



# Module 6

## Algorithmes de base en apprentissage machine

### Sommaire

<b>6.1</b>	<b>Principe de l'apprentissage machine</b>	<b>3</b>
6.1.1	Entraînement, test et validation du modèle	4
6.1.2	Validation par $K$ -fold	4
<b>6.2</b>	<b>Techniques d'apprentissage machine</b>	<b>7</b>
<b>6.3</b>	<b>Classification</b>	<b>9</b>
6.3.1	Réseau de neurones	9
6.3.2	Évaluation d'un système de classification	10
<b>6.4</b>	<b>Regroupement</b>	<b>13</b>
6.4.1	Algorithme $K$ -moyennes	14
6.4.2	Évaluation de la qualité du regroupement : Indice de silhouette	14
<b>6.5</b>	<b>Régression</b>	<b>16</b>
6.5.1	Modèle de régression univarié	16
6.5.2	Modèle de régression multivariée	18
<b>6.6</b>	<b>Règles d'association</b>	<b>18</b>

6.6.1	Support et confiance . . . . .	19
6.6.2	Algorithme Apriori . . . . .	19

---

Dernière mise à jour le 24 septembre 2019

# Introduction

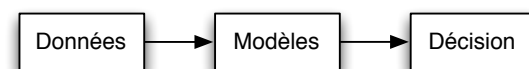
L'apprentissage machine (aussi appelé apprentissage artificiel ou automatique, en anglais *machine learning*) est ...

*le processus par lequel un ordinateur acquiert de nouvelles connaissances et améliore son mode de fonctionnement en tenant compte des résultats obtenus lors de traitements antérieurs. (Office québécois de la langue française)*

Il existe deux approches principales en apprentissage machine. La première est issue de l'intelligence artificielle *syntactique* ou *symbolique*. Elle est fondée sur la modélisation du raisonnement logique et sur la représentation et la manipulation de la connaissance par des symboles formels. La deuxième est issue l'intelligence artificielle *statistique*, aussi parfois appelée numérique parce que, souvent, la représentation des données est sous une forme numérique. Le cours SCI 1016 s'intéresse à l'apprentissage machine statistique.

## 6.1 Principe de l'apprentissage machine

La définition présentée en introduction permet d'enchaîner avec les éléments suivants sur le principe de l'apprentissage machine : il s'agit d'un ensemble de méthodes qui permettent de construire un modèle de la réalité à partir de données, soit en améliorant un modèle existant moins général, soit en créant un nouveau modèle représentatif de nouvelles données. Le choix du modèle dépend essentiellement des données à analyser. Les paramètres du modèle sont déterminés durant la phase d'apprentissage en utilisant l'algorithme qui lui est spécifique. Les modèles servent souvent à prendre des décisions. La décision à prendre dépend essentiellement de la problématique de l'apprentissage machine à résoudre. Par exemple, s'il s'agit d'un problème de classification de forme, la décision consiste souvent à reconnaître la classe de la forme. Par contre s'il s'agit d'un problème de regroupement, la décision consistera à identifier les différents regroupements.



**Figure 6.1:** Neurone formel.

### 6.1.1 Entraînement, test et validation du modèle

L'ensemble des données considérées pour une analyse par apprentissage machine peut être réparti en trois sous-ensembles :

- Un ensemble d'entraînement qui est utilisé pour l'apprentissage (entraînement) des paramètres du modèle.
- Un ensemble de validation. Il s'agit d'un ensemble de données qui permet d'évaluer le modèle pendant la phase d'entraînement. Cette étape, appelée validation du modèle peut être omise en passant à la phase de test directement.
- Un ensemble de test. Une fois le modèle construit à partir de l'ensemble d'entraînement, le modèle est évalué en utilisant un ensemble de test : un ensemble d'échantillons n'ayant pas servi pour l'apprentissage.

### 6.1.2 Validation par $K$ -fold

Il existe différentes méthodes pour réaliser le partage des données en un ensemble de données d'apprentissage et un ensemble de données de test. Parmi ces méthodes on trouve la validation par  $K$ -plis (plus connu sous le terme  $K$ -fold).

Le principe de cette méthode de validation consiste à diviser l'échantillon original en  $K$  échantillons de même taille et à prendre un échantillon pour procéder à la validation. Le processus est répété jusqu'à parcourir des  $K$  échantillons. En d'autres termes, (1) nous divisons l'échantillon original en  $K$  échantillons, (2) nous sélectionnons un des  $K$  échantillons comme ensemble de tests et (3) nous gardons les  $(K - 1)$  autres échantillons pour réaliser l'apprentissage et la conception du système de classification. À la fin, (4) nous calculons la moyenne des résultats de validations pour avoir un seul résultat.

La répartition la plus communément utilisée entre ces échantillons est une proportion de  $2/3$  pour l'apprentissage et  $1/3$  pour le test, c'est-à-dire avec un  $K = 3$  fold. Dans le cas de grande taille de base de données, la validation 10-fold est plus appropriée.

- **Exemple 6.1** Nous considérons ici la base de données IRIS qui contient 150 échantillons provenant de trois classes décrites par quatre variables. nous désirons faire une validation croisée 10-fold (9 fold pour l'entraînement et 1 pour le test). Les données étant ordonnés par classes, il sera donc nécessaires de les mélanger aléatoirement (Ligne 3 du code suivant) avant d'extraire les données d'entraînement et de test.

```
1  # Exemple 6.1
2  # Mélanger les données aléatoirement
3  MyData<-iris
4  MyData<-MyData[sample(nrow(MyData)),]
5  # Créer 10 folds égaux
6  folds <- cut(seq(1,nrow(MyData)),breaks=10,labels=FALSE, seed=10)
7  for(i in 1:10){
8      testIndexes <- which(folds==i,arr.ind=TRUE)
9      testData <- MyData[testIndexes , ]
10     trainData <- MyData[-testIndexes , ]
11 }
12 summary(testData)
13 summary(trainData)
```

Nous obtenons deux ensembles `testData` et `trainData` contenant respectivement 15 et 135 échantillons (soit un total de 150 échantillons).

Nous pouvons remarquer que le nombre d'échantillons par classe est différents d'une classe à une autre. Par exemple dans `testData` nous avons respectivement 6, 4 et 5 pour les classes `setosa` et `versicolor` et `virginica`. Ceci est dû à l'ordonnancement aléatoire des échantillons dans la base de données. Nous aurions pu faire une division par classe afin d'obtenir des échantillons équitables par classe. Il est à noter que deux exécutions successives pourraient entraîner des partitions différentes des données dû à l'ordonnancement aléatoire des échantillons. Pour obtenir les mêmes partitions entre deux exécutions successives, nous fixons la valeur du `seed`.

La fonction `summary` est utilisée pour avoir une description des données et la fonction `dim` pour avoir la dimension.

```

> summary(testData)
Sepal.Length Sepal.Width Petal.Length Petal.Width Species
Min.      :4.400 Min.      :2.200 Min.      :1.200 Min.      :0.20 setosa      :6
1st Qu.:5.250 1st Qu.:2.600 1st Qu.:1.400 1st Qu.:0.30 versicolor:4
Median :5.800 Median :2.900 Median :4.000 Median :1.30 virginica  :5
Mean   :5.713 Mean   :2.967 Mean   :3.433 Mean   :1.12
3rd Qu.:6.250 3rd Qu.:3.200 3rd Qu.:5.050 3rd Qu.:1.70
Max.   :6.900 Max.   :4.000 Max.   :5.500 Max.   :2.40

> summary(trainData)
Sepal.Length Sepal.Width Petal.Length Petal.Width Species
Min.      :4.300 Min.      :2.000 Min.      :1.000 Min.      :0.100 setosa      :44
1st Qu.:5.100 1st Qu.:2.800 1st Qu.:1.600 1st Qu.:0.300 versicolor:46
Median :5.800 Median :3.000 Median :4.400 Median :1.300 virginica  :45
Mean   :5.858 Mean   :3.067 Mean   :3.794 Mean   :1.208
3rd Qu.:6.400 3rd Qu.:3.350 3rd Qu.:5.100 3rd Qu.:1.800
Max.   :7.900 Max.   :4.400 Max.   :6.900 Max.   :2.500

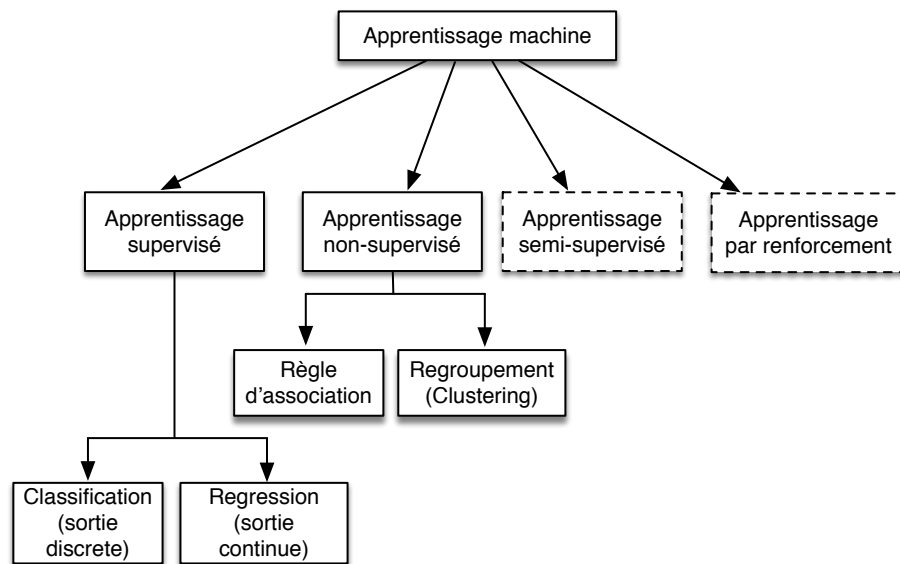
> dim(testData)
[1] 15  5

> dim(trainData)
[1] 135  5

```

## 6.2 Techniques d'apprentissage machine

Les algorithmes d'apprentissage sont catégorisés selon les techniques d'apprentissage qu'ils emploient. Nous en citons l'apprentissage supervisé, l'apprentissage non supervisé, l'apprentissage par renforcement et l'apprentissage semi-supervisé.



**Figure 6.2:** Principe de l'apprentissage machine.

## Apprentissage supervisé

Dans le cas d'un apprentissage supervisé (en anglais *supervised learning*), le système observe des couples de types entrée-sortie et apprend une fonction (un modèle) qui permet d'aboutir à la sortie à partir de l'entrée. Cette phase est appelée *phase d'apprentissage* ou *d'entraînement*. C'est en ce sens que l'apprentissage est appelé supervisé, métaphore qui signifie qu'un professeur apprend au système la sortie à fournir pour chaque entrée. Les données d'entrée et les données de sortie correspondantes, aussi appelées classes, sont connues (aussi dites labellisées ou étiquetées). Telles que décrits précédemment, elles sont regroupées dans un ensemble de données appelées *données d'apprentissage* ou *d'entraînement* qui se présentent sous la forme de couples  $(\mathbf{x}_i, y_i)_{1 \leq i \leq N}$  avec  $\mathbf{x}_i \in \mathbb{R}$  un

vecteur de caractéristique et  $y_i$  la classe (ou étiquette ou attribut) correspondante.  $N$  étant le nombre d'échantillons.

### **Apprentissage non supervisé**

Contrairement à l'apprentissage supervisé, dans le cas non supervisé les données de sortie ne sont pas connues. Le système apprend alors de lui-même à organiser les données ou à déterminer des structures dans les données. La tâche d'apprentissage la plus courante est le regroupement (*clustering* en anglais) qui consiste à regrouper les données d'entrées selon leurs caractéristiques communes. Ce type d'apprentissage est utilisé dans le but de visualiser ou explorer des données.

### **Apprentissage semi supervisé**

L'apprentissage semi-supervisé se base sur un mélange de données étiquetées et non-étiquetées. Ceci permet, d'une part, d'améliorer la qualité de l'apprentissage, et d'autre part, de réduire le temps de préparation des données pour leur étiquetage.

### **Apprentissage par renforcement**

L'apprentissage par renforcement se base sur des données d'entrée similaires à celles utilisées en apprentissage supervisé. Cependant dans ce cas, l'apprentissage est guidé par l'environnement sous la forme de récompenses (positive ou négative) calculées en fonction de l'erreur commise lors de l'apprentissage. En robotique, l'apprentissage par renforcement a permis de mettre au point des robots plus autonomes et adaptatifs à leurs environnements.

Le cours SCI1016 présente en détails les deux techniques les plus connues en apprentissage à savoir l'apprentissage supervisé et l'apprentissage non supervisé.

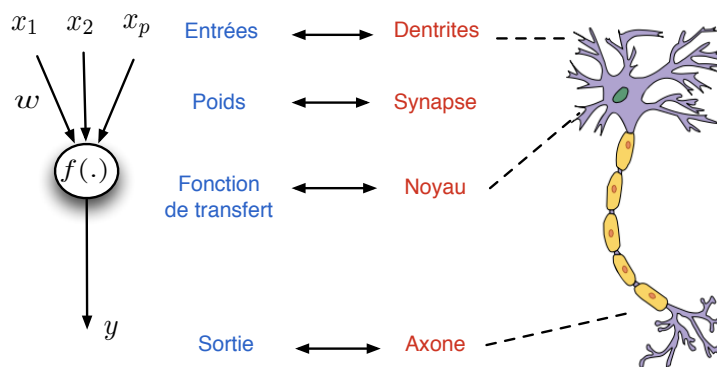


## 6.3 Classification

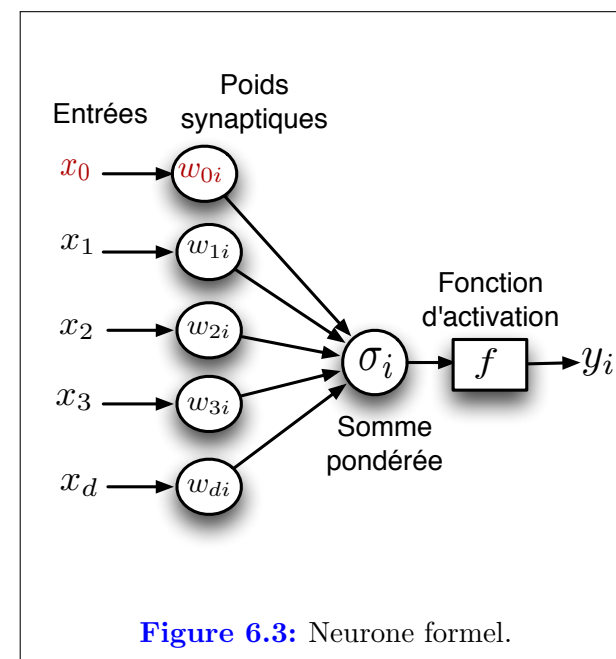
La reconnaissance ou classification de formes (en anglais, *pattern recognition*) est « l'ensemble des techniques permettant à l'ordinateur de détecter la présence de formes, en comparant leurs caractéristiques avec celles de motifs de référence » (Office québécois de la langue française, 2010). La classification de formes est l'une des branches de l'intelligence artificielle qui fait largement appel aux techniques d'apprentissage machine ; plus particulièrement aux réseaux de neurones.

### 6.3.1 Réseau de neurones

Un *réseau de neurones* aussi appelé *réseau de neurones artificiels* est un modèle mathématique très simple dérivé du neurone biologique dont les composantes élémentaires sont des unités de calcul qui reçoivent des entrées pour produire des sorties. Ce modèle reprend les principes du fonctionnement du neurone biologique, en particulier par la sommation des entrées et la pondération de la somme de ses entrées par des *poids synaptiques* (aussi appelés coefficients synaptiques) tels que illustrés à la Figure 6.4.



**Figure 6.4:** Analogie entre le neurone biologique et le neurone artificiel.



**Figure 6.3:** Neurone formel.

Un réseau de neurones est défini par sa structure (c'est-à-dire le nombre de couches, le nombre de noeuds par couche et le nombre de sorties) et par la règle d'apprentissage qui permet de calculer les poids synaptiques. La règle d'apprentissage dépend du type d'apprentissage comme nous l'avons décrit précédemment.

Il existe plusieurs types de réseau de neurones tels que le perceptron simple, le perceptron multi-couche, le réseau de Kohonen. Ces réseaux se différencient aux niveaux de leur architecture et de la technique d'apprentissage sur laquelle ils se basent.

### 6.3.2 Évaluation d'un système de classification

Lorsque l'apprentissage machine a pour objectif de développer un système de classification de formes, la validation se fait en calculant le taux de bonne classification et à travers la matrice de confusion.

#### Taux de classification

Soit  $S$  un ensemble d'échantillons d'apprentissage, et  $T$  un ensemble d'échantillons de test. L'estimation du taux de bonne classification est mesurée sur l'ensemble de test selon :

$$\tau = \frac{\text{nbr bien classifié}(T)}{\text{nbr}(T)}$$

C'est-à-dire le rapport entre le nombre d'échantillons de l'ensemble de test qui sont bien classifiés par rapport au nombre total d'échantillons de la base de test.

Le taux de bonne classification est généralement donné en pourcentage. Il lui correspond une valeur complémentaire à 100 correspondant au taux d'erreur ( $\epsilon(\%) = 100 - \tau(\%)$ ).

#### Matrice de confusion

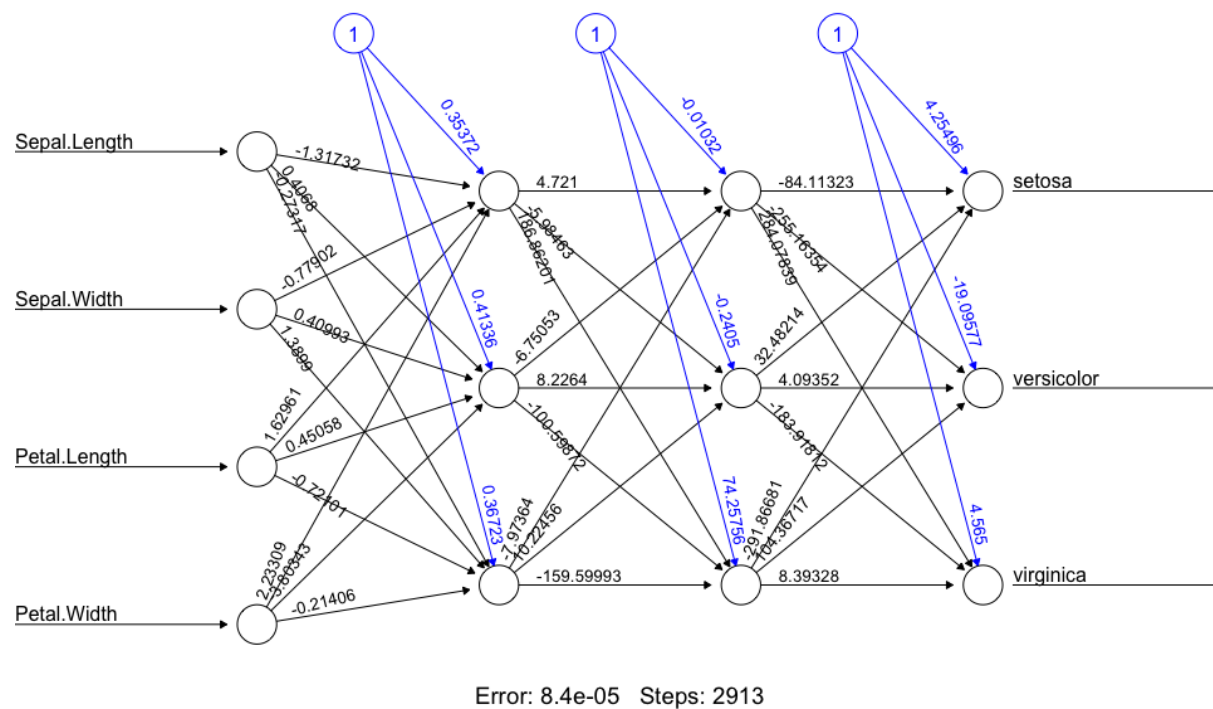
La matrice de confusion est aussi connue sous les termes matrice d'erreur, tableau de contingence ou matrice d'erreur de classification. C'est une matrice affichant les statistiques de la précision de classification et plus particulièrement les taux de classification

par classes. Généralement, L'information des lignes (données horizontales) correspond aux classes réelles des formes. Quant aux colonnes (données verticales), elles contiennent l'information prédite résultant de la classification.

Les valeurs de la diagonale de la matrice représentent le nombre de formes correctement classifié. La somme des valeurs par ligne correspond au nombre d'échantillons de test par classe. Le taux de classification par classes est donné par la valeur à la diagonale divisée par la somme des valeurs par ligne. Des exemples d'utilisation de la matrice de confusion pour l'évaluation de systèmes de classification seront donnés dans ce qui suit.

- **Exemple 6.2** Nous nous intéressons aux données IRIS. L'objectif de cet exemple est de développer un réseau de neurones qui permet la classification de données selon les différentes espèces en se basant les 4 variables descriptives : `Petal.Length`, `Petal.width`, `Setal.Length` et `Setal.width`. Pour ce faire, nous avons utilisé le package `neuralnet`. La commande `sample` (ligne 7) permet de générer des échantillons aléatoires. `hidden` permet de fixer le nombre de neurone cachés. Dans notre cas, nous avons choisi (3,3), c'est-à-dire deux couches contenant chacune trois neurones. Finalement la commande `plot` permet d'afficher l'architecture du réseau de neurone (Figure 6.5).

```
1 library(neuralnet)
2 data(iris)
3 iris$setosa <- iris$Species=="setosa"
4 iris$virginica <- iris$Species == "virginica"
5 iris$versicolor <- iris$Species == "versicolor"
6 iris.train.idx <- sample(x = nrow(iris), size = nrow(iris)*0.5)
7 iris.train <- iris[iris.train.idx,]
8 iris.test <- iris[-iris.train.idx,]
9 # Entraînement du réseau de neurone
10 iris.net <- neuralnet(setosa+virginica+versicolor~
11     Sepal.Length + Sepal.Width + Petal.Length + Petal.Width,
12     data=iris.train, hidden=c(3,3), rep = 5, err.fct = "ce",
13     linear.output = F, lifesign = "minimal", stepmax = 1000000,
14     threshold = 0.001)
15 plot(iris.net, rep="best")
```



**Figure 6.5:** Architecture du réseau de neurone

Une fois le réseau de neurone développé, nous procédons à une prédiction en utilisant le modèle développé (Ligne 2 du code suivant). La commande `summary` permet de voir le nombre d'échantillons dans la base de test par classe (`setosa`: 23, `versicolor` : 25, `virginica`: 27). La diagonale de la matrice de confusion donné par la commande `table`, permet de calculer le taux de classification par classe (`setosa`: 23/23, `versicolor` : 24/25, `virginica`: 26/27).

Notons que nous avons utilisé `iris.test[-5:-8]` pour supprimer les colonnes 5 à 8 qui contiennent les classes (étiquettes) des échantillons de test. L'objectif étant de tester les performances du réseau de neurones sur de nouvelles données.

```

1  # Prédiction
2  iris.prediction <- compute(iris.net, iris.test[-5:-8])
3  idx <- apply(iris.prediction$net.result, 1, which.max)
4  predicted <- c('setosa', 'virginica', 'versicolor')[idx]
5  summary(iris.test)
6  table(predicted, iris.test$Species)

```

```

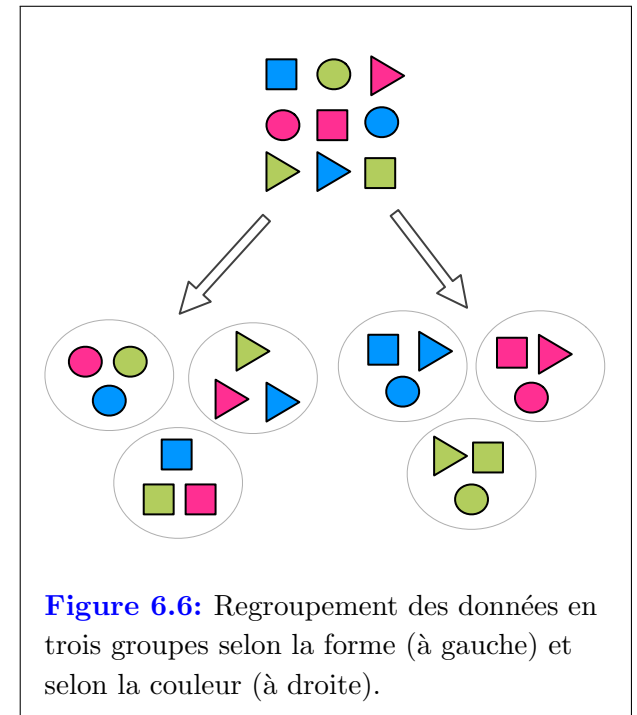
> summary(iris.test)
Species      setosa      virginica
setosa      :23      Mode :logical  Mode :logical
versicolor:25      FALSE:52      FALSE:48
virginica   :27      TRUE :23      TRUE :27
> table(predicted, iris.test$Species)
predicted      setosa versicolor virginica
setosa         23         0         0
versicolor      0         24         1
virginica       0         1         26

```

## 6.4 Regroupement

Le *regroupement*, aussi appelé *agrégation* (*clustering*), est une méthode d'apprentissage machine non supervisée, qui a pour objectif de former des groupes ou agrégats (*clusters*) d'objets similaires à partir d'un ensemble hétérogène d'objets (aussi appelées individus ou points).

Formellement, soit un ensemble de données formé de  $N$  objets décrits par  $P$  variables. Le regroupement consiste à trouver une répartition de ces données en  $K$  groupes de sorte que les objets à l'intérieur du même groupe soient similaires. Le nombre de groupes  $K$  peut être déterminé a priori.



### 6.4.1 Algorithme $K$ –moyennes

L'algorithme  $K$ –moyennes ( $K$ –*means*) utilise une technique d'affinement itératif qui consiste à améliorer progressivement la qualité des groupes selon un indice de similarité entre les différents individus. Le paramètre d'entrée de l'algorithme  $K$ –moyennes est le nombre de groupes désiré  $K$ . La distance euclidienne est souvent utilisée pour mesurer le rapprochement des groupes.

---

**Algorithme  $K$ –moyennes :** On cherche à déterminer  $K$  regroupements à partir d'un ensemble de  $N$  individus.

---

1. Choisir aléatoirement  $K$  centres initiaux  $G_i$  avec  $i = 1, 2, \dots, K$ .
  2. **Répéter**
  3. Répartir chaque individu dans le groupe  $G_i$  le plus proche.
  4. Recalculer les nouveaux centres  $G_i$ .
  5. **Jusqu'à ce que** les centres deviennent stables.
- 

### 6.4.2 Évaluation de la qualité du regroupement : Indice de silhouette

L'indice de silhouette est un indicateur efficace pour valider une méthode de regroupement. Pour tout individu  $i$  de l'ensemble des données, l'indice de silhouette est défini par la formule suivante :

$$s(i) = \frac{b_i - a_i}{\max(a_i, b_i)} \quad (6.1)$$

où  $a_i$  est la dissimilarité moyenne entre l'individu  $i$  et tous les autres individus du groupe  $G_j$  auquel il appartient et  $b_i$  est le minimum des dissimilarités moyennes entre l'individu  $i$  et tous les autres individus des groupes  $G_k$  ( $k = 1, 2, \dots, c$  avec  $k \neq j$ ).

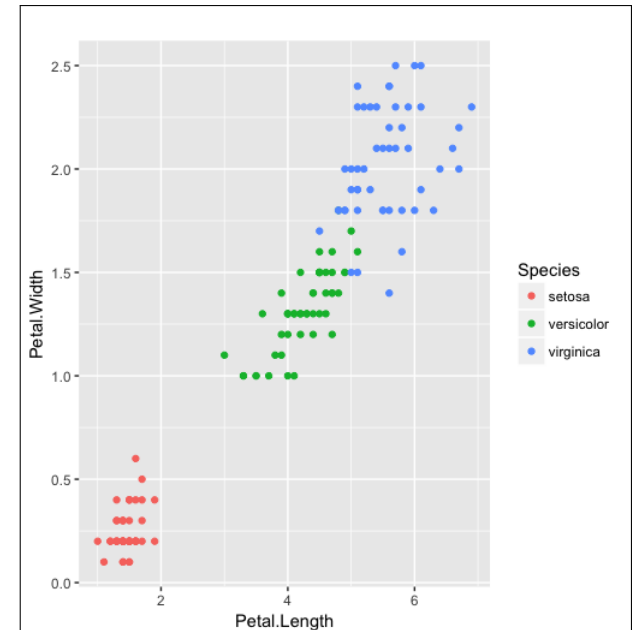
► **Exemple 6.3** Revenons aux données IRIS. L'objectif de cet exemple est de faire le regroupement des différentes espèces en se basant sur la longueur (`Petal.Length`) et la

largeur des pétales (`Petal.width`). La figure 6.7 représente la variation de la longueur (`Petal.Length`) en fonction de la largeur des pétales (`Petal.width`) pour les trois espèces. Nous pouvons déjà voir visuellement la présence d'au moins deux regroupements.

```

1  # Exemple 6.3
2  library(datasets)
3  library(ggplot2)
4  head(iris)
5  ggplot(iris, aes(Petal.Length, Petal.Width, color = Species)) + 
6    geom_point()
7  set.seed(20)
8  irisCluster <- kmeans(iris[, 3:4], 3, nstart = 20)
9  # Matrise de confusion
10 table(irisCluster$cluster, iris$Species)
11 # Representation des clusters
12 irisCluster$cluster <- as.factor(irisCluster$cluster)
13 ggplot(iris, aes(Petal.Length, Petal.Width, color = irisCluster$cluster)) + geom_point()
14 # Representation de la silhouette
15 library(cluster)
16 irisCluster <- kmeans(iris[, 3:4], 3, nstart = 20)
17 s <- silhouette(irisCluster$cluster, dist(iris[, 3:4]))
18 plot(s)

```



**Figure 6.7:** Variation de la longueur (`Petal.Length`) en fonction de la largeur des pétales (`Petal.width`).

La représentation graphique de l'indice de silhouette (Figure 6.8) montre une bonne séparation des données IRIS en trois groupes. Sur l'axe des ordonnées sont représentées les données des différents groupes selon le regroupement assigné. Sur l'axe des abscisses la valeur de l'indice de silhouette calculé selon l'équation 6.1. L'indice de silhouette est positif pour tous les échantillons. Ce qui démontre l'homogénéité des données au sein des différents groupes. Il est à noter la commande `set.seed` permet de fixer le paramètre `seed` et par conséquent d'avoir les mêmes groupes dans des simulations différentes.



**Figure 6.8:** Indice de silhouette d'un ensemble de données en 3 groupes.

## 6.5 Régression

Dans le cas d'un modèle de régression, le principe de l'apprentissage machine consiste à déterminer des modèles mathématiques qui permettent de représenter une variable à expliquer (ou variable dépendante ou variable endogène ou réponse)  $Y$ , en fonction d'une ou plusieurs variables explicatives (ou indépendantes ou variables exogènes)  $X = (x_1, x_2, \dots, x_p)$ . Le modèle de régression est la relation entre la variable à expliquer  $Y$  et les variables explicatives en entrée  $x_i$ . L'algorithme de régression permet de trouver le modèle en se basant sur les données d'entraînement. Le modèle calculé permet de donner une estimation sur les données de tests.

### 6.5.1 Modèle de régression univarié

Le modèle de régression est dit univarié ou simple s'il n'inclut qu'une seule variable explicative. Il consiste à trouver la meilleure droite qui s'approche le plus des données d'apprentissage. Le modèle obtenu est une droite de la forme :

$$Y = f(X) = \alpha + \beta x_1 \quad (6.2)$$

avec  $\alpha$  et  $\beta$  les coefficients de la droite.

► **Exemple 6.4** Nous traitons dans cet exemple la base de données disponible dans R dénommée `cars` (décrite dans l'Exemple 4.6). `cars` contient 50 voitures caractérisées chacune par deux variables : `dist` et `speed`. Nous désirons développer un modèle de régression qui permet de relier la variable prédictive `dist` et la variable à prédire `speed`.

Le résultat du modèle de régression contient le modèle `dist =  $\alpha$  +  $\beta$  speed` avec  $\alpha = -17.5791$  et  $\beta = 3.9324$ . Selon cette équation, la distance dépend linéairement de la vitesse : la distance augmente de 3.9324 pour chaque kilomètre/h de vitesse additionnelle. La valeur  $p = 1.49e^{-12}$  étant inférieure à 0.05, nous pouvons conclure qu'il y a donc une relation significative entre les variables du modèle de régression linéaire de l'ensemble de données. Le valeur du coefficient de détermination  $R^2$  permet d'évaluer la qualité de



prédiction du modèle. La valeur du coefficient de détermination ajusté  $R_{adj}^2$  prend en considération le nombre de variable. Dans le cas d'une régression simple, ce coefficient correspond au rapport entre la variance expliquée et la variance totale.

La valeur de R-squared de 0.651 est assez élevé; ce qui permet de conclure que le modèle de régression est de bonne qualité.

```

1 # Exemple 6.4
2 require(stats)
3 reg <- lm(cars$dist ~ cars$speed)
4 coeff=coefficients(reg)
5 # Equation de la droite de regression :
6 eq = paste0("y = ", round(coeff[2],1), "*x ", round(coeff[1],1))
7 # Graphes
8 plot(cars, main=eq, col="red")
9 abline(reg, col="blue")

```

```

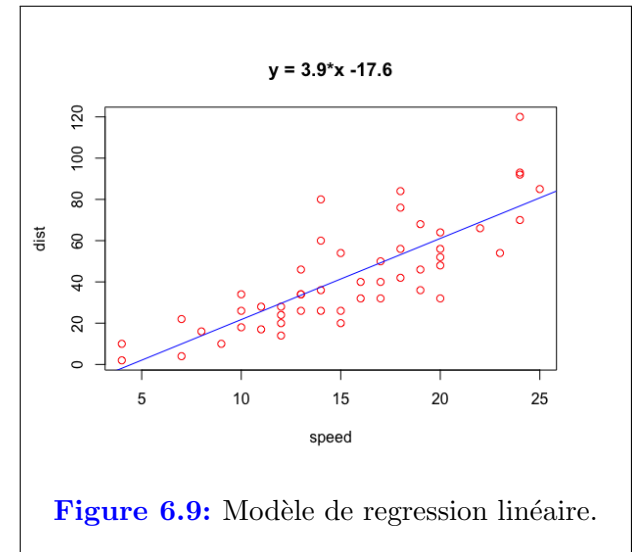
> summary(reg)

Call:
lm(formula = cars$dist ~ cars$speed)

Residuals:
Min      1Q  Median      3Q      Max
-29.069  -9.525  -2.272   9.215  43.201

Coefficients:
(Intercept) -17.5791      6.7584  -2.601   0.0123 *
cars$speed    3.9324      0.4155   9.464 1.49e-12 ***
Signif. codes:  0 *** 0.001 ** 0.01 * 0.05 . 0.1 1
Residual standard error: 15.38 on 48 degrees of freedom
Multiple R-squared:  0.6511, Adjusted R-squared:  0.6438
F-statistic: 89.57 on 1 and 48 DF,  p-value: 1.49e-12

```



**Figure 6.9:** Modèle de regression linéaire.

Nous ouvrons une parenthèse sur la relation entre le test de significativité et le modèle de régression qui permet de tester l'apport de chaque variable dans le modèle. Le modèle à estimer a pour équation  $Y = f(X) = \alpha + \beta * X$ . Le test de significativité pour le coefficient  $\beta$  est le suivant :

$$\begin{cases} H_0 : \beta = 0 \\ H_1 : \beta \neq 0 \end{cases} \quad (6.3)$$

La règle de décision est la suivante : si la valeur absolue de la statistique observée est supérieure à la valeur théorique de la Student on rejette au seuil de  $\alpha$  l'hypothèse nulle en faveur de l'hypothèse alternative. La variable est alors supprimé du modèle de régression.

### 6.5.2 Modèle de régression multivariée

Le modèle de régression est dit multivarié ou multiple s'il inclut plusieurs variables explicatives. Dans le cas de deux variables explicatives, la relation obtenue est de la forme :

$$Y = f(X) = \alpha + \beta_1 x_1 + \beta_2 x_2 \quad (6.4)$$

$\alpha$  est une constante et  $\beta_1$  et  $\beta_2$  sont les coefficients des variables explicatives.

## 6.6 Règles d'association

La recherche de règles d'association est une approche importante en apprentissage machine non-supervisé qui consiste à extraire des informations à partir de la coïncidence dans les données. Ce qui permet d'identifier des relations implicites qui sont cachés dans d'importante quantité de données. L'application la plus connue des règles d'association est le "panier de la ménagère" qui consiste à analyser des tickets de caisse dans les supermarchés afin de découvrir des corrélations d'événements transactionnels. Ces corrélations sont des relations d'implication conditionnelles entres les items d'un ensemble de données

qui se présentent sous la forme :

”Si l’item **X** existe, Alors il est possible que l’item **B** existe aussi”. Le tableau suivant illustre un exemple d’évènements transactionnels.

Client	Liste des items achetés
1	{ pain, beurre, lait }
2	{ pain, viande }
.	.....
.	.....
<i>n</i>	{ Jus de fruit, pain, beurre, fraise }

### 6.6.1 Support et confiance

Il existe deux mesures statistiques importantes associés aux règles d’association *le support* et la *confiance*.

Le support d’un ensemble d’items est la fréquence d’apparition simultanée des items figurant dans l’ensemble des données.

La confiance dans une règle  $R : X \rightarrow Y$  est définie par :

$$\begin{aligned} Conf(R) &= p(Y \subseteq T | X \subseteq T) \\ &= \frac{support(X \cup Y)}{support(X)} \end{aligned}$$

ou  $X$  et  $Y$  sont des items disjoints.

### 6.6.2 Algorithme Apriori

L’algorithme Apriori est un algorithme itératif de recherche des itemsets les plus fréquents par niveaux. Cet algorithme fonctionne en deux phases : La première consiste à déterminer

les ensembles d'items fréquents (EIF); Ensuite, on utilise ces EIF pour déterminer les règles d'association dont la confiance est supérieure au seuil fixé.

► **Exemple 6.5** Nous traitons dans cet exemple la base de données disponible dans R dénommée **Groceries** qui contient des informations sur les ventes d'une épicerie.

La base de données comprend 9835 transactions et 169 articles (**items**). L'article **whole milk** étant le plus populaire avec une vente de 2513 unités. L'utilisation de l'algorithme Apriori (**arules**) permet de ressortir 5668 règles d'association (**rules**). Les trois règles retenus sont obtenus par la commande **inspect** en fixant le niveau **n** (Ligne 7 du code).

```
1 library(arules)
2 library("arulesViz")
3 data("Groceries")
4 summary(Groceries)
5 rules <- apriori(Groceries, parameter=list(support=0.001, confidence=0.5))
6 rules
7 inspect(head(rules, n = 3, by = "lift"))
```

```

> summary(Groceries)
transactions as itemMatrix in sparse format with
9835 rows (elements/itemsets/transactions) and
169 columns (items) and a density of 0.02609146

most frequent items:
whole milk other vegetables    rolls/buns    soda    yogurt    (Other)
2513          1903          1809          1715    1372    34055

> rules
set of 5668 rules

> inspect(head(rules, n = 3, by = "lift"))
lhs                rhs                support    confidence lift    count
[1] {Instant food products,soda} => {hamburger meat} 0.001220132 0.6315789 18.99565 12
[2] {soda,popcorn}           => {salty snack}    0.001220132 0.6315789 16.69779 12
[3] {flour,baking powder}    => {sugar}         0.001016777 0.5555556 16.40807 10

```

