b , we validate on $\sim 36.78\%$ of the dataset.



Cas des séries chronologiques

Time Series

A time series is a sequence of observations (y_t) ordered in time.

Write $y_t = s_t + u_t$, with systematic part s_t (signal / trend) and 'residual' term u_t (u_t) is supposed to be a strictly stationary time series (s_t) might be a 'linear' trend, plus a seasonal cycle

Buys-Ballot (1847, Les changements périodiques de température, dépendants de la nature du soleil et de la lune, mis en rapport avec le pronostic du temps, déduits d'observations néerlandaises de 1729 à 1846) - original probably in Dutch.

TABLEAU REPRÉSENTANT LA MARCHÉ DE LA TEMPÉRATURE PENDANT L'ANNÉE.

Date.	Temp.	Temp.	Diffé.	Temp.	Diffé.	Date.	Temp.	Temp.	Diffé.	Temp.	Diffé.
		calculée.		calculée.				calculée.		calculée.	
10 Janv.	23.50	23.53	0	35.00	0	17 Juill.	0.50	+ 0.24	+ 0.24	93.00	+ 0.00
15 "	+ 1.51	+ 0.78	0	0.23	+ 0.80	0	+ 0.12	0	- 0.50	64.33	+ 0.00
20 "	0.00	0.01	+ 0.01	- 0.01	- 0.01	0	0.00	0.01	+ 0.01	64.37	+ 0.01
25 "	0.00	0.01	+ 0.01	- 0.01	- 0.01	0	0.00	0.01	+ 0.01	64.31	+ 0.01
30 "	0.00	0.01	+ 0.01	- 0.01	- 0.01	0	0.00	0.01	+ 0.01	64.31	+ 0.01
4 Fév.	0.00	0.01	+ 0.01	- 0.01	- 0.01	0	0.00	0.01	+ 0.01	64.31	+ 0.01
9 "	0.00	0.01	+ 0.01	- 0.01	- 0.01	0	0.00	0.01	+ 0.01	64.31	+ 0.01
14 "	0.00	0.01	+ 0.01	- 0.01	- 0.01	0	0.00	0.01	+ 0.01	64.31	+ 0.01
19 "	0.00	0.01	+ 0.01	- 0.01	- 0.01	0	0.00	0.01	+ 0.01	64.31	+ 0.01
24 "	0.00	0.01	+ 0.01	- 0.01	- 0.01	0	0.00	0.01	+ 0.01	64.31	+ 0.01
1 Mars	0.00	0.01	+ 0.01	- 0.01	- 0.01	0	0.00	0.01	+ 0.01	64.31	+ 0.01
6 "	0.00	0.01	+ 0.01	- 0.01	- 0.01	0	0.00	0.01	+ 0.01	64.31	+ 0.01
11 "	0.00	0.01	+ 0.01	- 0.01	- 0.01	0	0.00	0.01	+ 0.01	64.31	+ 0.01
16 "	0.00	0.01	+ 0.01	- 0.01	- 0.01	0	0.00	0.01	+ 0.01	64.31	+ 0.01
21 "	0.00	0.01	+ 0.01	- 0.01	- 0.01	0	0.00	0.01	+ 0.01	64.31	+ 0.01
26 "	0.00	0.01	+ 0.01	- 0.01	- 0.01	0	0.00	0.01	+ 0.01	64.31	+ 0.01
31 Mars	0.00	0.01	+ 0.01	- 0.01	- 0.01	0	0.00	0.01	+ 0.01	64.31	+ 0.01
5 Avril	0.00	0.01	+ 0.01	- 0.01	- 0.01	0	0.00	0.01	+ 0.01	64.31	+ 0.01
10 "	0.00	0.01	+ 0.01	- 0.01	- 0.01	0	0.00	0.01	+ 0.01	64.31	+ 0.01
15 "	0.00	0.01	+ 0.01	- 0.01	- 0.01	0	0.00	0.01	+ 0.01	64.31	+ 0.01
20 "	0.00	0.01	+ 0.01	- 0.01	- 0.01	0	0.00	0.01	+ 0.01	64.31	+ 0.01
25 "	0.00	0.01	+ 0.01	- 0.01	- 0.01	0	0.00	0.01	+ 0.01	64.31	+ 0.01
30 "	0.00	0.01	+ 0.01	- 0.01	- 0.01	0	0.00	0.01	+ 0.01	64.31	+ 0.01
4 Mai	0.00	0.01	+ 0.01	- 0.01	- 0.01	0	0.00	0.01	+ 0.01	64.31	+ 0.01
9 "	0.00	0.01	+ 0.01	- 0.01	- 0.01	0	0.00	0.01	+ 0.01	64.31	+ 0.01
14 "	0.00	0.01	+ 0.01	- 0.01	- 0.01	0	0.00	0.01	+ 0.01	64.31	+ 0.01
19 "	0.00	0.01	+ 0.01	- 0.01	- 0.01	0	0.00	0.01	+ 0.01	64.31	+ 0.01
24 "	0.00	0.01	+ 0.01	- 0.01	- 0.01	0	0.00	0.01	+ 0.01	64.31	+ 0.01
29 "	0.00	0.01	+ 0.01	- 0.01	- 0.01	0	0.00	0.01	+ 0.01	64.31	+ 0.01
31 Mai	0.00	0.01	+ 0.01	- 0.01	- 0.01	0	0.00	0.01	+ 0.01	64.31	+ 0.01
5 Juin	0.00	0.01	+ 0.01	- 0.01	- 0.01	0	0.00	0.01	+ 0.01	64.31	+ 0.01
10 "	0.00	0.01	+ 0.01	- 0.01	- 0.01	0	0.00	0.01	+ 0.01	64.31	+ 0.01
15 "	0.00	0.01	+ 0.01	- 0.01	- 0.01	0	0.00	0.01	+ 0.01	64.31	+ 0.01
20 "	0.00	0.01	+ 0.01	- 0.01	- 0.01	0	0.00	0.01	+ 0.01	64.31	+ 0.01
25 "	0.00	0.01	+ 0.01	- 0.01	- 0.01	0	0.00	0.01	+ 0.01	64.31	+ 0.01
30 "	0.00	0.01	+ 0.01	- 0.01	- 0.01	0	0.00	0.01	+ 0.01	64.31	+ 0.01
31 Juin	0.00	0.01	+ 0.01	- 0.01	- 0.01	0	0.00	0.01	+ 0.01	64.31	+ 0.01
5 Juil.	0.00	0.01	+ 0.01	- 0.01	- 0.01	0	0.00	0.01	+ 0.01	64.31	+ 0.01
10 "	0.00	0.01	+ 0.01	- 0.01	- 0.01	0	0.00	0.01	+ 0.01	64.31	+ 0.01
15 "	0.00	0.01	+ 0.01	- 0.01	- 0.01	0	0.00	0.01	+ 0.01	64.31	+ 0.01

Cas des séries chronologiques: Lissage Exponentiel

Exponential smoothing - Simple

From time series (y_t) define a smooth version

$$s_t = \alpha \cdot y_t + (1 - \alpha) \cdot s_{t-1} = s_{t-1} + \alpha \cdot (y_t - s_{t-1})$$

for some $\alpha \in (0, 1)$ and starting point $s_0 = y_1$ Forecast is ${}_t\hat{y}_{t+h} = s_t$

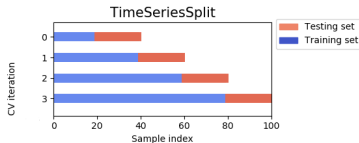
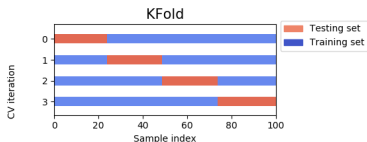
It is called **exponential smoothing** since

$$\begin{aligned} s_t &= \alpha y_t + (1 - \alpha) s_{t-1} \\ &= \alpha y_t + \alpha(1 - \alpha) y_{t-1} + (1 - \alpha)^2 s_{t-2} \\ &= \alpha [y_t + (1 - \alpha) y_{t-1} + (1 - \alpha)^2 y_{t-2} + \cdots + (1 - \alpha)^{t-1} y_1] + (1 - \alpha) \end{aligned}$$

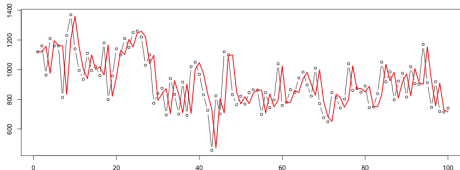
corresponding to **exponentially weighted moving average**

Need to adapt cross-validation techniques,

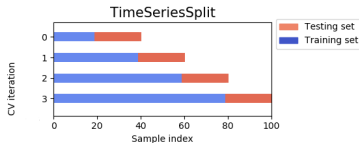
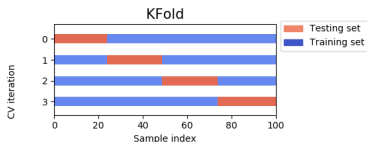
Cas des séries chronologiques: Lissage Exponentiel



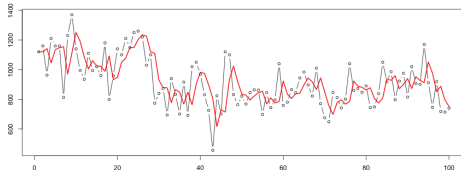
Optimal α ? $\alpha^* \in \operatorname{argmin} \left\{ \sum_{t=2}^T \ell_2(y_t - {}_{t-1}\hat{y}_t) \right\}$ (leave-one-out strategy)



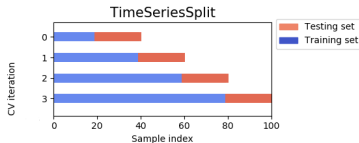
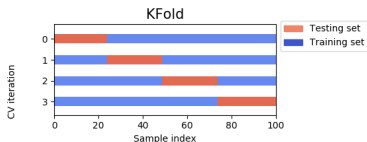
Cas des séries chronologiques: Lissage Exponentiel



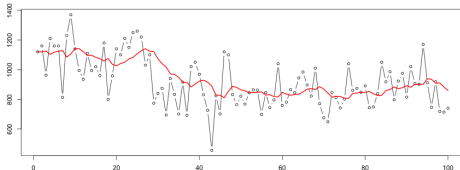
Optimal α ? $\alpha^* \in \operatorname{argmin} \left\{ \sum_{t=2}^T \ell_2(y_t - t_{-1} \hat{y}_t) \right\}$ (leave-one-out strategy)



Cas des séries chronologiques: Lissage Exponentiel



Optimal α ? $\alpha^* \in \operatorname{argmin} \left\{ \sum_{t=2}^T \ell_2(y_t - t_{-1}\hat{y}_t) \right\}$ (leave-one-out strategy)



Cas des séries chronologiques: Lissage Exponentiel

See Hyndman *et al.* (2008, [Forecasting with Exponential Smoothing](#))

Exponential smoothing - Double

From time series (y_t) define a smooth version

$$\begin{cases} s_t = \alpha y_t + (1 - \alpha)(s_{t-1} + b_{t-1}) & s_t = \alpha y_t + (1 - \alpha)(s_{t-1} + b_{t-1}) \\ b_t = \beta(s_t - s_{t-1}) + (1 - \beta)b_{t-1} & b_t = \beta(s_t - s_{t-1}) + (1 - \beta)b_{t-1} \end{cases}$$

for some $\alpha \in (0, 1)$, some trend $\beta \in (0, 1)$ and starting points s_0 and b_0 , $s_0 = y_0$ and $b_0 = y_1 - y_0$. Forecast is ${}_t\hat{y}_{t+h} = s_t + h \cdot b_t$.

Cas des séries chronologiques: Lissage Exponentiel

See Hyndman *et al.* (2008, [Forecasting with Exponential Smoothing](#))

Exponential smoothing - Double

From time series (y_t) define a smooth version

$$\begin{cases} s_t = \alpha y_t + (1 - \alpha)(s_{t-1} + b_{t-1}) & s_t = \alpha y_t + (1 - \alpha)(s_{t-1} + b_{t-1}) \\ b_t = \beta(s_t - s_{t-1}) + (1 - \beta)b_{t-1} & b_t = \beta(s_t - s_{t-1}) + (1 - \beta)b_{t-1} \end{cases}$$

for some $\alpha \in (0, 1)$, some trend $\beta \in (0, 1)$ and starting points s_0 and b_0 , $s_0 = y_0$ and $b_0 = y_1 - y_0$. Forecast is ${}_t\hat{y}_{t+h} = s_t + h \cdot b_t$.

Cas des séries chronologiques: Lissage Exponentiel

See Hyndman *et al.* (2008, [Forecasting with Exponential Smoothing](#))

Exponential smoothing - Double

From time series (y_t) define a smooth version

$$\begin{cases} s_t = \alpha y_t + (1 - \alpha)(s_{t-1} + b_{t-1}) & s_t = \alpha y_t + (1 - \alpha)(s_{t-1} + b_{t-1}) \\ b_t = \beta(s_t - s_{t-1}) + (1 - \beta)b_{t-1} & b_t = \beta(s_t - s_{t-1}) + (1 - \beta)b_{t-1} \end{cases}$$

for some $\alpha \in (0, 1)$, some trend $\beta \in (0, 1)$ and starting points s_0 and b_0 , $s_0 = y_0$ and $b_0 = y_1 - y_0$. Forecast is ${}_t\hat{y}_{t+h} = s_t + h \cdot b_t$.

Cas des séries chronologiques: Lissage Exponentiel

Exponential smoothing - Seasonal with lag L (Holt-Winters)

From time series (y_t) define a smooth version

$$\begin{cases} s_t = \alpha \frac{X_t}{c_{t-L}} + (1 - \alpha)(s_{t-1} + b_{t-1}) \\ b_t = \beta(s_t - s_{t-1}) + (1 - \beta)b_{t-1} \\ c_t = \gamma \frac{y_t}{s_t} + (1 - \gamma)c_{t-L} \end{cases}$$

for some $\alpha \in (0, 1)$, some trend $\beta \in (0, 1)$, some seasonal change smoothing factor, $\gamma \in (0, 1)$ and starting points $s_0 = y_0$. Forecast is ${}_t\hat{y}_{t+h} = (s_t + hb_t)c_{t-L+1+(h-1) \bmod L}$.

```
1 stats::HoltWinters()
```