

Biais et Variance

August 21, 2019

Mots clefs. Après avoir lu le cours, ces mots doivent vous être familiers.

- Sur-apprentissage.
- Biais / Variance
- Flexibilité / interprétabilité

1 Biais-Variance

1.1 Un exemple de modèle linéaire avec transformation des input

Considérons un input $X \in [0, 1]$ quantitatif et un output $Y \in \mathbb{R}$ quantitatif. Considérons ν (=fréquence maximale) un paramètre. On définit le modèle

suivant :

$$\begin{aligned} Y &= w_0 + \sum_{n=1}^{\nu} w_{2n} \cos(2\pi nX) + w_{2n-1} \sin(2\pi nX) \\ &= f_w(X) \end{aligned}$$

Plus ν est grand, et plus on a de chance de trouver un \hat{w} telle que

$$f_{\hat{w}} \simeq f?$$

Mais, est-il toujours souhaitable de prendre un ν très grand ?

1.2 Flexibilité d'un modèle

Nous définissons la flexibilité d'un modèle par des exemples.

- un modèle linéaire pure et simple est très peu flexible. Il ne s'adapte qu'à des données proches d'un hyper-plan.
- Le modèle SinCos avec $\nu = 30$ est très flexible.
- quand \hat{f} est obtenu en minimisant

$$\hat{f} = \operatorname{argmin}_{f \in \mathcal{F}} \sum_{Train} \operatorname{dist}(f(X_i) - Y_i) + \lambda \operatorname{Penalisation}(f)$$

Plus \mathcal{F} est grand et plus le modèle est flexible.

Plus λ est grand, plus on oblige le modèle à avoir des fonctions régulières, et moins le modèle est flexible.

- la technique des k -plus proches voisins

$$\hat{f}(x) = \frac{1}{k} \sum_{X_i \in V_k(x)} Y_i$$

est très flexible quand $k = 1$, et peu flexible quand k est grand. Cas extrême : si k est égal au nombre de données dans $Train$, \hat{f} sera une fonction constante, partout égale à la moyenne des Y_i : aucune flexibilité.

- Quand \hat{f} est obtenue avec des noyaux gaussiens :

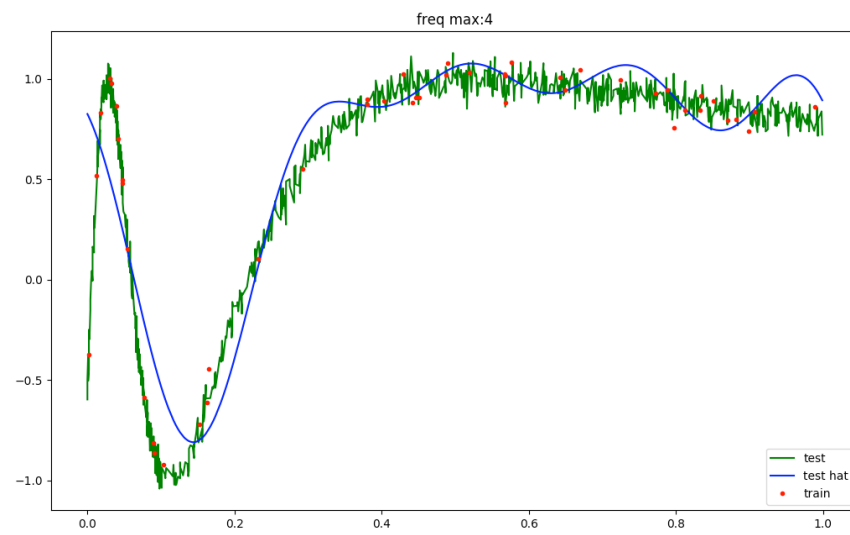
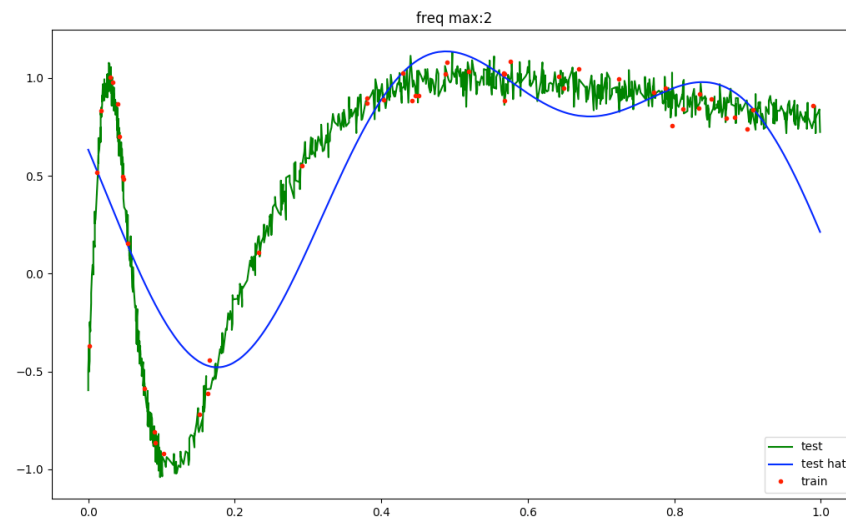
$$\hat{f}(x) = \frac{\sum_{Train} N(x, X_i) Y_i}{\sum_{Train} N(x, X_i)}, \quad N_\sigma(x, y) = e^{-\frac{1}{2} \left(\frac{x-y}{\sigma} \right)^2}$$

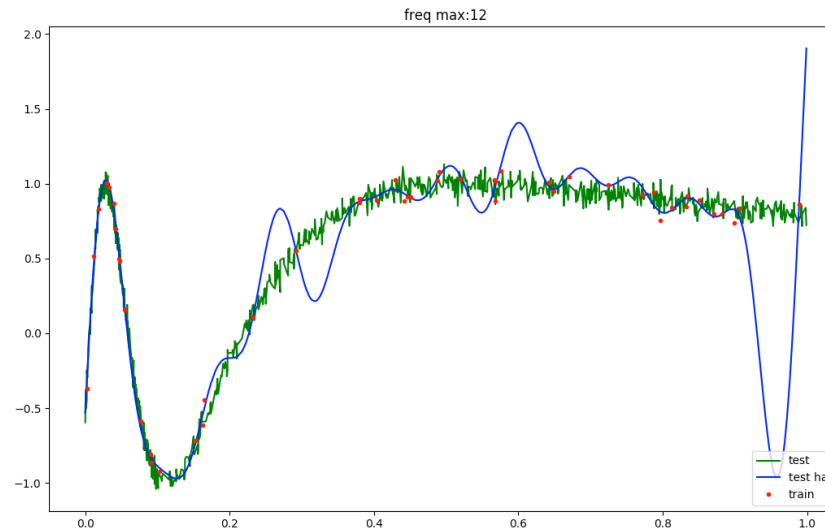
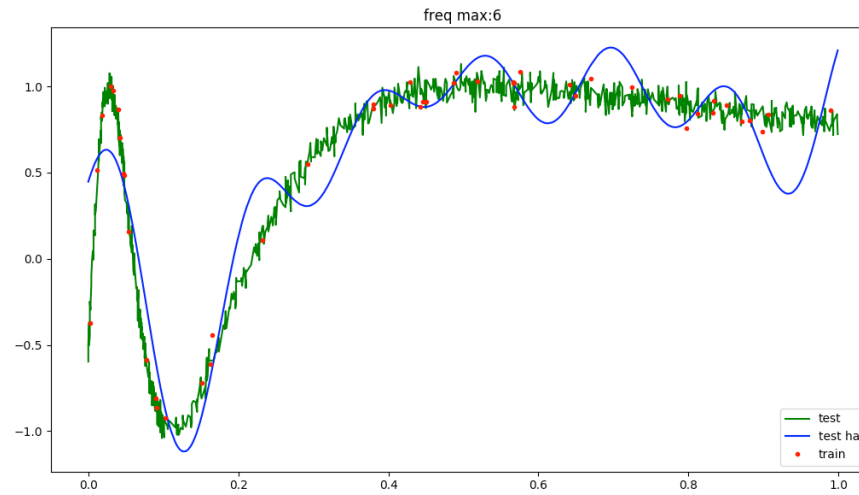
Plus σ est grand, plus la fonction $y \rightarrow N_\sigma(x, y)$ s'étale autour de x , et moins le modèle est flexible.

- De manière générale, plus un modèle s'adapte localement aux données, et plus il est flexible.
- Plus un modèle est flexible, plus il est complexe et moins il est interprétable.

1.3 Sur-apprentissage

Observons notre modèle sin-cos pour différents ν . On dispose de données $Train$ avec lesquelles on entraîne le modèle (= on calcule \hat{w}) et on dispose de données $Test$ très nombreuse qui nous permette de juger ce modèle.





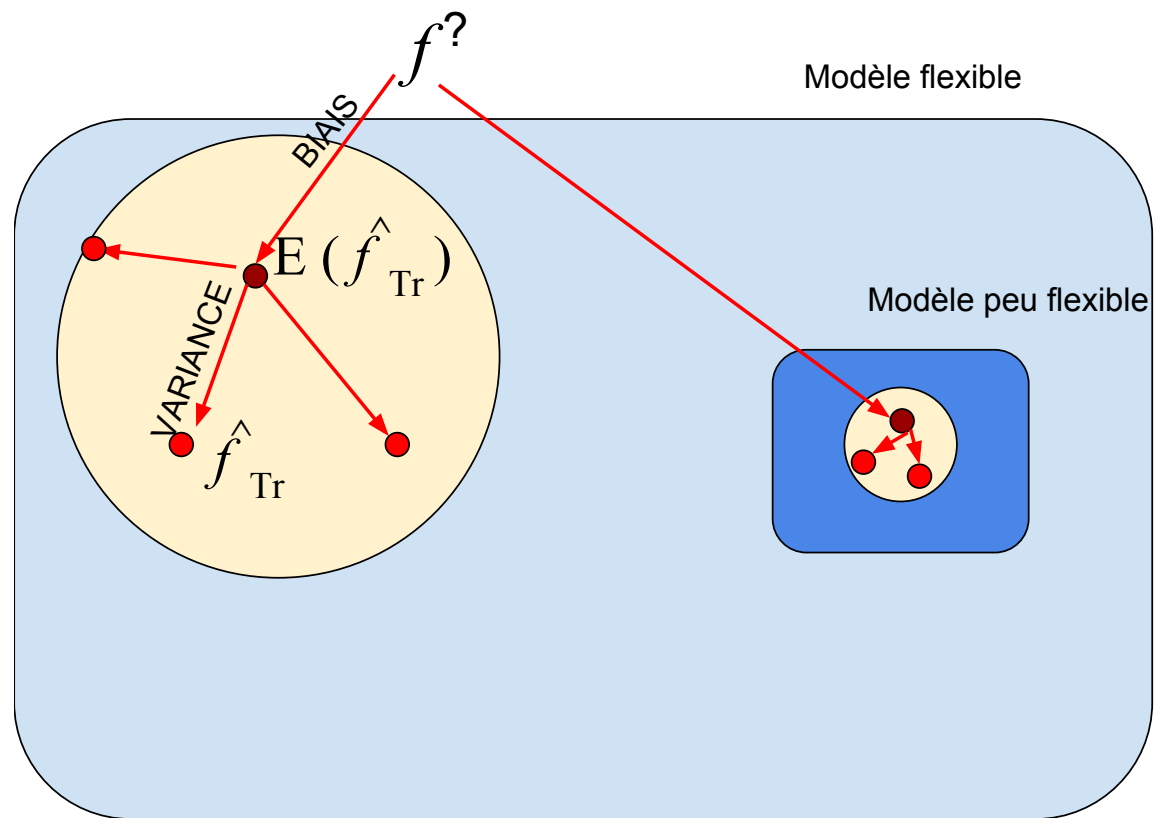
On remarque que plus ν est grand et plus l'estimation colle à $Train$ mais plus elle s'éloigne de $Test$ par endroit. Le modèle a appris par coeur les données

Train, et du coup il est mauvais pour estimer des données *Test*. Il y a eu sur-apprentissage.

1.4 Biais-Variance

Ainsi, un modèle flexible a une forte "probabilité" de s'éloigner de *Test*. Insistons sur le mot "probabilité". En effet, la fonction apprise \hat{f} dépend de l'échantillon d'apprentissage *Train* qui est aléatoire. Pour montrer cette dépendance, par la suite, nous noterons $T_r = \text{Train}$ et $\hat{f}_{T_r} = \hat{f}$ (la fonction d'estimation construite à partir de T_r).

Avec $\nu = 30$, la fonction \hat{f}_{T_r} est très sensible aux variations de T_r . A l'inverse avec $\nu = 1$, \hat{f}_{T_r} sera quasi la même d'un tirage à l'autre de T_r . On dit qu'un modèle flexible, a une grande "variance".



1.5 Quantifions

Considérons un x et sa prédiction $\widehat{f}_{T_r}(x)$. L'erreur quadratique pour l'observation (x, y) se décompose comme ceci :

$$\begin{aligned}\mathbf{E}\left(\widehat{f}_{T_r}(x) - y\right)^2 &= \mathbf{E}\left(\widehat{f}_{T_r}(x) - \mathbf{E}[\widehat{f}_{T_r}(x)] + \mathbf{E}[\widehat{f}_{T_r}(x)] - y\right)^2 \\ &= \underbrace{\mathbf{E}\left(\widehat{f}_{T_r}(x) - \mathbf{E}[\widehat{f}_{T_r}(x)]\right)^2}_{\text{variance}} + \underbrace{\left(\mathbf{E}[\widehat{f}_{T_r}(x)] - y\right)^2}_{\text{biais}}\end{aligned}$$

(la seule chose aléatoire ci-dessus c'est T_r)

Le premier terme s'appelle la variance en (x, y) , le second terme s'appelle le biais de en (x, y) .

Mais il est plus naturel d'estimer une erreur quadratique moyenne, une variance moyenne et un biais moyen, ce qui revient à intégrer sur tous les x, y :

$$\begin{aligned}&\int \mathbf{E}\left(\widehat{f}_{T_r}(x) - y\right)^2 \mathbf{P}[X \in dx, Y \in dy] \\ &= \int \mathbf{E}\left(\widehat{f}_{T_r}(x) - \mathbf{E}[\widehat{f}_{T_r}(x)]\right)^2 \mathbf{P}[X \in dx] \\ &+ \int \left(\mathbf{E}[\widehat{f}_{T_r}(x)] - y\right)^2 \mathbf{P}[X \in dx, Y \in dy]\end{aligned}$$

En plus court:

$$\mathbf{E}\left(\widehat{f}_{T_r}(X) - Y\right)^2 = \mathbf{E}\left(\widehat{f}_{T_r}(X) - \mathbf{E}[\widehat{f}_{T_r}(X)]\right)^2 + \mathbf{E}\left(\mathbf{E}[\widehat{f}_{T_r}(X)] - Y\right)^2$$

(Cette dernière formule est difficile à comprendre car l'espérance porte sur plusieurs niveaux d'aléatoire)

1.6 Estimons numériquement

Estimons maintenant ces quantités à l'aide de plusieurs réalisations d'échantillons d'apprentissage $(tr_i)_{i \in 1..I}$ et d'une réalisation d'un échantillon test $(x_j, y_j)_{j \in 1..J}$.

Ainsi la décomposition biais variance pour une valeur (x, y) :

$$\mathbf{E}\left(\widehat{f}_{T_r}(x) - y\right)^2 = \mathbf{E}\left(\widehat{f}_{T_r}(x) - \mathbf{E}[\widehat{f}_{T_r}(x)]\right)^2 + \left(\mathbf{E}[\widehat{f}_{T_r}(x)] - y\right)^2$$

S'estime par :

$$\begin{aligned} \text{mean}_i \left[(\widehat{f}_{tr_i}(x) - y)^2 \right] &\simeq \text{mean}_i \left[(\widehat{f}_{tr_i}(x) - \text{mean}_k[\widehat{f}_{tr_k}(x)])^2 \right] + \left(\text{mean}_i[\widehat{f}_{tr_i}(x)] - y \right)^2 \\ &\simeq \text{std}_i^2 \left[\widehat{f}_{tr_i}(x) \right] + \left(\text{mean}_i[\widehat{f}_{tr_i}(x)] - y \right)^2 \end{aligned}$$

La décomposition biais-variance moyenne :

$$\mathbf{E}\left(\widehat{f}_{T_r}(X) - Y\right)^2 = \mathbf{E}\left(\widehat{f}_{T_r}(X) - \mathbf{E}[\widehat{f}_{T_r}(X)]\right)^2 + \mathbf{E}\left(\mathbf{E}[\widehat{f}_{T_r}(X)] - Y\right)^2$$

s'estime par :

$$\text{mean}_j \text{ mean}_i \left[(\widehat{f}_{tr_i}(x_j) - y_j)^2 \right] \simeq \text{mean}_j \text{ std}_i^2 \left[\widehat{f}_{tr_i}(x_j) \right] + \text{mean}_j \left(\text{mean}_i [\widehat{f}_{tr_i}(x_j)] - y_j \right)^2$$

1.7 Remarque sur la "forme"

Ecrire du pseudo code aide beaucoup pour simplémenter des formules ou des algorithmes. Au dessus, on a utilisé des notations proches de numpy :

- $\text{mean}_i[u]$ pour `np.mean(u,axis=i)`
- $\text{std}_i^2[u]$ pour `np.std(u,axis=i)**2`

Mathématiquement, la formule d'estimation aurait plutôt été :

$$\frac{1}{IJ} \sum_{i,j} \left(\widehat{f}_{tr_i}(x_j) - y_j \right)^2 \simeq \frac{1}{IJ} \sum_{i,j} \left(\widehat{f}_{tr_i}(x_j) - \frac{1}{I} \sum_k \widehat{f}_{tr_k}(x_j) \right)^2 + \frac{1}{J} \sum_j \left(\frac{1}{I} \sum_k \widehat{f}_{tr_k}(x_j) - y_j \right)^2$$

Ce qui est plus compliqué à transformer en code !

1.8 La décomposition biais-variance avec une courbe

