

INTRODUCTION

Jérémie Cabessa
Laboratoire DAVID, UVSQ

INTELLIGENCE ARTIFICIELLE (IA)

- ▶ **L'intelligence artificielle (IA)** désigne l'ensemble des théories et techniques visant à simuler certains aspects de l'intelligence humaine par des systèmes informatiques.
- ▶ Ces aspects de l'intelligence comprennent entre autres: le raisonnement, l'apprentissage, la perception, etc.



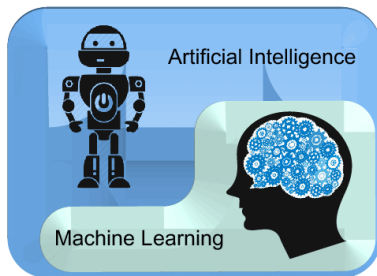
INTELLIGENCE ARTIFICIELLE (IA)

- ▶ L'**intelligence artificielle (IA)** désigne l'ensemble des théories et techniques visant à simuler certains aspects de l'intelligence humaine par des systèmes informatiques.
- ▶ Ces aspects de l'intelligence comprennent entre autres: le raisonnement, l'apprentissage, la perception, etc.



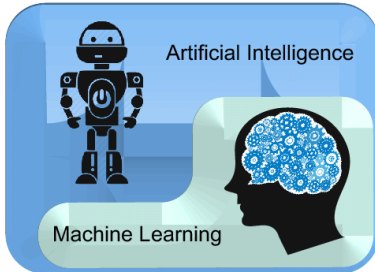
MACHINE LEARNING (ML)

- ▶ Le **machine learning (ML)** désigne l'ensemble des théories et techniques algorithmiques qui permettent d'apprendre à des systèmes informatiques à résoudre des problèmes, exécuter des tâches, etc., en se basant sur des données (data).
- ▶ Le machine learning (ML) est un sous-domaine de l'intelligence artificielle (IA).



MACHINE LEARNING (ML)

- ▶ Le **machine learning (ML)** désigne l'ensemble des théories et techniques algorithmiques qui permettent d'apprendre à des systèmes informatiques à résoudre des problèmes, exécuter des tâches, etc., en se basant sur des données (data).
- ▶ Le machine learning (ML) est un sous-domaine de l'intelligence artificielle (IA).



MACHINE LEARNING (ML)

- ▶ Le machine learning constitue un **changement de paradigme** dans la résolution de problèmes.
- ▶ **Exemple:** on aimerait créer un algorithme de reconnaissance d'images qui détecte les images contenant des fleurs.
 - ▶ Approche "hard-coded": se base sur notre compréhension du concept de "fleur", programmé en dur.
 - On code les formes, couleurs, etc., que peuvent avoir les différentes fleurs, et on crée un algorithme qui détecte les fleurs à partir de ces caractéristiques.
 - ▶ Approche "machine learning": on base sur des données d'exemples de "fleurs" et de "non-fleurs" et l'algorithme apprend à détecter les fleurs à partir des données.

MACHINE LEARNING (ML)

- ▶ Le machine learning constitue un **changement de paradigme** dans la résolution de problèmes.
- ▶ **Exemple:** on aimerait créer un algorithme de reconnaissance d'images qui détecte les images contenant des fleurs.

- ▶ **Approche “hard-coded”:** se base sur notre compréhension du concept de “fleur”, programmé en dur.

On code les formes, couleurs, etc., que peuvent avoir les différentes fleurs, et on crée un algorithme qui détecte les fleurs à partir de ces caractéristiques.

- ▶ **Approche “machine learning”:** se base sur des data, i.e., des images de “fleurs” et de “non-fleurs” utilisées comme exemples.

On crée **un algorithme qui apprend par lui-même** à détecter les fleurs à partir des data.

MACHINE LEARNING (ML)

- ▶ Le machine learning constitue un **changement de paradigme** dans la résolution de problèmes.
- ▶ **Exemple:** on aimerait créer un algorithme de reconnaissance d'images qui détecte les images contenant des fleurs.

- ▶ **Approche “hard-coded”:** se base sur notre compréhension du concept de “fleur”, programmé en dur.

On code les formes, couleurs, etc., que peuvent avoir les différentes fleurs, et on crée un algorithme qui détecte les fleurs à partir de ces caractéristiques.

- ▶ **Approche “machine learning”:** se base sur des data, i.e., des images de “fleurs” et de “non-fleurs” utilisées comme exemples.

On crée **un algorithme qui apprend par lui-même** à détecter les fleurs à partir des data.

MACHINE LEARNING (ML)

- ▶ Le machine learning constitue un **changement de paradigme** dans la résolution de problèmes.
- ▶ **Exemple:** on aimerait créer un algorithme de reconnaissance d'images qui détecte les images contenant des fleurs.
 - ▶ **Approche “hard-coded”:** se base sur notre compréhension du concept de “fleur”, programmé en dur.

On code les formes, couleurs, etc., que peuvent avoir les différentes fleurs, et on crée un algorithme qui détecte les fleurs à partir de ces caractéristiques.
 - ▶ **Approche “machine learning”:** se base sur des data, i.e., des images de “fleurs” et de “non-fleurs” utilisées comme exemples.

On crée **un algorithme qui apprend par lui-même** à détecter les fleurs à partir des data.

APPRENTISSAGE SUPERVISÉ ET NON-SUPERVISÉ

- ▶ Deux types de méthodes d'apprentissage:
 - apprentissage supervisé (supervised learning)
 - apprentissage non-supervisé (unsupervised learning)

- ▶ Apprentissage supervisé

Les données contiennent des variables d'input et des variables d'output. Il s'agit de modéliser la relation entre les variables d'input et la variable d'output (prediction, classification).

- ▶ Apprentissage non-supervisé

Les données ne contiennent que des variables d'input, et pas de variable d'output. Il s'agit d'extraire des classes ou groupes de données présentant des caractéristiques communes (clustering).

APPRENTISSAGE SUPERVISÉ ET NON-SUPERVISÉ

- ▶ Deux types de méthodes d'apprentissage:
 - apprentissage supervisé (supervised learning)
 - apprentissage non-supervisé (unsupervised learning)

- ▶ **Apprentissage supervisé**

Les données contiennent des variables d'input et des variables d'output. Il s'agit de modéliser la relation entre les variables d'input et la variable d'output (prediction, classification).

- ▶ Apprentissage non-supervisé

Les données ne contiennent que des variables d'input, et pas de variable d'output. Il s'agit d'extraire des classes ou groupes de données présentant des caractéristiques communes (clustering).

APPRENTISSAGE SUPERVISÉ ET NON-SUPERVISÉ

- ▶ Deux types de méthodes d'apprentissage:
 - apprentissage supervisé (supervised learning)
 - apprentissage non-supervisé (unsupervised learning)

- ▶ **Apprentissage supervisé**

Les données contiennent des variables d'input et des variables d'output. Il s'agit de modéliser la relation entre les variables d'input et la variable d'output (prediction, classification).

- ▶ **Apprentissage non-supervisé**

Les données ne contiennent que des variables d'input, et pas de variable d'output. Il s'agit d'extraire des classes ou groupes de données présentant des caractéristiques communes (clustering).

APPRENTISSAGE SUPERVISÉ

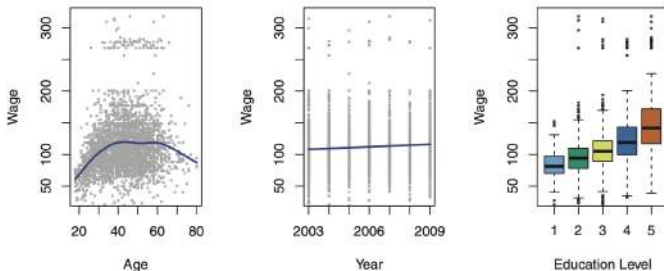


FIGURE 1.1. *Wage data, which contains income survey information for males from the central Atlantic region of the United States. Left: **wage** as a function of **age**. On average, **wage** increases with **age** until about 60 years of age, at which point it begins to decline. Center: **wage** as a function of **year**. There is a slow but steady increase of approximately \$10,000 in the average **wage** between 2003 and 2009. Right: Boxplots displaying **wage** as a function of **education**, with 1 indicating the lowest level (no high school diploma) and 5 the highest level (an advanced graduate degree). On average, **wage** increases with the level of education.*

Figure taken from [James et al., 2013].

APPRENTISSAGE SUPERVISÉ

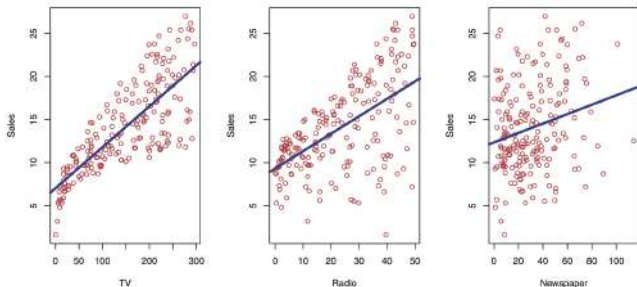


FIGURE 2.1. The Advertising data set. The plot displays sales, in thousands of units, as a function of TV, radio, and newspaper budgets, in thousands of dollars, for 200 different markets. In each plot we show the simple least squares fit of sales to that variable, as described in Chapter 3. In other words, each blue line represents a simple model that can be used to predict sales using TV, radio, and newspaper, respectively.

Figure taken from [James et al., 2013].

APPRENTISSAGE SUPERVISÉ

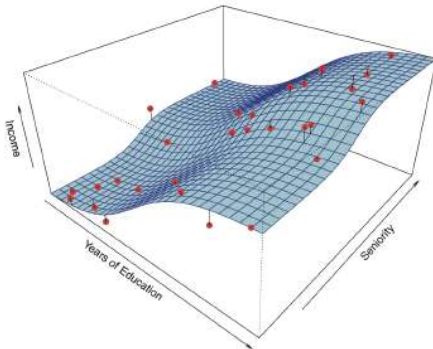


FIGURE 2.3. The plot displays **income** as a function of **years of education** and **seniority** in the **Income** data set. The blue surface represents the true underlying relationship between **income** and **years of education** and **seniority**, which is known since the data are simulated. The red dots indicate the observed values of these quantities for 30 individuals.

Figure taken from [James et al., 2013].

APPRENTISSAGE NON-SUPERVISÉ

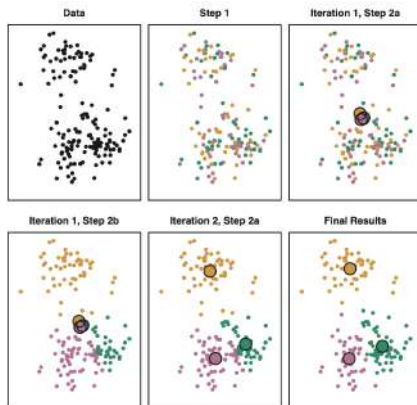


FIGURE 10.6. The progress of the K-means algorithm on the example of Figure 10.5 with $K=3$. Top left: the observations are shown. Top center: in Step 1 of the algorithm, each observation is randomly assigned to a cluster. Top right: in Step 2(a), the cluster centroids are computed. These are shown as large colored disks. Initially the centroids are almost completely overlapping because the initial cluster assignments were chosen at random. Bottom left: in Step 2(b), each observation is assigned to the nearest centroid. Bottom center: Step 2(a) is once again performed, leading to new cluster centroids. Bottom right: the results obtained after ten iterations.

Figure taken from [James et al., 2013].

RÉGRESSION ET CLASSIFICATION

- ▶ Dans le cadre de l'**apprentissage supervisé**, on distingue deux types de méthodes:
- ▶ Méthodes de régression
La variable d'output est **quantitative**.
- ▶ Méthodes de classification
La variable d'output est **qualitative**.

RÉGRESSION ET CLASSIFICATION

- ▶ Dans le cadre de l'**apprentissage supervisé**, on distingue deux types de méthodes:
- ▶ **Méthodes de régression**
La variable d'output est **quantitative**.
- ▶ Méthodes de classification
La variable d'output est **qualitative**.

RÉGRESSION ET CLASSIFICATION

- ▶ Dans le cadre de l'**apprentissage supervisé**, on distingue deux types de méthodes:
- ▶ **Méthodes de régression**
La variable d'output est **quantitative**.
- ▶ **Méthodes de classification**
La variable d'output est **qualitative**.

RÉGRESSION

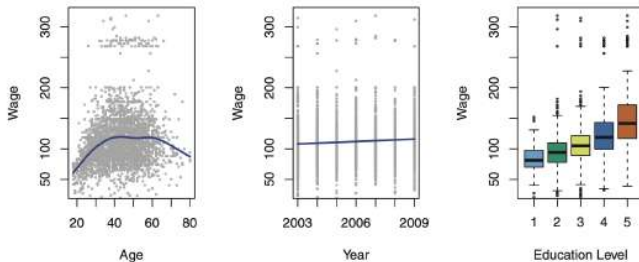


FIGURE 1.1. Wage data, which contains income survey information for males from the central Atlantic region of the United States. Left: wage as a function of age. On average, wage increases with age until about 60 years of age, at which point it begins to decline. Center: wage as a function of year. There is a slow but steady increase of approximately \$10,000 in the average wage between 2003 and 2009. Right: Boxplots displaying wage as a function of education, with 1 indicating the lowest level (no high school diploma) and 5 the highest level (an advanced graduate degree). On average, wage increases with the level of education.

Figure taken from [James et al., 2013].

CLASSIFICATION

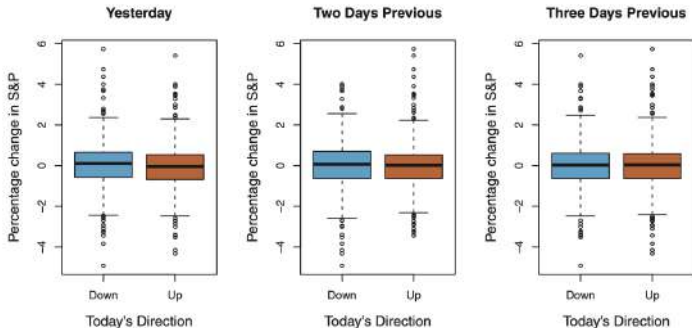


FIGURE 1.2. Left: *Boxplots of the previous day's percentage change in the S&P index for the days for which the market increased or decreased, obtained from the Smarket data.* Center and Right: *Same as left panel, but the percentage changes for 2 and 3 days previous are shown.*

Figure taken from [James et al., 2013].

CLASSIFICATION

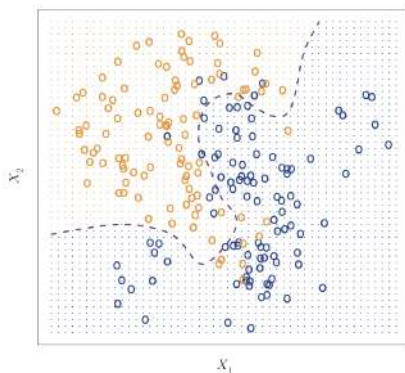


FIGURE 2.13. A simulated data set consisting of 100 observations in each of two groups, indicated in blue and in orange. The purple dashed line represents the Bayes decision boundary. The orange background grid indicates the region in which a test observation will be assigned to the orange class, and the blue background grid indicates the region in which a test observation will be assigned to the blue class.

Figure taken from [James et al., 2013].

APPRENTISSAGE SUPERVISÉ: FORMULATION DU PROBLÈME

- ▶ 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).

APPRENTISSAGE SUPERVISÉ: FORMULATION DU PROBLÈME

- ▶ 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)**.

APPRENTISSAGE SUPERVISÉ: FORMULATION DU PROBLÈME

- ▶ 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)**.

APPRENTISSAGE SUPERVISÉ: FORMULATION DU PROBLÈME

- ▶ 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 .

APPRENTISSAGE SUPERVISÉ: FORMULATION DU PROBLÈME

- ▶ 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 .

APPRENTISSAGE SUPERVISÉ: FORMULATION DU PROBLÈME

- 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).

APPRENTISSAGE SUPERVISÉ: FORMULATION DU PROBLÈME

- 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).

APPRENTISSAGE SUPERVISÉ: FORMULATION DU PROBLÈME

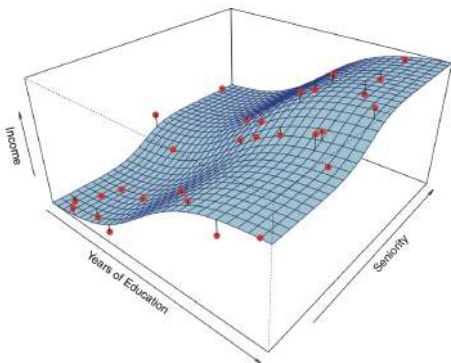


Figure taken from [James et al., 2013]

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.

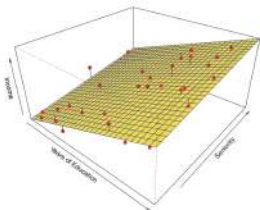


FIGURE 2.4. A linear model fit by least squares to the `Income` data from Figure 2.3. The observations are shown in red, and the yellow plane indicates the least squares fit to the data.

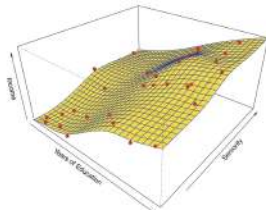


FIGURE 2.5. A smooth thin-plate spline fit to the `Income` data from Figure 2.3 is shown in yellow; the observations are displayed in red. Splines are discussed in Chapter 7.

Figure taken from [James et al., 2013]

ENSEMBLES DE TRAIN ET DE TEST

(TRAIN/TEST SETS)

- ▶ Le machine learning se base sur des données (data). L'ensemble de toutes les données dont on dispose au départ d'appelle un **dataset**
- ▶ Les données utilisées pour entraîner le modèle s'appellent **ensemble de train** ou **train set**.
- ▶ Pour évaluer le modèle, on utilise un ensemble de données distinct du train set, des données que le modèle n'a jamais vues auparavant, qu'on appelle **ensemble de test** ou **test set**.

ENSEMBLES DE TRAIN ET DE TEST

(TRAIN/TEST SETS)

- ▶ Le machine learning se base sur des données (data). L'ensemble de toutes les données dont on dispose au départ s'appelle un **dataset**
- ▶ Les données utilisées pour entraîner le modèle s'appellent **ensemble de train** ou **train set**.
- ▶ Pour évaluer le modèle, on utilise un ensemble de données distinct du train set, des données que le modèle n'a jamais vues auparavant, qu'on appelle **ensemble de test** ou **test set**.

ENSEMBLES DE TRAIN ET DE TEST

(TRAIN/TEST SETS)

- ▶ Le machine learning se base sur des données (data). L'ensemble de toutes les données dont on dispose au départ d'appelle un **dataset**
- ▶ Les données utilisées pour entraîner le modèle s'appellent **ensemble de train** ou **train set**.
- ▶ Pour évaluer le modèle, on utilise un ensemble de données distinct du train set, des données que le modèle n'a jamais vues auparavant, qu'on appelle **ensemble de test** ou **test set**.

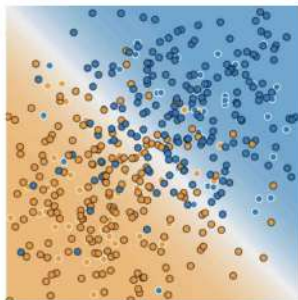
ENSEMBLES DE TRAIN ET DE TEST (TRAIN/TEST SETS)

- ▶ En général, on sépare de **dataset** de départ en prenant 80% des données pour le **train set** et 20% pour le **test set** (train-test split).

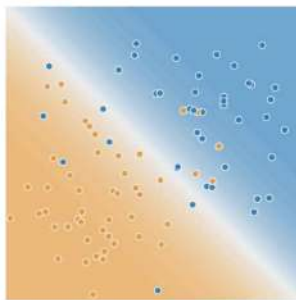
Dataset



ENSEMBLES DE TRAIN ET DE TEST (TRAIN/TEST SETS)



Training Data



Test Data

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 moindre 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

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

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

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

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, le **test set**.
- ▶ 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, le **test set**.
- ▶ 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_{test} élevée).
- ▶ Le modèle a donc "sur-appri" (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_{test} 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_{test} 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

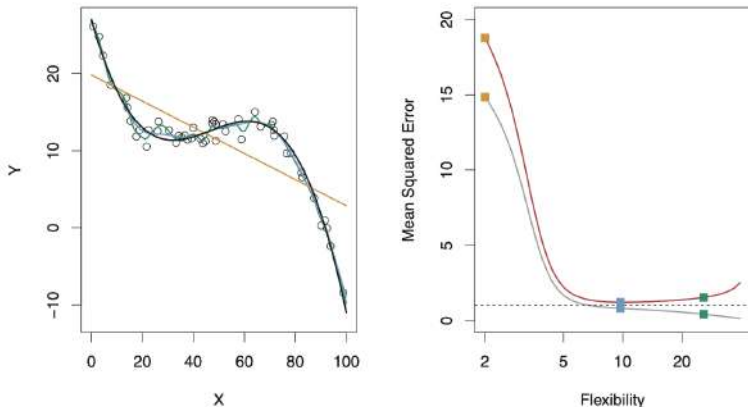


FIGURE 2.11. Details are as in Figure 2.9, using a different f that is far from linear. In this setting, linear regression provides a very poor fit to the data.

Figure taken from [James et al., 2013]

BIBLIOGRAPHIE

Most images taken from [[James et al., 2013](#)].



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.