

OVERFITTING AND BIAS-VARIANCE TRADE-OFF

Jérémie Cabessa
Laboratoire DAVID, UVSQ

FORMULATION DU PROBLÈME (APPRENTISSAGE SUPERVISÉ)

- ▶ Soient X_1, \dots, X_p et Y des variables aléatoires.
- ▶ X_1, \dots, X_p sont appelées variables d'inputs, variables indépendantes, variables explicatives, prédicteurs, (features).
- ▶ Y est appelée variable d'output, variable dépendante, réponse, (response, target).

FORMULATION DU PROBLÈME (APPRENTISSAGE SUPERVISÉ)

- ▶ Soient X_1, \dots, X_p et Y des variables aléatoires.
- ▶ X_1, \dots, X_p sont appelées **variables d'inputs**, **variables indépendantes**, **variables explicatives**, **prédicteurs**, (**features**).
- ▶ Y est appelée **variable d'output**, **variable dépendante**, **réponse**, (**response**, **target**).

FORMULATION DU PROBLÈME (APPRENTISSAGE SUPERVISÉ)

- ▶ Soient X_1, \dots, X_p et Y des variables aléatoires.
- ▶ X_1, \dots, X_p sont appelées **variables d'inputs**, **variables indépendantes**, **variables explicatives**, **prédicteurs**, (**features**).
- ▶ Y est appelée **variable d'output**, **variable dépendante**, **réponse**, (**response**, **target**).

FORMULATION DU PROBLÈME (APPRENTISSAGE SUPERVISÉ)

- ▶ On suppose qu'il existe une (vraie) **relation** f entre X_1, \dots, X_p et Y de la forme

$$Y = f(X_1, \dots, X_p) + \epsilon$$

où f est une fonction inconnue et ϵ est une variable aléatoire indépendante de X_1, \dots, X_p et de moyenne 0, le bruit.

- ▶ On aimerait apprendre une (bonne) **estimation** \hat{f} de f . On aura alors

$$\hat{Y} = \hat{f}(X_1, \dots, X_p)$$

où \hat{f} est l'estimation de f et \hat{Y} est la **prediction** de Y .

FORMULATION DU PROBLÈME (APPRENTISSAGE SUPERVISÉ)

- ▶ On suppose qu'il existe une (vraie) **relation** f entre X_1, \dots, X_p et Y de la forme

$$Y = f(X_1, \dots, X_p) + \epsilon$$

où f est une fonction inconnue et ϵ est une variable aléatoire indépendante de X_1, \dots, X_p et de moyenne 0, le bruit.

- ▶ On aimerait apprendre une (bonne) **estimation** \hat{f} de f . On aura alors

$$\hat{Y} = \hat{f}(X_1, \dots, X_p)$$

où \hat{f} est l'**estimation** de f et \hat{Y} est la **prediction** de Y .

FORMULATION DU PROBLÈME (APPRENTISSAGE SUPERVISÉ)

- Pour apprendre l'estimation \hat{f} de f , on dispose de données (data)

$$S_{\text{train}} = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n)\}$$

où $\mathbf{x}_i = (x_{i1}, \dots, x_{ip})$ pour tout $i = 1, \dots, n$.

- Ces données constituent le “training set” (training data).

FORMULATION DU PROBLÈME (APPRENTISSAGE SUPERVISÉ)

- Pour apprendre l'estimation \hat{f} de f , on dispose de données (data)

$$S_{\text{train}} = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n)\}$$

où $\mathbf{x}_i = (x_{i1}, \dots, x_{ip})$ pour tout $i = 1, \dots, n$.

- Ces données constituent le “**training set**” (training data).

FORMULATION DU PROBLÈME (APPRENTISSAGE SUPERVISÉ)

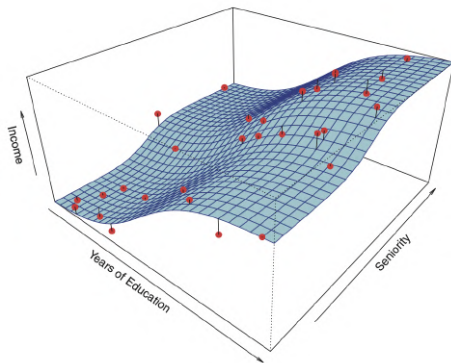


Figure taken from [James et al., 2013]

ERREUR RÉDUCTIBLE ET IRRÉDUCTIBLE

► On a donc

$$\begin{aligned} Y &= f(X_1, \dots, X_p) + \epsilon && \text{vraie relation} \\ \hat{Y} &= \hat{f}(X_1, \dots, X_p) && \text{estimation} \end{aligned}$$

► On peut facilement montrer que

$$\begin{aligned} \mathbb{E}[Y - \hat{Y}]^2 &= \mathbb{E}[f(X_1, \dots, X_p) + \epsilon - \hat{f}(X_1, \dots, X_p)]^2 \\ &= \mathbb{E}[f(\mathbf{X}) + \epsilon - \hat{f}(\mathbf{X})]^2 \\ &= \underbrace{(f(\mathbf{X}) - \hat{f}(\mathbf{X}))^2}_{\text{erreur réductible}} + \underbrace{\text{Var}[\epsilon]}_{\text{erreur irréductible}} \end{aligned}$$

où $\mathbf{X} = (X_1, \dots, X_p)$.

ERREUR RÉDUCTIBLE ET IRRÉDUCTIBLE

- On a donc

$$\begin{aligned} Y &= f(X_1, \dots, X_p) + \epsilon && \text{vraie relation} \\ \hat{Y} &= \hat{f}(X_1, \dots, X_p) && \text{estimation} \end{aligned}$$

- On peut facilement montrer que

$$\begin{aligned} \mathbb{E}[Y - \hat{Y}]^2 &= \mathbb{E}[f(X_1, \dots, X_p) + \epsilon - \hat{f}(X_1, \dots, X_p)]^2 \\ &= \mathbb{E}[f(\mathbf{X}) + \epsilon - \hat{f}(\mathbf{X})]^2 \\ &= \underbrace{(f(\mathbf{X}) - \hat{f}(\mathbf{X}))^2}_{\text{erreur réductible}} + \underbrace{\text{Var}[\epsilon]}_{\text{erreur irréductible}} \end{aligned}$$

où $\mathbf{X} = (X_1, \dots, X_p)$.

ERREUR RÉDUCTIBLE ET IRRÉDUCTIBLE

- ▶ **Erreur réductible (reducible error):** peut être réduite; plus \hat{f} est une bonne estimation de f , plus cette erreur sera faible.
- ▶ **Erreur irréductible (irreducible error):** ne peut être réduite; par le biais de notre estimation \hat{f} , nous n'avons aucune prise sur le "bruit" réel intrinsèque au modèle.

ERREUR RÉDUCTIBLE ET IRRÉDUCTIBLE

- ▶ **Erreur réductible (reducible error):** peut être réduite; plus \hat{f} est une bonne estimation de f , plus cette erreur sera faible.
- ▶ **Erreur irréductible (irreducible error):** ne peut être réduite; par le biais de notre estimation \hat{f} , nous n'avons aucune prise sur le “bruit” réel intrinsèque au modèle.

ERREUR RÉDUCTIBLE ET IRRÉDUCTIBLE

- Exemple de deux estimations \hat{f} . La deuxième estimation est meilleure car elle est associée à une erreur réductible plus faible.

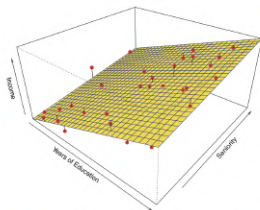


FIGURE 2.4. A linear model fit by least squares to the **Income** data from Figure 2.3. The observations are shown in red, and the yellow plane indicates the least squares fit to the data.

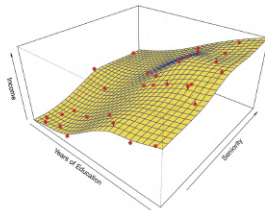


FIGURE 2.5. A smooth thin-plate spline fit to the **Income** data from Figure 2.3 is shown in yellow; the observations are displayed in red. Splines are discussed in Chapter 7.

Figure taken from [James et al., 2013]

FONCTION DE COÛT (COST FUNCTION)

- ▶ Comment mesurer la qualité d'un modèle \hat{f} ?
- ▶ On utilise une fonction de coût (cost or loss function).
- ▶ La plus célèbre est l'erreur des moindres carrés (mean squared error) **MSE**. Étant donné un training set

$$S_{\text{train}} = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$$

on définit

$$\text{MSE}_{\text{train}} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{f}(x_i))^2$$

FONCTION DE COÛT (COST FUNCTION)

- ▶ Comment mesurer la qualité d'un modèle \hat{f} ?
- ▶ On utilise une **fonction de coût (cost or loss function)**.
- ▶ La plus célèbre est l'erreur des moindres carrés (mean squared error) **MSE**. Étant donné un training set

$$S_{\text{train}} = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n)\}$$

on définit

$$\text{MSE}_{\text{train}} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{f}(\mathbf{x}_i))^2$$

FONCTION DE COÛT (COST FUNCTION)

- ▶ Comment mesurer la qualité d'un modèle \hat{f} ?
- ▶ On utilise une **fonction de coût** (cost or loss function).
- ▶ La plus célèbre est l'**erreur des moindres carrés** (mean squared error) **MSE**. Étant donné un training set

$$S_{\text{train}} = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n)\}$$

on définit

$$\text{MSE}_{\text{train}} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{f}(\mathbf{x}_i))^2$$

OVERFITTING

- ▶ **Problème:** Le modèle \hat{f} a été construit sur la base du training set S_{train} . Ainsi, il peut être très performant lorsqu'il est évalué sur le training set, mais bien moins performant lorsqu'il est évalué sur des données qu'il n'a jamais vues.
- ▶ Ainsi, il convient d'utiliser un **training set** S_{train} pour construire le modèle \hat{f} , et un autre ensemble de données disjoint, appelé **test set** S_{test} , pour évaluer la performance du modèle \hat{f} .
- ▶ Lorsque le modèle est performant sur le **training set** ($\text{MSE}_{\text{train}}$ basse), mais qu'il est bien moins performant sur le **test set** (MSE_{test} plus élevée), on est dans une situation d'**overfitting**.

OVERFITTING

- ▶ **Problème:** Le modèle \hat{f} a été construit sur la base du training set S_{train} . Ainsi, il peut être très performant lorsqu'il est évalué sur le training set, mais bien moins performant lorsqu'il est évalué sur des données qu'il n'a jamais vues.
- ▶ Ainsi, il convient d'utiliser un **training set** S_{train} pour construire le modèle \hat{f} , et un autre ensemble de données disjoint, appelé **test set** S_{test} , pour évaluer la performance du modèle \hat{f} .
- ▶ Lorsque le modèle est performant sur le **training set** ($\text{MSE}_{\text{train}}$ basse), mais qu'il est bien moins performant sur le **test set** (MSE_{test} plus élevée), on est dans une situation d'**overfitting**.

OVERFITTING

- ▶ **Problème:** Le modèle \hat{f} a été construit sur la base du training set S_{train} . Ainsi, il peut être très performant lorsqu'il est évalué sur le training set, mais bien moins performant lorsqu'il est évalué sur des données qu'il n'a jamais vues.
- ▶ Ainsi, il convient d'utiliser un **training set** S_{train} pour construire le modèle \hat{f} , et un autre ensemble de données disjoint, appelé **test set** S_{test} , pour évaluer la performance du modèle \hat{f} .
- ▶ Lorsque le modèle est performant sur le **training set** ($\text{MSE}_{\text{train}}$ basse), mais qu'il est bien moins performant sur le **test set** (MSE_{test} plus élevée), on est dans une situation d'**overfitting**.

OVERFITTING

- ▶ **Overfitting:** Le modèle est très performant sur le **training set** (i.e., MSE_{train} basse) alors qu'il est bien moins performant sur le **test set** (i.e., MSE_{train} élevée).
- ▶ Le modèle a donc "sur-appris" (overfit) les données d'apprentissage (train set), de sorte que ses performances ne se généralisent pas bien sur des données inconnues (test set).
- ▶ En fait, le modèle a appris le bruit des données d'apprentissage, au lieu de l'ignorer.

OVERFITTING

- ▶ **Overfitting:** Le modèle est très performant sur le **training set** (i.e., MSE_{train} basse) alors qu'il est bien moins performant sur le **test set** (i.e., MSE_{train} élevée).
- ▶ Le modèle a donc “sur-appris” (overfit) les données d'apprentissage (train set), de sorte que ses performances ne se généralisent pas bien sur des données inconnues (test set).
- ▶ En fait, le modèle a appris le bruit des données d'apprentissage, au lieu de l'ignorer.

OVERFITTING

- ▶ **Overfitting:** Le modèle est très performant sur le **training set** (i.e., MSE_{train} basse) alors qu'il est bien moins performant sur le **test set** (i.e., MSE_{train} élevée).
- ▶ Le modèle a donc “sur-appris” (overfit) les données d'apprentissage (train set), de sorte que ses performances ne se généralisent pas bien sur des données inconnues (test set).
- ▶ En fait, le modèle a appris le bruit des données d'apprentissage, au lieu de l'ignorer.

OVERFITTING

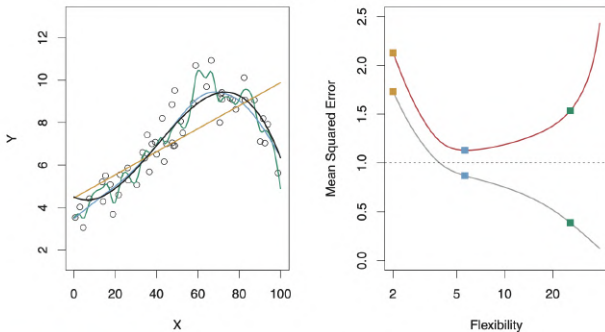


FIGURE 2.9. Left: Data simulated from f , shown in black. Three estimates of f are shown: the linear regression line (orange curve), and two smoothing spline fits (blue and green curves). Right: Training MSE (grey curve), test MSE (red curve), and minimum possible test MSE over all methods (dashed line). Squares represent the training and test MSEs for the three fits shown in the left-hand panel.

Figure taken from [James et al., 2013]

OVERFITTING

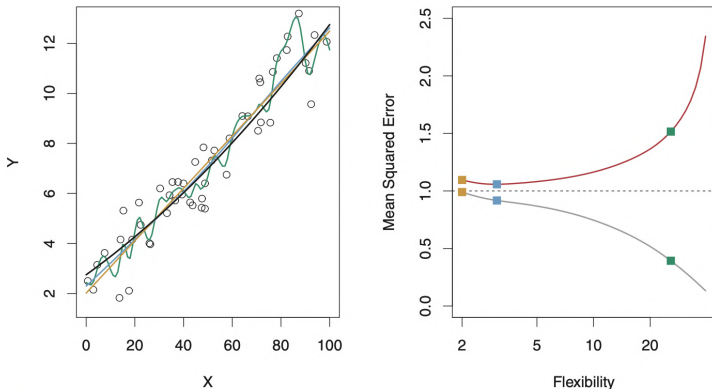


FIGURE 2.10. Details are as in Figure 2.9, using a different true f that is much closer to linear. In this setting, linear regression provides a very good fit to the data.

Figure taken from [James et al., 2013]

OVERFITTING

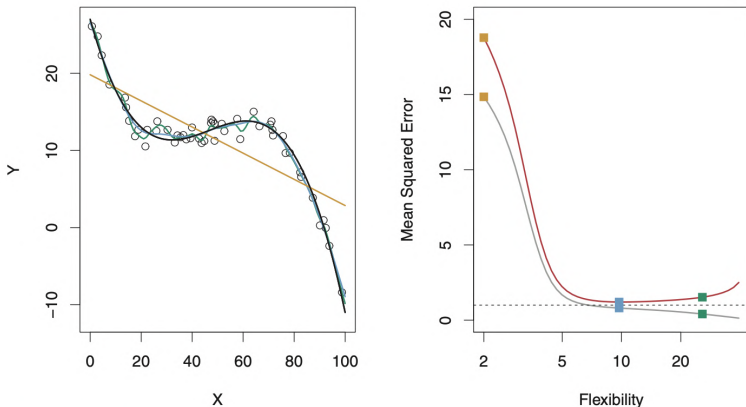


FIGURE 2.11. Details are as in Figure 2.9, using a different f that is far from linear. In this setting, linear regression provides a very poor fit to the data.

Figure taken from [James et al., 2013]

BIAS-VARIANCE TRADE-OFF

- Soit une relation fonctionnelle

$$Y = f(\mathbf{X}) + \epsilon$$

entre des variables d'input $\mathbf{X} = (X_1, \dots, X_p)$ et d'output Y .
Le bruit ϵ satisfait $E[\epsilon] = 0$.

- On cherche à obtenir un modèle

$$\hat{Y} = \hat{f}(\mathbf{X})$$

qui soit le plus performant possible *sur le test set* !

BIAS-VARIANCE TRADE-OFF

- Soit une relation fonctionnelle

$$Y = f(\mathbf{X}) + \epsilon$$

entre des variables d'input $\mathbf{X} = (X_1, \dots, X_p)$ et d'output Y .
Le bruit ϵ satisfait $E[\epsilon] = 0$.

- On cherche à obtenir un modèle

$$\hat{Y} = \hat{f}(\mathbf{X})$$

qui soit le plus performant possible *sur le test set* !

BIAS-VARIANCE TRADE-OFF

- Soit un dataset

$$S = \{(\mathbf{x}_i, y_i) \in \mathbb{R}^M \times \mathbb{R} : i = 1, \dots, N\}$$

et soit $(\mathbf{x}, y) \in S$ un point du dataset.

- On note $E_{S, \epsilon}[\dots] := E[\dots]$.
- Par hypothèse, on a: $E[\epsilon] = 0$. Et puisque la "vraie" relation fonctionnelle f est déterministe, on a: $E[f] = f$.
- On rappelle que

$$\text{Var}[X] := E[(X - E[X])^2] = E[X^2] - E[X]^2.$$

BIAS-VARIANCE TRADE-OFF

- Soit un dataset

$$S = \{(\mathbf{x}_i, y_i) \in \mathbb{R}^M \times \mathbb{R} : i = 1, \dots, N\}$$

et soit $(\mathbf{x}, y) \in S$ un point du dataset.

- On note $E_{S, \epsilon}[\dots] := E[\dots]$.
- Par hypothèse, on a: $E[\epsilon] = 0$. Et puisque la “vraie” relation fonctionnelle f est déterministe, on a: $E[f] = f$.
- On rappelle que

$$\text{Var}[X] := E[(X - E[X])^2] = E[X^2] - E[X]^2.$$

BIAS-VARIANCE TRADE-OFF

- Soit un dataset

$$S = \{(\mathbf{x}_i, y_i) \in \mathbb{R}^M \times \mathbb{R} : i = 1, \dots, N\}$$

et soit $(\mathbf{x}, y) \in S$ un point du dataset.

- On note $E_{S, \epsilon}[\dots] := E[\dots]$.
- Par hypothèse, on a: $E[\epsilon] = 0$. Et puisque la “vraie” relation fonctionnelle f est déterministe, on a: $E[f] = f$.
- On rappelle que

$$\text{Var}[X] := E[(X - E[X])^2] = E[X^2] - E[X]^2.$$

BIAS-VARIANCE TRADE-OFF

- Soit un dataset

$$S = \{(\mathbf{x}_i, y_i) \in \mathbb{R}^M \times \mathbb{R} : i = 1, \dots, N\}$$

et soit $(\mathbf{x}, y) \in S$ un point du dataset.

- On note $E_{S, \epsilon}[\dots] := E[\dots]$.
- Par hypothèse, on a: $E[\epsilon] = 0$. Et puisque la “vraie” relation fonctionnelle f est déterministe, on a: $E[f] = f$.
- On rappelle que

$$\text{Var}[X] := E[(X - E[X])^2] = E[X^2] - E[X]^2.$$

BIAS-VARIANCE TRADE-OFF

On a:

$$\begin{aligned}
\mathbb{E}[(y - \hat{f}(x))^2] &= \mathbb{E}[(f(x) + \epsilon - \hat{f}(x))^2] \\
&= \mathbb{E}\left[\underbrace{(f(x) - \mathbb{E}[\hat{f}(x)])}_{=A} + \underbrace{\epsilon}_{=B} + \underbrace{(\mathbb{E}[\hat{f}(x)] - \hat{f}(x))}_{=C}\right]^2 \\
&= \mathbb{E}[A^2 + B^2 + C^2 + 2AB + 2BC + 2CA] \\
(\mathbb{E}[B] = \mathbb{E}[C] = 0) \quad &= \mathbb{E}[A^2] + \mathbb{E}[B^2] + \mathbb{E}[C^2] \\
(\mathbb{E}[\epsilon] = 0) \quad &= \mathbb{E}[(f(x) - \mathbb{E}[\hat{f}(x)])^2] + \mathbb{E}[(\epsilon - \mathbb{E}[\epsilon])^2] + \mathbb{E}[(\hat{f}(x) - \mathbb{E}[\hat{f}(x)])^2] \\
(\mathbb{E}[f(x)] = f(x)) \quad &= \left(\mathbb{E}[\hat{f}(x)] - f(x)\right)^2 + \text{Var}[\hat{f}(x)] + \text{Var}[\epsilon] \\
&= \text{Biais}[\hat{f}(x)]^2 + \text{Var}[\hat{f}(x)] + \text{Var}[\epsilon]
\end{aligned}$$

où

$$\text{Biais}[\hat{f}(x)] := \mathbb{E}[\hat{f}(x)] - f(x) \quad \text{Var}[\hat{f}(x)] := \mathbb{E}[(\hat{f}(x) - \mathbb{E}[\hat{f}(x)])^2]$$

BIAS-VARIANCE TRADE-OFF

On a:

$$\begin{aligned}
\mathbb{E}[(y - \hat{f}(x))^2] &= \mathbb{E}[(f(x) + \epsilon - \hat{f}(x))^2] \\
&= \mathbb{E}\left[\underbrace{(f(x) - \mathbb{E}[\hat{f}(x)])}_{=A} + \underbrace{\epsilon}_{=B} + \underbrace{(\mathbb{E}[\hat{f}(x)] - \hat{f}(x))}_{=C}\right]^2 \\
&= \mathbb{E}[A^2 + B^2 + C^2 + 2AB + 2BC + 2CA] \\
(\mathbb{E}[B] = \mathbb{E}[C] = 0) \quad &= \mathbb{E}[A^2] + \mathbb{E}[B^2] + \mathbb{E}[C^2] \\
(\mathbb{E}[\epsilon] = 0) \quad &= \mathbb{E}[(f(x) - \mathbb{E}[\hat{f}(x)])^2] + \mathbb{E}[(\epsilon - \mathbb{E}[\epsilon])^2] + \mathbb{E}[(\hat{f}(x) - \mathbb{E}[\hat{f}(x)])^2] \\
(\mathbb{E}[f(x)] = f(x)) \quad &= \left(\mathbb{E}[\hat{f}(x)] - f(x)\right)^2 + \text{Var}[\hat{f}(x)] + \text{Var}[\epsilon] \\
&= \text{Biais}[\hat{f}(x)]^2 + \text{Var}[\hat{f}(x)] + \text{Var}[\epsilon]
\end{aligned}$$

où

$$\text{Biais}[\hat{f}(x)] := \mathbb{E}[\hat{f}(x)] - f(x) \quad \text{Var}[\hat{f}(x)] := \mathbb{E}[(\hat{f}(x) - \mathbb{E}[\hat{f}(x)])^2]$$

BIAS-VARIANCE TRADE-OFF

On a:

$$\begin{aligned}
\mathbb{E}[(y - \hat{f}(x))^2] &= \mathbb{E}[(f(x) + \epsilon - \hat{f}(x))^2] \\
&= \mathbb{E}\left[\underbrace{(f(x) - \mathbb{E}[\hat{f}(x)])}_{=A} + \underbrace{\epsilon}_{=B} + \underbrace{(\mathbb{E}[\hat{f}(x)] - \hat{f}(x))}_{=C}\right]^2 \\
&= \mathbb{E}[A^2 + B^2 + C^2 + 2AB + 2BC + 2CA] \\
(\mathbb{E}[B] = \mathbb{E}[C] = 0) \quad &= \mathbb{E}[A^2] + \mathbb{E}[B^2] + \mathbb{E}[C^2] \\
(\mathbb{E}[\epsilon] = 0) \quad &= \mathbb{E}[(f(x) - \mathbb{E}[\hat{f}(x)])^2] + \mathbb{E}[(\epsilon - \mathbb{E}[\epsilon])^2] + \mathbb{E}[(\hat{f}(x) - \mathbb{E}[\hat{f}(x)])^2] \\
(\mathbb{E}[f(x)] = f(x)) \quad &= \left(\mathbb{E}[\hat{f}(x)] - f(x)\right)^2 + \text{Var}[\hat{f}(x)] + \text{Var}[\epsilon] \\
&= \text{Biais}[\hat{f}(x)]^2 + \text{Var}[\hat{f}(x)] + \text{Var}[\epsilon]
\end{aligned}$$

où

$$\text{Biais}[\hat{f}(x)] := \mathbb{E}[\hat{f}(x)] - f(x) \quad \text{Var}[\hat{f}(x)] := \mathbb{E}[(\hat{f}(x) - \mathbb{E}[\hat{f}(x)])^2]$$

BIAS-VARIANCE TRADE-OFF

On a:

$$\begin{aligned}
\mathbb{E}[(y - \hat{f}(x))^2] &= \mathbb{E}[(f(x) + \epsilon - \hat{f}(x))^2] \\
&= \mathbb{E}\left[\underbrace{(f(x) - \mathbb{E}[\hat{f}(x)])}_{=A} + \underbrace{\epsilon}_{=B} + \underbrace{(\mathbb{E}[\hat{f}(x)] - \hat{f}(x))}_{=C}\right]^2 \\
&= \mathbb{E}[A^2 + B^2 + C^2 + 2AB + 2BC + 2CA] \\
(\mathbb{E}[B] = \mathbb{E}[C] = 0) \quad &= \mathbb{E}[A^2] + \mathbb{E}[B^2] + \mathbb{E}[C^2] \\
(\mathbb{E}[\epsilon] = 0) \quad &= \mathbb{E}[(f(x) - \mathbb{E}[\hat{f}(x)])^2] + \mathbb{E}[(\epsilon - \mathbb{E}[\epsilon])^2] + \mathbb{E}[(\hat{f}(x) - \mathbb{E}[\hat{f}(x)])^2] \\
(\mathbb{E}[f(x)] = f(x)) \quad &= \left(\mathbb{E}[\hat{f}(x)] - f(x)\right)^2 + \text{Var}[\hat{f}(x)] + \text{Var}[\epsilon] \\
&= \text{Biais}[\hat{f}(x)]^2 + \text{Var}[\hat{f}(x)] + \text{Var}[\epsilon]
\end{aligned}$$

où

$$\text{Biais}[\hat{f}(x)] := \mathbb{E}[\hat{f}(x)] - f(x) \quad \text{Var}[\hat{f}(x)] := \mathbb{E}[(\hat{f}(x) - \mathbb{E}[\hat{f}(x)])^2]$$

BIAS-VARIANCE TRADE-OFF

On a:

$$\begin{aligned}
\mathbb{E}[(y - \hat{f}(\mathbf{x}))^2] &= \mathbb{E}[(f(\mathbf{x}) + \epsilon - \hat{f}(\mathbf{x}))^2] \\
&= \mathbb{E}\left[\underbrace{(f(\mathbf{x}) - \mathbb{E}[\hat{f}(\mathbf{x})])}_{=A} + \underbrace{\epsilon}_{=B} + \underbrace{(\mathbb{E}[\hat{f}(\mathbf{x})] - \hat{f}(\mathbf{x}))}_{=C}\right]^2 \\
&= \mathbb{E}[A^2 + B^2 + C^2 + 2AB + 2BC + 2CA] \\
(\mathbb{E}[B] = \mathbb{E}[C] = 0) \quad &= \mathbb{E}[A^2] + \mathbb{E}[B^2] + \mathbb{E}[C^2] \\
(\mathbb{E}[\epsilon] = 0) \quad &= \mathbb{E}[(f(\mathbf{x}) - \mathbb{E}[\hat{f}(\mathbf{x})])^2] + \mathbb{E}[(\epsilon - \mathbb{E}[\epsilon])^2] + \mathbb{E}[(\hat{f}(\mathbf{x}) - \mathbb{E}[\hat{f}(\mathbf{x})])^2] \\
(\mathbb{E}[f(\mathbf{x})] = f(\mathbf{x})) \quad &= \left(\mathbb{E}[\hat{f}(\mathbf{x})] - f(\mathbf{x})\right)^2 + \text{Var}[\hat{f}(\mathbf{x})] + \text{Var}[\epsilon] \\
&= \text{Biais}[\hat{f}(\mathbf{x})]^2 + \text{Var}[\hat{f}(\mathbf{x})] + \text{Var}[\epsilon]
\end{aligned}$$

où

$$\text{Biais}[\hat{f}(\mathbf{x})] := \mathbb{E}[\hat{f}(\mathbf{x})] - f(\mathbf{x}) \quad \text{Var}[\hat{f}(\mathbf{x})] := \mathbb{E}[(\hat{f}(\mathbf{x}) - \mathbb{E}[\hat{f}(\mathbf{x})])^2]$$

BIAS-VARIANCE TRADE-OFF

On a:

$$\begin{aligned}
\mathbb{E}[(y - \hat{f}(\mathbf{x}))^2] &= \mathbb{E}[(f(\mathbf{x}) + \epsilon - \hat{f}(\mathbf{x}))^2] \\
&= \mathbb{E}\left[\underbrace{(f(\mathbf{x}) - \mathbb{E}[\hat{f}(\mathbf{x})])}_{=A} + \underbrace{\epsilon}_{=B} + \underbrace{(\mathbb{E}[\hat{f}(\mathbf{x})] - \hat{f}(\mathbf{x}))}_{=C}\right]^2 \\
&= \mathbb{E}[A^2 + B^2 + C^2 + 2AB + 2BC + 2CA] \\
(\mathbb{E}[B] = \mathbb{E}[C] = 0) \quad &= \mathbb{E}[A^2] + \mathbb{E}[B^2] + \mathbb{E}[C^2] \\
(\mathbb{E}[\epsilon] = 0) \quad &= \mathbb{E}[(f(\mathbf{x}) - \mathbb{E}[\hat{f}(\mathbf{x})])^2] + \mathbb{E}[(\epsilon - \mathbb{E}[\epsilon])^2] + \mathbb{E}[(\hat{f}(\mathbf{x}) - \mathbb{E}[\hat{f}(\mathbf{x})])^2] \\
(\mathbb{E}[f(\mathbf{x})] = f(\mathbf{x})) \quad &= (\mathbb{E}[\hat{f}(\mathbf{x})] - f(\mathbf{x}))^2 + \text{Var}[\hat{f}(\mathbf{x})] + \text{Var}[\epsilon] \\
&= \text{Biais}[\hat{f}(\mathbf{x})]^2 + \text{Var}[\hat{f}(\mathbf{x})] + \text{Var}[\epsilon]
\end{aligned}$$

où

$$\text{Biais}[\hat{f}(\mathbf{x})] := \mathbb{E}[\hat{f}(\mathbf{x})] - f(\mathbf{x}) \quad \text{Var}[\hat{f}(\mathbf{x})] := \mathbb{E}[(\hat{f}(\mathbf{x}) - \mathbb{E}[\hat{f}(\mathbf{x})])^2]$$

BIAS-VARIANCE TRADE-OFF

On a:

$$\begin{aligned}
 \mathbb{E}[(y - \hat{f}(\mathbf{x}))^2] &= \mathbb{E}[(f(\mathbf{x}) + \epsilon - \hat{f}(\mathbf{x}))^2] \\
 &= \mathbb{E}\left[\underbrace{(f(\mathbf{x}) - \mathbb{E}[\hat{f}(\mathbf{x})])}_{=A} + \underbrace{\epsilon}_{=B} + \underbrace{(\mathbb{E}[\hat{f}(\mathbf{x})] - \hat{f}(\mathbf{x}))}_{=C}\right]^2 \\
 &= \mathbb{E}[A^2 + B^2 + C^2 + 2AB + 2BC + 2CA] \\
 (\mathbb{E}[B] = \mathbb{E}[C] = 0) \quad &= \mathbb{E}[A^2] + \mathbb{E}[B^2] + \mathbb{E}[C^2] \\
 (\mathbb{E}[\epsilon] = 0) \quad &= \mathbb{E}[(f(\mathbf{x}) - \mathbb{E}[\hat{f}(\mathbf{x})])^2] + \mathbb{E}[(\epsilon - \mathbb{E}[\epsilon])^2] + \mathbb{E}[(\hat{f}(\mathbf{x}) - \mathbb{E}[\hat{f}(\mathbf{x})])^2] \\
 (\mathbb{E}[f(\mathbf{x})] = f(\mathbf{x})) \quad &= \left(\mathbb{E}[\hat{f}(\mathbf{x})] - f(\mathbf{x})\right)^2 + \text{Var}[\hat{f}(\mathbf{x})] + \text{Var}[\epsilon] \\
 &= \text{Biais}[\hat{f}(\mathbf{x})]^2 + \text{Var}[\hat{f}(\mathbf{x})] + \text{Var}[\epsilon]
 \end{aligned}$$

où

$$\text{Biais}[\hat{f}(\mathbf{x})] := \mathbb{E}[\hat{f}(\mathbf{x})] - f(\mathbf{x}) \quad \text{Var}[\hat{f}(\mathbf{x})] := \mathbb{E}[(\hat{f}(\mathbf{x}) - \mathbb{E}[\hat{f}(\mathbf{x})])^2]$$

BIAS-VARIANCE TRADE-OFF

On a:

$$\begin{aligned}
 \mathbb{E}[(y - \hat{f}(\mathbf{x}))^2] &= \mathbb{E}[(f(\mathbf{x}) + \epsilon - \hat{f}(\mathbf{x}))^2] \\
 &= \mathbb{E}\left[\underbrace{(f(\mathbf{x}) - \mathbb{E}[\hat{f}(\mathbf{x})])}_{=A} + \underbrace{\epsilon}_{=B} + \underbrace{(\mathbb{E}[\hat{f}(\mathbf{x})] - \hat{f}(\mathbf{x}))}_{=C}\right]^2 \\
 &= \mathbb{E}[A^2 + B^2 + C^2 + 2AB + 2BC + 2CA] \\
 (\mathbb{E}[B] = \mathbb{E}[C] = 0) \quad &= \mathbb{E}[A^2] + \mathbb{E}[B^2] + \mathbb{E}[C^2] \\
 (\mathbb{E}[\epsilon] = 0) \quad &= \mathbb{E}[(f(\mathbf{x}) - \mathbb{E}[\hat{f}(\mathbf{x})])^2] + \mathbb{E}[(\epsilon - \mathbb{E}[\epsilon])^2] + \mathbb{E}[(\hat{f}(\mathbf{x}) - \mathbb{E}[\hat{f}(\mathbf{x})])^2] \\
 (\mathbb{E}[f(\mathbf{x})] = f(\mathbf{x})) \quad &= \left(\mathbb{E}[\hat{f}(\mathbf{x})] - f(\mathbf{x})\right)^2 + \text{Var}[\hat{f}(\mathbf{x})] + \text{Var}[\epsilon] \\
 &= \text{Biais}[\hat{f}(\mathbf{x})]^2 + \text{Var}[\hat{f}(\mathbf{x})] + \text{Var}[\epsilon]
 \end{aligned}$$

où

$$\text{Biais}[\hat{f}(\mathbf{x})] := \mathbb{E}[\hat{f}(\mathbf{x})] - f(\mathbf{x}) \quad \text{Var}[\hat{f}(\mathbf{x})] := \mathbb{E}[(\hat{f}(\mathbf{x}) - \mathbb{E}[\hat{f}(\mathbf{x})])^2]$$

BIAS-VARIANCE TRADE-OFF

- ▶ Le biais $\text{Biais}[\hat{f}(x)]$ représente **l'erreur due à la complexité du modèle \hat{f}** .
- ▶ La variance $\text{Var}[\hat{f}(x)]$ représente **la sensibilité du modèle \hat{f} par rapport à son training set**, i.e., la variation moyenne de $\hat{f}(x)$ autour de sa moyenne $E[\hat{f}(x)]$ si le modèle \hat{f} était estimé à partir de différents training sets.
- ▶ $\text{Var}[\epsilon]$ représente **l'erreur irréductible liée au bruit inhérent à la vraie relation fonctionnelle f** .
- ▶ **Dilemme biais-variance (bias-variance trade-off)**: plus le modèle \hat{f} est complexe, plus le biais sera faible, mais plus la variance sera élevée.

BIAS-VARIANCE TRADE-OFF

- ▶ Le biais $\text{Biais}[\hat{f}(x)]$ représente **l'erreur due à la complexité du modèle \hat{f}** .
- ▶ La variance $\text{Var}[\hat{f}(x)]$ représente **la sensibilité du modèle \hat{f} par rapport à son training set**, i.e., la variation moyenne de $\hat{f}(x)$ autour de sa moyenne $E[\hat{f}(x)]$ si le modèle \hat{f} était estimé à partir de différents training sets.
- ▶ $\text{Var}[\epsilon]$ représente l'erreur irréductible liée au bruit inhérent à la vraie relation fonctionnelle f .
- ▶ Dilemme biais-variance (bias-variance trade-off): plus le modèle \hat{f} est complexe, plus le biais sera faible, mais plus la variance sera élevée.

BIAS-VARIANCE TRADE-OFF

- ▶ Le biais $\text{Biais}[\hat{f}(x)]$ représente **l'erreur due à la complexité du modèle \hat{f}** .
- ▶ La variance $\text{Var}[\hat{f}(x)]$ représente **la sensibilité du modèle \hat{f} par rapport à son training set**, i.e., la variation moyenne de $\hat{f}(x)$ autour de sa moyenne $E[\hat{f}(x)]$ si le modèle \hat{f} était estimé à partir de différents training sets.
- ▶ $\text{Var}[\epsilon]$ représente **l'erreur irréductible** liée au bruit inhérent à la vraie relation fonctionnelle f .
- ▶ Dilemme biais-variance (bias-variance trade-off): plus le modèle \hat{f} est complexe, plus le biais sera faible, mais plus la variance sera élevée.

BIAS-VARIANCE TRADE-OFF

- ▶ Le biais $\text{Biais}[\hat{f}(x)]$ représente **l'erreur due à la complexité du modèle \hat{f}** .
- ▶ La variance $\text{Var}[\hat{f}(x)]$ représente **la sensibilité du modèle \hat{f} par rapport à son training set**, i.e., la variation moyenne de $\hat{f}(x)$ autour de sa moyenne $E[\hat{f}(x)]$ si le modèle \hat{f} était estimé à partir de différents training sets.
- ▶ $\text{Var}[\epsilon]$ représente **l'erreur irréductible** liée au bruit inhérent à la vraie relation fonctionnelle f .
- ▶ **Dilemme biais-variance (bias-variance trade-off)**: plus le modèle \hat{f} est complexe, plus le biais sera faible, mais plus la variance sera élevée.

BIAS-VARIANCE TRADE-OFF

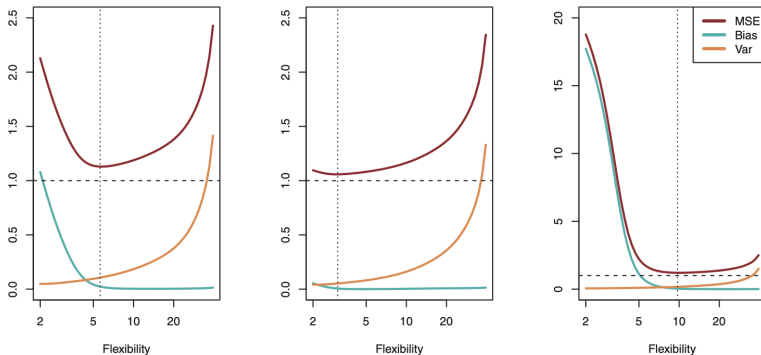


FIGURE 2.12. Squared bias (blue curve), variance (orange curve), $\text{Var}(\epsilon)$ (dashed line), and test MSE (red curve) for the three data sets in Figures 2.9–2.11. The vertical dotted line indicates the flexibility level corresponding to the smallest test MSE.

Figure taken from [James et al., 2013]

BIBLIOGRAPHIE



Fleuret, F. (2022).
Deep Learning Course.



James, G., Witten, D., Hastie, T., and Tibshirani, R. (2013).
An Introduction to Statistical Learning: with Applications in R, volume 103 of
Springer Texts in Statistics.
Springer, New York.



Wikipedia contributors (2022).
Bias–variance tradeoff — Wikipedia, the free encyclopedia.