

Méthodes d'apprentissage
IFT603-712

Théorie de la décision
Par
Pierre-Marc Jodoin

Régression linéaire

RAPPEL

- Le modèle de **régression linéaire** est le suivant :

$$y_{\vec{w}}(\vec{x}) = w_0 + w_1x_1 + w_2x_2 + \dots + w_dx_d$$

$$\text{où } \vec{x} = (x_1, x_2, \dots, x_d)^T$$

- La prédiction correspond donc à

- Une **droite** pour $d=1$
- Un **plan** pour $d=2$
- Un **hyperplan** pour $d>2$

2

Régression linéaire

RAPPEL

$$y_{\vec{w}}(\vec{x}) = w_0 + w_1x_1 + w_2x_2 + \dots + w_dx_d$$

$$y_{\vec{w}}(\vec{x}) = \vec{w}^T \vec{x}$$

$$\text{où } \vec{x}' = (1, x_1, x_2, \dots, x_d)^T$$

3

Problème à résoudre

RAPPEL

$$\bar{w} = \arg \min_{\bar{w}} \sum_{n=1}^N (y_{\bar{w}}(\bar{x}_n) - t_n)^2$$

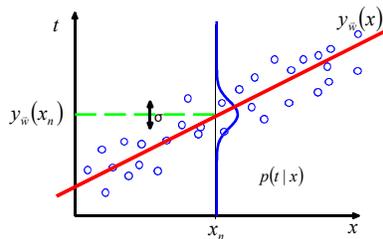
Il est très bien connu en technique d'apprentissage que cette solution est optimale lorsque le **bruit est gaussien**.



Formulation probabiliste

Loi conditionnelle (*maximum de vraisemblance*)

RAPPEL



Formulation probabiliste

RAPPEL

Pour entraîner le modèle $y_w(\bar{x})$ nous passerons par une formulation probabiliste :

$$p(t | \bar{x}, \bar{w}, \Sigma) = N(t | y_w(\bar{x}), \Sigma)$$



➤ Revient à supposer que les **cibles** sont des **versions bruitées** du vrai modèle

$$t_n = y_w(\bar{x}_n) + \varepsilon$$

Bruit gaussien de moyenne 0 et de variance σ^2

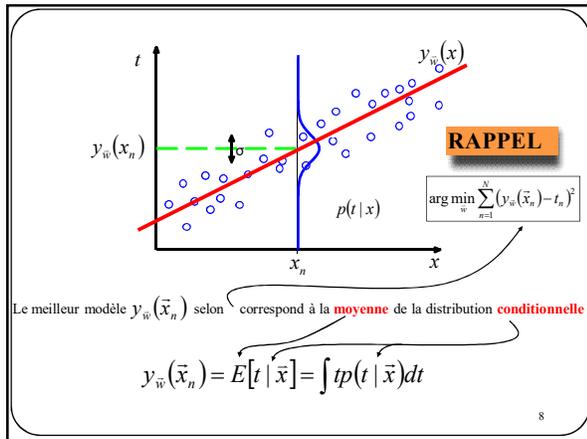
Maximum de vraisemblance

RAPPEL

$$\bar{w} = \arg \min_{\bar{w}} \sum_{n=1}^N (y_{\bar{w}}(\bar{x}_n) - t_n)^2$$

$$\bar{w}_{MV} = (X^T X)^{-1} X^T T$$

Et bien sûr, on peut utiliser une **fonction de base** pour rendre le modèle non linéaire.



Comment prouver cette affirmation?

$$\arg \min_y \sum_{n=1}^N (y(x_n) - t_n)^2$$

Meilleure solution

$$y(x) = E[t | x] = \int tp(t | x) dt$$

Preuve 1 (1.5.5 Bishop)

$$\arg \min_y \underbrace{\sum_{n=1}^N (y(\bar{x}_n) - t_n)^2}_{L(y(\bar{x}), t_n)} = \arg \min_y \underbrace{\frac{1}{N} \sum_{n=1}^N L(y(\bar{x}_n), t_n)}_{\substack{\text{Erreur moyenne:} \\ N \rightarrow \infty, E[L]}}$$

puisque (\bar{x}, t) est i.i.d de $p(\bar{x}, t)$

$$\begin{aligned} E[L] &= \iint L p(\bar{x}, t) dx dt \\ &= \iint (y(\bar{x}) - t)^2 p(\bar{x}, t) dx dt \end{aligned}$$

10

Preuve 1

$$\arg \min_y \underbrace{\iint (y(\bar{x}) - t)^2 p(\bar{x}, t) dx dt}_{E[L]}$$

$$\frac{\delta E[L]}{\delta y} = 0 \Rightarrow 2 \int (y(\bar{x}) - t) p(\bar{x}, t) dt = 0$$



11

Preuve 1

$$\int (y(\bar{x}) - t) p(\bar{x}, t) dt = 0$$

$$\int (y(\bar{x}) p(\bar{x}, t) - t p(\bar{x}, t)) dt = 0$$

$$\int y(\bar{x}) p(\bar{x}, t) dt - \int t p(\bar{x}, t) dt = 0$$

$$y(\bar{x}) \int p(\bar{x}, t) dt - \int t p(\bar{x}, t) dt = 0$$

Marginalisation de t

$$y(\bar{x}) p(\bar{x}) - \int t p(\bar{x}, t) dt = 0$$

$$y(\bar{x}) p(\bar{x}) - \int t p(t | \bar{x}) p(\bar{x}) dt = 0 \quad \text{car } p(x, t) = p(t | x) p(x)$$

Preuve 1

$$y(\bar{x})p(\bar{x}) - \int tp(t|\bar{x})p(\bar{x})dt = 0$$

$$y(\bar{x})p(\bar{x}) - \underbrace{p(\bar{x}) \int tp(t|\bar{x})dt}_{\text{Expérance mathématique conditionnelle}} = 0$$

$$y(\bar{x}) - E[t|\bar{x}] = 0$$

$$\boxed{y(\bar{x}) = E[t|\bar{x}]}$$

Preuve 2 (1.5.5 Bishop)

$$L = (y(\bar{x}) - t)^2 = (y(\bar{x}) - E[t|\bar{x}] + E[t|\bar{x}] - t)^2$$

$$\begin{aligned} &= (y(\bar{x}) - E[t|\bar{x}])^2 \\ &\quad + 2(y(\bar{x}) - E[t|\bar{x}])(E[t|\bar{x}] - t) \\ &\quad + (E[t|\bar{x}] - t)^2 \end{aligned}$$

14

Preuve 2

$$E[L] = \iint (y(\bar{x}) - t)^2 p(\bar{x}, t) dx dt$$

$$= \iint (y(\bar{x}) - E[t|\bar{x}])^2 p(\bar{x}, t) dx dt$$

$$+ \iint 2(y(\bar{x}) - E[t|\bar{x}])(E[t|\bar{x}] - t) p(\bar{x}, t) dx dt$$

$$+ \iint (E[t|\bar{x}] - t)^2 p(\bar{x}, t) dx dt$$

$$\iint 2(y(\bar{x}) - E[t|\bar{x}])(E[t|\bar{x}] - t) p(\bar{x}, t) dx dt$$

15

Preuve 2

$$\begin{aligned}
 E[L] &= \iint (y(\bar{x}) - t)^2 p(x, t) dx dt \\
 &= \iint (y(\bar{x}) - E[t | \bar{x}])^2 p(\bar{x}, t) dx dt \\
 &\quad + \iint 2(y(\bar{x}) - E[t | \bar{x}])(E[t | \bar{x}] - t) p(\bar{x}, t) dx dt \\
 &\quad + \iint (E[t | \bar{x}] - t)^2 p(\bar{x}, t) dx dt \\
 &= \iint 2(y(\bar{x}) - E[t | \bar{x}])(E[t | \bar{x}] - t) p(t | \bar{x}) p(\bar{x}) dt dx
 \end{aligned}$$

16

Preuve 2

$$\begin{aligned}
 E[L] &= \iint (y(\bar{x}) - t)^2 p(x, t) dx dt \\
 &= \iint (y(\bar{x}) - E[t | \bar{x}])^2 p(\bar{x}, t) dx dt \\
 &\quad + \iint 2(y(\bar{x}) - E[t | \bar{x}])(E[t | \bar{x}] - t) p(\bar{x}, t) dx dt \\
 &\quad + \iint (E[t | \bar{x}] - t)^2 p(\bar{x}, t) dx dt \\
 &= \int 2(y(\bar{x}) - E[t | \bar{x}]) p(\bar{x}) \left\{ \int (E[t | \bar{x}] - t) p(t | \bar{x}) dt \right\} dx
 \end{aligned}$$

17

Preuve 2

$$\begin{aligned}
 E[L] &= \iint (y(\bar{x}) - t)^2 p(\bar{x}, t) dx dt \\
 &= \iint (y(\bar{x}) - E[t | \bar{x}])^2 p(\bar{x}, t) dx dt \\
 &\quad + \iint 2(y(\bar{x}) - E[t | \bar{x}])(E[t | \bar{x}] - t) p(\bar{x}, t) dx dt \\
 &\quad + \iint (E[t | \bar{x}] - t)^2 p(\bar{x}, t) dx dt \\
 &= \int 2(y(\bar{x}) - E[t | \bar{x}]) p(\bar{x}) \left\{ \int E[t | \bar{x}] p(t | \bar{x}) dt - \int t p(t | \bar{x}) dt \right\} dx
 \end{aligned}$$

18

Preuve 2

$$\begin{aligned}
 E[L] &= \iint (y(\bar{x}) - t)^2 p(x, t) dx dt \\
 &= \iint (y(\bar{x}) - E[t | \bar{x}])^2 p(\bar{x}, t) dx dt \\
 &\quad + \iint 2(y(\bar{x}) - E[t | \bar{x}])(E[t | \bar{x}] - t) p(\bar{x}, t) dx dt \\
 &\quad + \iint (E[t | \bar{x}] - t)^2 p(\bar{x}, t) dx dt \\
 &= \int 2(y(\bar{x}) - E[t | \bar{x}]) p(\bar{x}) \left\{ E[t | \bar{x}] \int p(t | \bar{x}) dt - E[t | \bar{x}] \right\} dx
 \end{aligned}$$

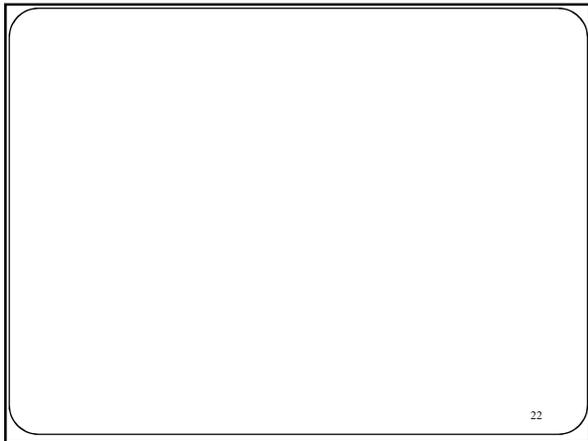
Preuve 2

$$\begin{aligned}
 E[L] &= \iint (y(\bar{x}) - t)^2 p(x, t) dx dt \\
 &= \iint (y(\bar{x}) - E[t | \bar{x}])^2 p(\bar{x}, t) dx dt \\
 &\quad + \iint 2(y(\bar{x}) - E[t | \bar{x}])(E[t | \bar{x}] - t) p(\bar{x}, t) dx dt \\
 &\quad + \iint (E[t | \bar{x}] - t)^2 p(\bar{x}, t) dx dt \\
 &= \int 2(y(\bar{x}) - E[t | \bar{x}]) p(\bar{x}) \left\{ E[t | \bar{x}] - E[t | \bar{x}] \right\} dx
 \end{aligned}$$

Preuve 2

$$\begin{aligned}
 E[L] &= \iint (y(\bar{x}) - t)^2 p(\bar{x}, t) dx dt \\
 &= \iint (y(\bar{x}) - E[t | \bar{x}])^2 p(\bar{x}, t) dx dt \\
 &\quad + \iint (E[t | \bar{x}] - t)^2 p(\bar{x}, t) dx dt
 \end{aligned}$$

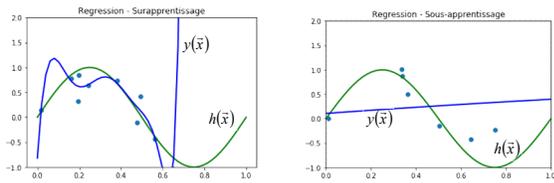
$E[L]$ est minimum lorsque $y(\bar{x}) = E[t | \bar{x}]$



22

Décomposition biais-variance (Bishop 3.2)

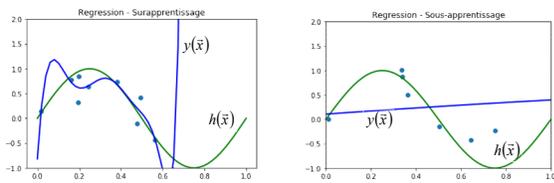
On a vu les concepts de « sur-apprentissage » et de « sous-apprentissage ».



$y(\vec{x})$: modèle trouvé par apprentissage
 $h(\vec{x})$: le meilleur modèle représentant les données

Décomposition biais-variance (Bishop 3.2)

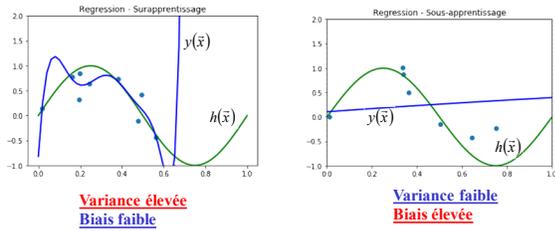
On a vu les concepts de « sur-apprentissage » et de « sous-apprentissage ».



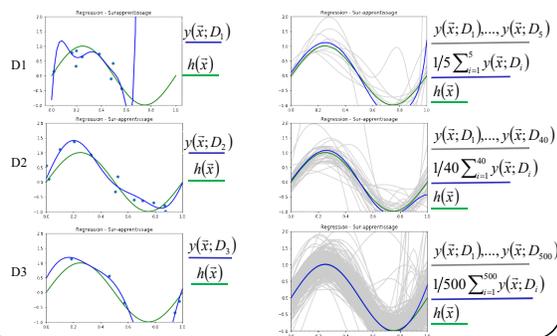
Le but d'un bon modèle est de trouver le **bon compromis biais-variance**

Décomposition biais-variance (Bishop 3.2)

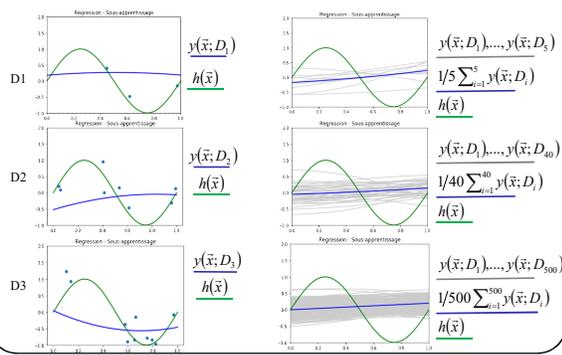
On a vu les concepts de « sur-apprentissage » et de « sous-apprentissage ».



Variance élevée, Biais faible

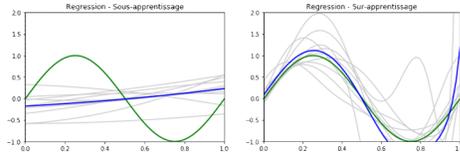


Variance faible, Biais élevé



Variance

Peut être vu intuitivement comme la **écart moyen entre les « courbes grises »** $y(\bar{x}; D_i)$



Variance faible

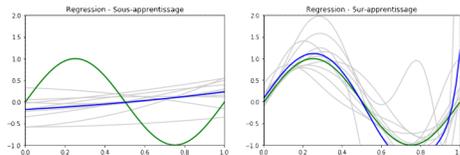
Variance élevée

Dans le cas à 5 courbes:

$$\text{variance} = \frac{1}{5} \sum_{i=1}^5 \left(y(\bar{x}; D_i) - \underbrace{1/5 \sum_{i=1}^5 y(\bar{x}; D_i)}_{\text{Courbe moyenne (bleue)}} \right)^2$$

Biais

Peut être vu intuitivement comme la **différence entre les « courbes verte et bleue »**



Biais élevé

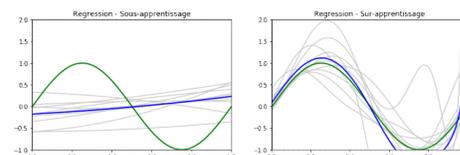
Variance faible

Dans le cas à 5 courbes:

$$\text{biais} = \frac{1}{5} \sum_{i=1}^5 \left(\underbrace{h(\bar{x})}_{\text{Courbe optimale (verte)}} - \underbrace{1/5 \sum_{i=1}^5 y(\bar{x}; D_i)}_{\text{Courbe moyenne (bleue)}} \right)^2$$

29

Compromis biais-variance



Biais élevé
Variance faible

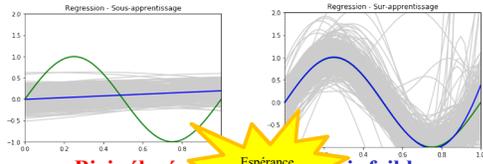
Biais faible
Variance élevée

$$\text{variance} = \frac{1}{5} \sum_{i=1}^5 \left(y(\bar{x}; D_i) - 1/5 \sum_{i=1}^5 y(\bar{x}; D_i) \right)^2$$

$$\text{biais} = \frac{1}{5} \sum_{i=1}^5 \left(h(\bar{x}) - 1/5 \sum_{i=1}^5 y(\bar{x}; D_i) \right)^2$$

30

Analyse quand $N \rightarrow \infty$



Biais élevé
Variance faible

Espérance mathématique

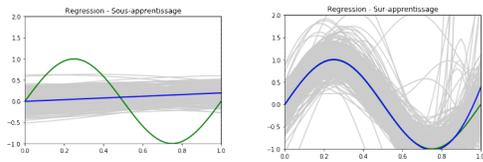
Biais faible
Variance élevée

$$\text{variance} = \frac{1}{N} \sum_{i=1}^N (y(\bar{x}; D_i) - 1/N \sum_{i=1}^N y(\bar{x}; D_i))^2$$

$$\text{biais} = \frac{1}{N} \sum_{i=1}^N (h(\bar{x}) - 1/N \sum_{i=1}^N y(\bar{x}; D_i))^2$$

31

Analyse quand $N \rightarrow \infty$



Biais élevé
Variance faible

Biais faible
Variance élevée

$$\text{variance} = E_D \left[\left(y(\bar{x}; D) - E_D [y(\bar{x}; D)] \right)^2 \right]$$

$$\text{biais} = E_D \left[\left(h(\bar{x}) - E_D [y(\bar{x}; D)] \right)^2 \right]$$

32

Analyse formelle

On a démontré précédemment que

$$E[L] = \iint \{y(\bar{x}) - E[t | \bar{x}]\}^2 p(x, t) dx dt + \iint \{E[t | \bar{x}] - t\}^2 p(x, t) dx dt$$

où

$y(\bar{x})$: modèle entraîné à l'aide de D

$E[t | \bar{x}]$: modèle théorique optimal

(courbe verte)

33

Analyse formelle

Récriture

$$E[L] = \iint \{y(\bar{x}; D) - h(\bar{x})\}^2 p(x, t) dx dt \quad \left. \vphantom{\iint} \right\} \begin{array}{l} \text{Mesure la performance} \\ \text{du modèle } y \end{array}$$

$$+ \iint \{h(\bar{x}) - t\}^2 p(x, t) dx dt \quad \left. \vphantom{\iint} \right\} \begin{array}{l} \text{Mesure la magnitude} \\ \text{du bruit dans les données} \end{array}$$

34

Analyse formelle

$$\{y(\bar{x}; D) - h(\bar{x})\}^2 = \{y(\bar{x}; D) - E_D[y(\bar{x}; D)] + E_D[y(\bar{x}; D)] - h(\bar{x})\}^2$$

$$= \{y(\bar{x}; D) - E_D[y(\bar{x}; D)]\}^2$$

$$+ 2\{y(\bar{x}; D) - E_D[y(\bar{x}; D)]\}\{E_D[y(\bar{x}; D)] - h(\bar{x})\}$$

$$+ \{E_D[y(\bar{x}; D)] - h(\bar{x})\}^2$$

On peut démontrer que

$$E[\{y(\bar{x}; D) - h(\bar{x})\}^2] = E_D[\underbrace{\{y(\bar{x}; D) - E_D[y(\bar{x}; D)]\}^2}_{\text{Variance}}]$$

$$+ E_D[\underbrace{\{h(\bar{x}) - E_D[y(\bar{x}; D)]\}^2}_{\text{Biais}}]$$

36
