

Intelligence-Artificielle



- Dispensé par **MWAMBA KASONGO Dahouda**
- Docteur en génie logiciel et systèmes d'information
- Machine and Deep Learning Engineer

- Assisté par Ass. **Jason MUSA**

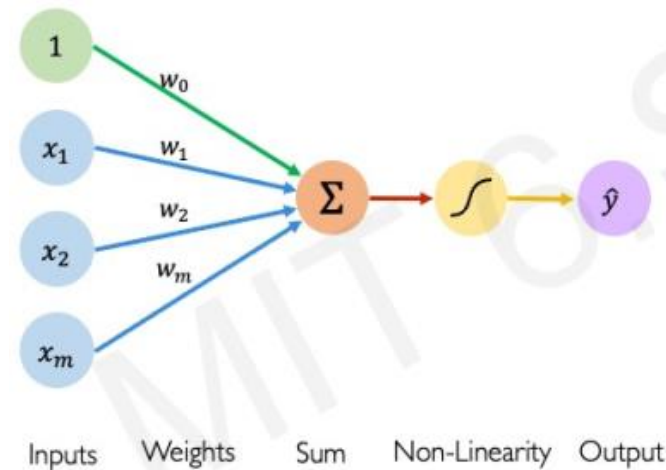
Mardi : 8H00 – 12H00

Mercredi : 13H00 – 17H00

Vendredi : 13H00 – 13H00

- E-mail : dahouda37@gmail.com
- Tel.: +243 99 66 55 265

The Perceptron: Forward Propagation

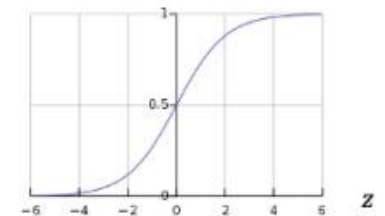


Activation Functions

$$\hat{y} = g(w_0 + \mathbf{X}^T \mathbf{W})$$

- Example: sigmoid function

$$g(z) = \sigma(z) = \frac{1}{1 + e^{-z}}$$



PLAN DU COURS

CHAPITRE 2 Concepts clés sur l'apprentissage automatique

2.1. Terminologies d'apprentissage automatique

- **Features (Caractéristiques), Label (étiquettes) et Dataset (ensembles de données)**

2.2 Types de données

- **Catégorielles, Numériques, Textuelles, Images, audio**
- Training set (Données d'entraînement), Validation set (Données de validation), Test set (Données de test)

2.3 Introduction aux algorithmes (Type de Machine Learning)

- **Qu'est-ce qu'un algorithme ?**
- **Les algorithmes de Machine Learning**

2.4 Exemple Pratique de la Régression linéaire simple

2.5 Exemple Pratique de la Régression linéaire multiple

2.6. Introduction à l'apprentissage automatique sur AWS



CHAPITRE 2 MACHINE LEARNING

2.1. Terminologies d'apprentissage automatique

✓ **Ensemble de données [Dataset]** : est un ensemble de données organisées de manière structurée ou non structurée.

Les ensembles de données peuvent se présenter sous différents formats et peuvent contenir différents types de données, tels que du texte, des nombres, des images, des vidéos ou de l'audio. Exemple d'une Dataset: <https://archive.ics.uci.edu/ml/datasets/Bank+Marketing>



| age | job | marital | education | default | balance | housing | loan | contact | day | month | duration | campaign | pdays | previous | poutcome | y |
|-----|-------------|---------|-----------|---------|---------|---------|------|----------|-----|-------|----------|----------|-------|----------|----------|-----|
| 30 | unemployed | married | primary | no | 1787 | no | no | cellular | 19 | oct | 79 | 1 | -1 | 0 | unknown | no |
| 33 | services | married | secondary | no | 4789 | yes | yes | cellular | 11 | may | 220 | 1 | 339 | 4 | failure | no |
| 35 | managemen | single | tertiary | no | 1350 | yes | no | cellular | 16 | apr | 185 | 1 | 330 | 1 | failure | no |
| 30 | managemen | married | tertiary | no | 1476 | yes | yes | unknown | 3 | jun | 199 | 4 | -1 | 0 | unknown | no |
| 59 | blue-collar | married | secondary | no | 0 | yes | no | unknown | 5 | may | 226 | 1 | -1 | 0 | unknown | no |
| 35 | managemen | single | tertiary | no | 747 | no | no | cellular | 23 | feb | 141 | 2 | 176 | 3 | failure | no |
| 36 | self-employ | married | tertiary | no | 307 | yes | no | cellular | 14 | may | 341 | 1 | 330 | 2 | other | no |
| 39 | technician | married | secondary | no | 147 | yes | no | cellular | 6 | may | 151 | 2 | -1 | 0 | unknown | no |
| 41 | entrepreneu | married | tertiary | no | 221 | yes | no | unknown | 14 | may | 57 | 2 | -1 | 0 | unknown | no |
| 43 | services | married | primary | no | -88 | yes | yes | cellular | 17 | apr | 313 | 1 | 147 | 2 | failure | no |
| 39 | services | married | secondary | no | 9374 | yes | no | unknown | 20 | may | 273 | 1 | -1 | 0 | unknown | no |
| 43 | admin. | married | secondary | no | 264 | yes | no | cellular | 17 | apr | 113 | 2 | -1 | 0 | unknown | no |
| 36 | technician | married | tertiary | no | 1109 | no | no | cellular | 13 | aug | 328 | 2 | -1 | 0 | unknown | no |
| 20 | student | single | secondary | no | 502 | no | no | cellular | 30 | apr | 261 | 1 | -1 | 0 | unknown | yes |
| 31 | blue-collar | married | secondary | no | 360 | yes | yes | cellular | 29 | jan | 89 | 1 | 241 | 1 | failure | no |

Categorical Variables

Numerical Variables

CHAPITRE 2 CONCEPTS CLES SUR MACHINE LEARNING



2.1. Terminologies d'apprentissage automatique

- ✓ **Modèle [Model]** : un modèle est une représentation spécifique apprise à partir de données en appliquant un algorithme de ML.
Un modèle est également appelé hypothèse.
- ✓ **Fonctionnalité [Feature]** : une caractéristique est une propriété individuelle mesurable de nos données.
Un ensemble de caractéristiques numériques peut être décrit de manière pratique par un vecteur de caractéristiques.

- ❖ Par exemple, pour prédire un fruit, il peut y avoir des caractéristiques comme la couleur, l'odeur, le goût, etc.
Remarque : le choix de caractéristiques informatives, discriminantes et indépendantes est une étape cruciale pour des algorithmes efficaces.

- ✓ **Cible ou étiquette [Target ou Label]** : une variable cible ou une étiquette est la valeur à prédire par notre modèle.
Pour l'exemple de fruit abordé dans la section sur les fonctionnalités, l'étiquette de chaque ensemble d'entrées serait le nom du fruit comme la pomme, l'orange, la banane, etc.
- ✓ **Entraînement [Training]** : l'idée est de fournir un ensemble d'entrées (caractéristiques) et ses sorties attendues (étiquettes), de sorte qu'après l'entraînement, nous aurons un modèle (hypothèse) qui mapperait ensuite les nouvelles données à l'une des catégories sur lesquelles l'entraînement a été effectué.
- ✓ **Prédiction** : une fois que notre modèle est prêt, il peut être alimenté par un ensemble d'entrées auxquelles il fournira une sortie prédite (étiquette). Mais assurez-vous que si la machine fonctionne bien sur des données invisibles (Test data), alors seulement nous pouvons dire que le modèle fonctionne bien.

CHAPITRE 2 CONCEPTS CLES SUR MACHINE LEARNING



2.2. Types de données

1. **Ensemble de données tabulaires** : organisé en lignes et en colonnes.

Couramment utilisé dans les affaires, les soins de santé et la finance (par exemple, les feuilles de calcul).

Out[9]:

| | price | resid_area | air_qual | room_num | age | dist1 | dist2 | dist3 | dist4 | teachers | poor_prop | n_hos_beds | n_hot_rooms | rainfall | parks | Sold |
|---|-------|------------|----------|----------|------|-------|-------|-------|-------|----------|-----------|------------|-------------|----------|----------|------|
| 0 | 24.0 | 32.31 | 0.538 | 6.575 | 65.2 | 4.35 | 3.81 | 4.18 | 4.01 | 24.7 | 4.98 | 5.480 | 11.1920 | 23 | 0.049347 | 0 |
| 1 | 21.6 | 37.07 | 0.469 | 6.421 | 78.9 | 4.99 | 4.70 | 5.12 | 5.06 | 22.2 | 9.14 | 7.332 | 12.1728 | 42 | 0.046146 | 1 |
| 2 | 34.7 | 37.07 | 0.469 | 7.185 | 61.1 | 5.03 | 4.86 | 5.01 | 4.97 | 22.2 | 4.03 | 7.394 | 101.1200 | 38 | 0.045764 | 0 |
| 3 | 33.4 | 32.18 | 0.458 | 6.998 | 45.8 | 6.21 | 5.93 | 6.16 | 5.96 | 21.3 | 2.94 | 9.268 | 11.2672 | 45 | 0.047151 | 0 |
| 4 | 36.2 | 32.18 | 0.458 | 7.147 | 54.2 | 6.16 | 5.86 | 6.37 | 5.86 | 21.3 | 5.33 | 8.824 | 11.2896 | 55 | 0.039474 | 0 |

In [10]: `df.shape`

Out[10]: (506, 16)

Variables numériques

CHAPITRE 2 CONCEPTS CLES SUR MACHINE LEARNING



2.2. Types de données

1. **Ensemble de données tabulaires** : organisé en lignes et en colonnes.

Couramment utilisé dans les affaires, les soins de santé et la finance (par exemple, les feuilles de calcul).

| | id | bin_0 | bin_1 | bin_2 | bin_3 | bin_4 | nom_0 | nom_1 | nom_2 | nom_3 | ... | nom_8 | nom_9 | ord_0 | ord_1 | ord_2 | ord_3 | ord_4 | ord_5 | day | month |
|---|--------|-------|-------|-------|-------|-------|-------|-----------|---------|---------|-----|-----------|-----------|-------|-------------|----------|-------|-------|-------|-----|-------|
| 0 | 300000 | 0 | 0 | 1 | T | Y | Blue | Triangle | Axolotl | Finland | ... | 9d117320c | 3c49b42b8 | 2 | Novice | Warm | j | P | be | 5 | 11 |
| 1 | 300001 | 0 | 0 | 0 | T | N | Red | Square | Lion | Canada | ... | 46ae3059c | 285771075 | 1 | Master | Lava Hot | l | A | RP | 7 | 5 |
| 2 | 300002 | 1 | 0 | 1 | F | Y | Blue | Square | Dog | China | ... | b759e21f0 | 6f323c53f | 2 | Expert | Freezing | a | G | tP | 1 | 12 |
| 3 | 300003 | 0 | 0 | 1 | T | Y | Red | Star | Cat | China | ... | 0b6ec68ff | b5de3dcc4 | 1 | Contributor | Lava Hot | b | Q | ke | 2 | 3 |
| 4 | 300004 | 0 | 1 | 1 | F | N | Red | Trapezoid | Dog | China | ... | f91f3b1ee | 967cfa9c9 | 3 | Grandmaster | Lava Hot | l | W | qK | 4 | 11 |

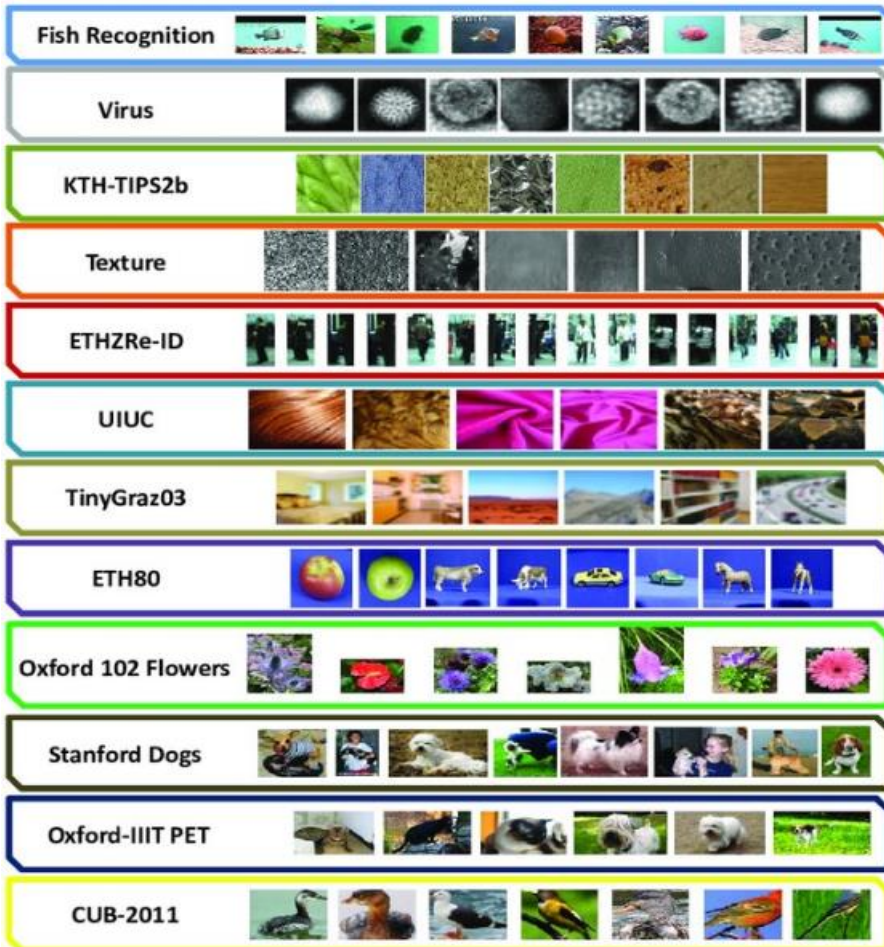
Variables numériques et catégorielles

CHAPITRE 2 CONCEPTS CLES SUR MACHINE LEARNING



2.2. Types de données

2. Ensemble de données d'images : Contient des images, souvent utilisées dans les tâches de vision par ordinateur (MNIST, CIFAR-10).



airplane

automobile

bird

cat

deer

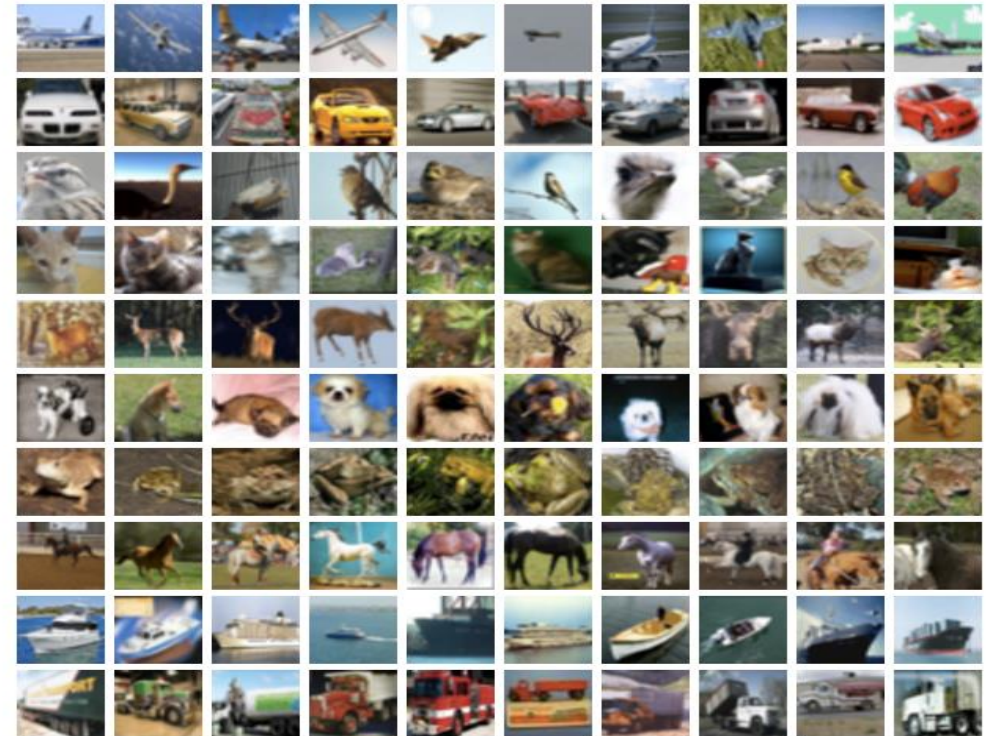
dog

frog

horse

ship

truck



CHAPITRE 2 CONCEPTS CLES SUR MACHINE LEARNING



2.2. Types de données

3. Ensemble de données textuelles : contient des données textuelles pour le traitement du langage naturel (par exemple, l'analyse des sentiments, la traduction linguistique).

| | review | sentiment |
|---|---------------------------------------------------|-----------|
| 0 | One of the other reviewers has mentioned that ... | positive |
| 1 | A wonderful little production. The... | positive |
| 2 | I thought this was a wonderful way to spend ti... | positive |
| 3 | Basically there's a family where a little boy ... | negative |
| 4 | Petter Mattei's "Love in the Time of Money" is... | positive |
| 5 | Probably my all-time favorite movie, a story o... | positive |
| 6 | I sure would like to see a resurrection of a u... | positive |
| 7 | This show was an amazing, fresh & innovative i... | negative |
| 8 | Encouraged by the positive comments about this... | negative |
| 9 | If you like original gut wrenching laughter yo... | positive |

L'ensemble de données IMDB contient 50 000 critiques de films étiquetées comme des sentiments « positifs » ou « négatifs ».

L'analyse des sentiments est une tâche cruciale de traitement du langage naturel (NLP) qui consiste à déterminer le sentiment ou l'émotion exprimé dans un texte.

CHAPITRE 2 CONCEPTS CLES SUR MACHINE LEARNING



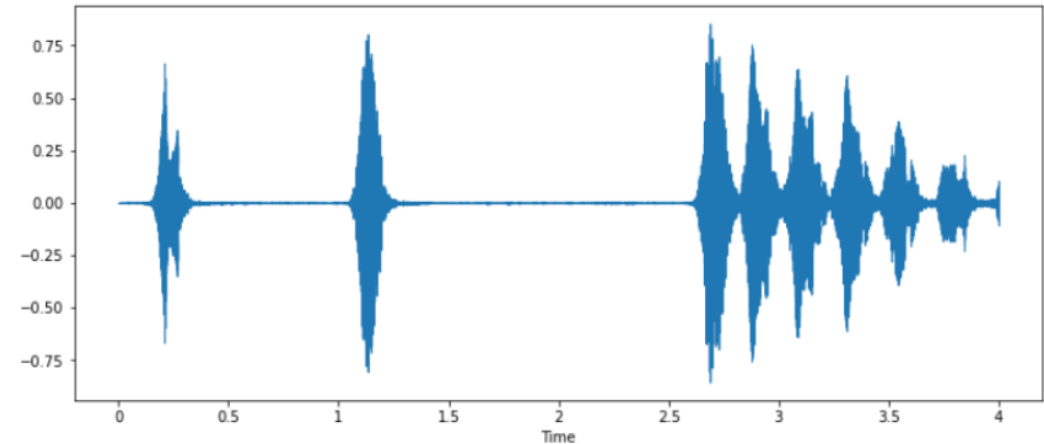
2.2. Types de données

4. Ensemble de données audio : contient des données sonores ou vocales (par exemple, utilisées pour des tâches de reconnaissance vocale).



Military Audio Dataset

Un ensemble de données audio militaires pour la connaissance de la situation et la surveillance



Urban Sound 8k Dataset

L'ensemble de données Urban Sound 8k. L'ensemble de données contient 8732 fichiers sonores de 10 classes différentes et est répertorié ci-dessous : 1. Air Conditioner, 2. Car Horn, 3. Children Playing, 4. Dog Bark, 5. Drilling Machine, 6. Engine Idling, 7. Gun Shot, 8. Jackhammer, 9. Siren, 10. Street Music

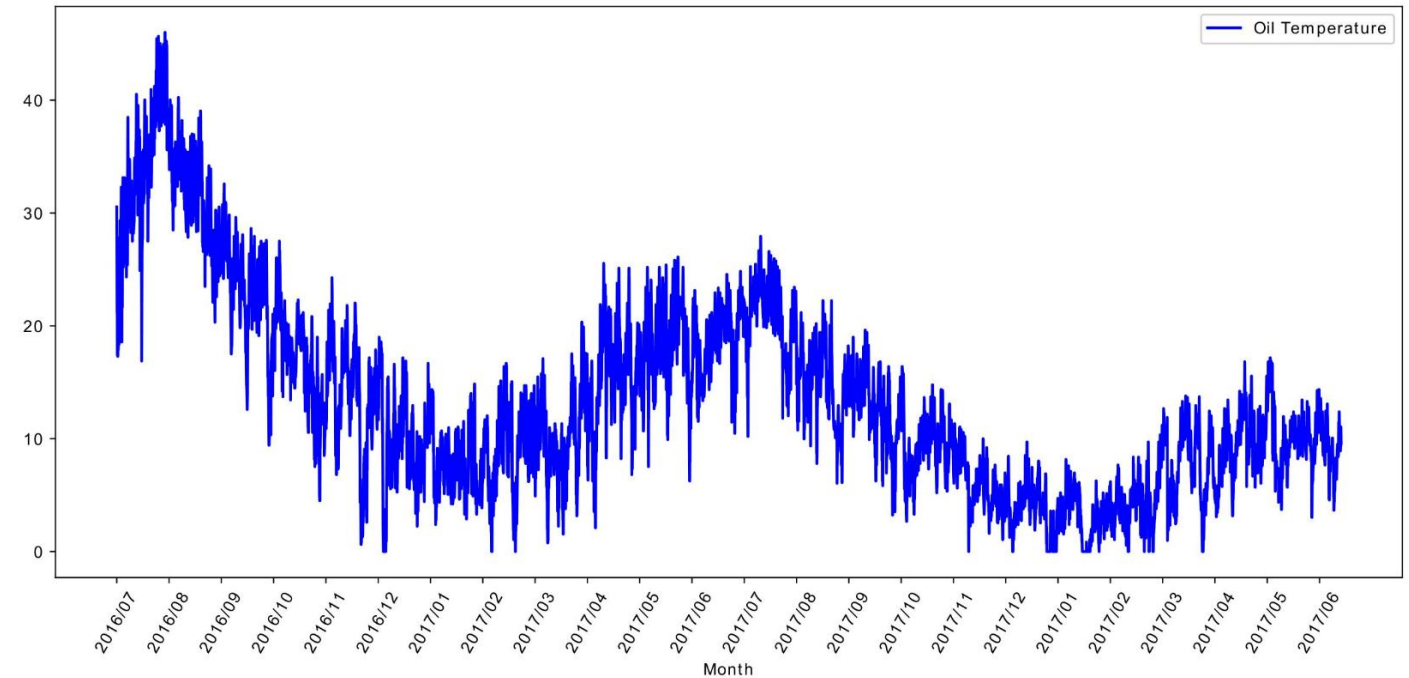
CHAPITRE 2 CONCEPTS CLES SUR MACHINE LEARNING



2.2. Types de données

5. Ensemble de données de séries chronologiques : Points de données indexés dans le temps, utilisés dans les prévisions ou les analyses de tendances (par exemple, cours des actions, données météorologiques).

| A | B | C | D | E | F |
|------------|--------------|----------|----------|--------------|-------|
| date | product_type | location | discount | weather_temp | sales |
| 2011-01-01 | A | X | 0.2 | 25 | 50000 |
| 2011-01-01 | A | Y | 0.15 | 27 | 6000 |
| 2011-01-01 | A | Z | 0.1 | 26 | 70000 |
| 2011-01-01 | B | X | 0.3 | 25 | 60000 |
| 2011-01-01 | B | Y | 0.25 | 27 | 8000 |
| 2011-01-01 | B | Z | 0.2 | 26 | 9000 |
| 2011-01-01 | C | X | 0.13 | 25 | 10000 |
| 2011-01-01 | C | Y | 0.14 | 27 | 65000 |
| 2011-01-01 | C | Z | 0.16 | 26 | 30000 |
| 2011-01-02 | A | X | 0.2 | 25 | 50000 |
| 2011-01-02 | A | Y | 0.15 | 27 | 6000 |
| 2011-01-02 | A | Z | 0.1 | 26 | 70000 |
| 2011-01-02 | B | X | 0.3 | 25 | 60000 |
| 2011-01-02 | B | Y | 0.25 | 27 | 8000 |
| 2011-01-02 | B | Z | 0.2 | 26 | 9000 |
| 2011-01-02 | C | X | 0.13 | 25 | 10000 |
| 2011-01-02 | C | Y | 0.14 | 27 | 65000 |
| 2011-01-02 | C | Z | 0.16 | 26 | 30000 |



ETT (Electricity Transformer Temperature) Dataset

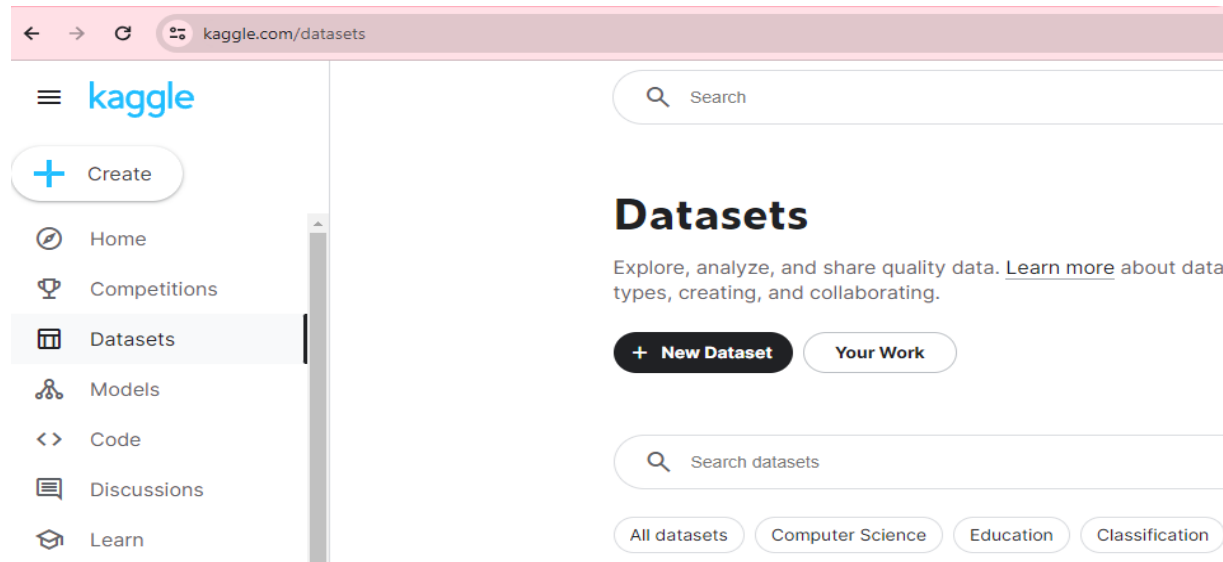
CHAPITRE 2 CONCEPTS CLES SUR MACHINE LEARNING



2.3. Source des données Publiques

Il existe de nombreuses sources où vous pouvez obtenir des données pour l'analyse de la science des données, en fonction de vos intérêts spécifiques et du type d'analyse que vous souhaitez effectuer. Voici quelques options populaires :

1. Kaggle : Kaggle est une plate-forme de concours de science des données, mais elle héberge également un grand nombre d'ensembles de données disponibles gratuitement pour l'exploration et l'analyse. Lien du Kaggle : <https://www.kaggle.com/datasets>



CHAPITRE 2 CONCEPTS CLES SUR MACHINE LEARNING



2.3. Source des données Publiques

2. Référentiel UCI Machine Learning : le référentiel UCI Machine Learning est un ensemble de bases de données, de théories de domaine et de générateurs de données largement utilisés par la communauté de Machine Learning. Lien de UCI : <https://archive.ics.uci.edu/>


[Datasets](#) [Contribute Dataset](#) [About Us](#)

Welcome to the UC Irvine Machine Learning Repository

We currently maintain 664 datasets as a service to the machine learning community. Here, you can donate and find datasets used by millions of people all around the world!

[VIEW DATASETS](#) [CONTRIBUTE A DATASET](#)


Popular Datasets



Iris

A small classic dataset from Fisher, 1936. One of the earliest known datasets used for ev...

[Classification](#) [150 Instances](#) [4 Features](#)




Dry Bean

Images of 13,611 grains of 7 different registered dry beans were taken with a high-resol...

[Classification](#) [13.61K Instances](#) [16 Features](#)


New Datasets



PhiUSIIL Phishing URL (Website)

PhiUSIIL Phishing URL Dataset is a substantial dataset comprising 134,850 legitimate an...

[Classification](#) [235.8K Instances](#) [54 Features](#)



RT-IoT2022

The RT-IoT2022, a proprietary dataset derived from a real-time IoT infrastructure, is intro...

[Classification, Regressi...](#) [123.12K Instances](#) [84 Features](#)

CHAPITRE 2 CONCEPTS CLES SUR MACHINE LEARNING



2.3. Source des données Publiques

3. Recherche d'ensembles de données Google : Google Dataset Search vous aide à trouver des ensembles de données stockés sur le Web.

C'est un outil utile pour découvrir des ensembles de données provenant de diverses sources : <https://datasetsearch.research.google.com/>

Dataset Search

The screenshot shows the Google Dataset Search interface. At the top, there is a search bar with the text 'covid' entered. Below the search bar, a list of search results is displayed. The first result is 'coronavirus covid-19'. The second result, 'covid', is highlighted. Other results include 'covid 19', 'covid-19', 'Vaccines.gov: COVID-19 vaccinating provider locations', 'COVID-19 Hospital Data Coverage Report', 'COVID-19-Open-Research-Dataset-Challenge--CORD-19-', 'Assessment of Effectiveness of COVID 19 Pandemic Scheduling Triage in an Academic Dermatology Clinic', and 'Number of COVID-19 people killed by age'.

CHAPITRE 2 CONCEPTS CLES SUR MACHINE LEARNING



2.3. Source des données Publiques

4. Portails de données ouvertes gouvernementaux : de nombreux gouvernements proposent des portails de données ouvertes où vous pouvez trouver des ensembles de données liés à la démographie, à l'économie, à la santé, etc. Lien : <https://data.gov/>



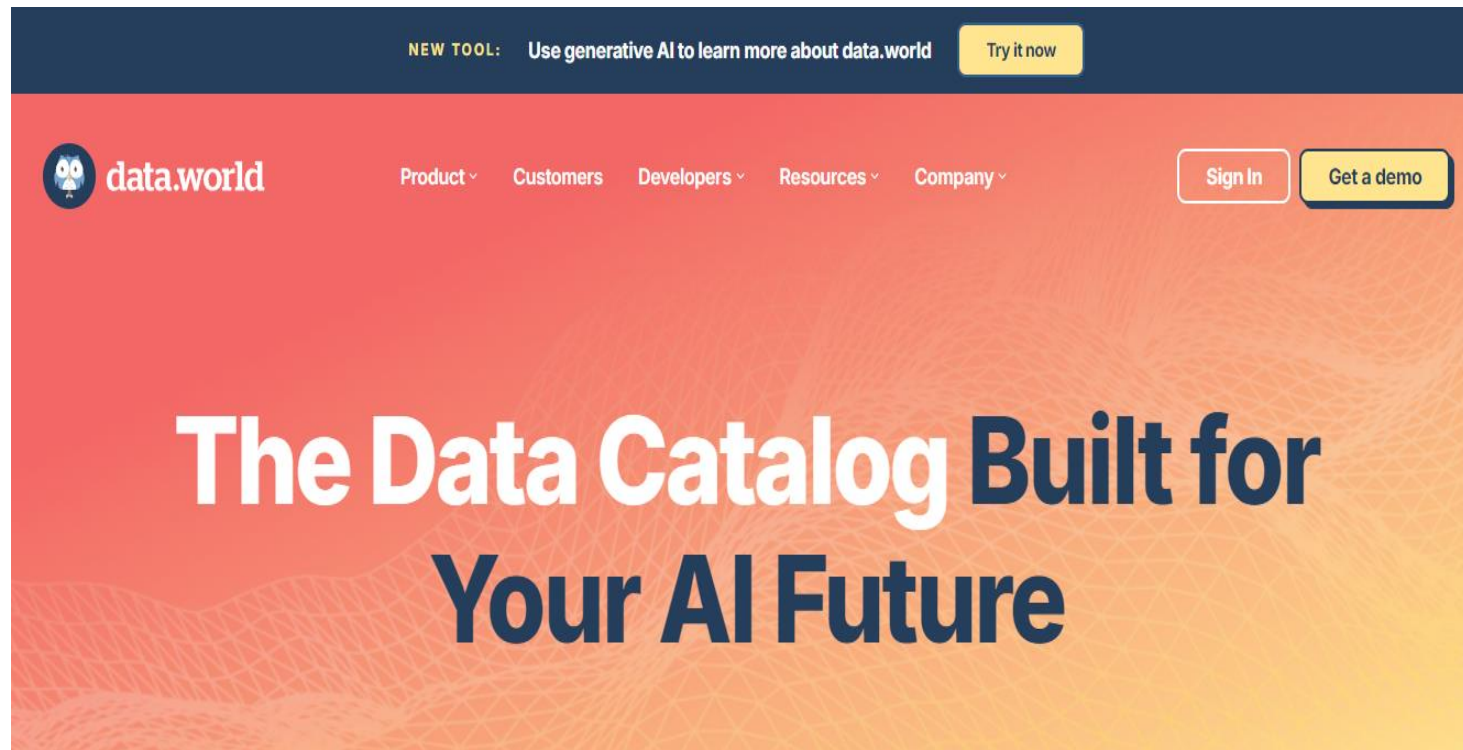
CHAPITRE 2 CONCEPTS CLES SUR MACHINE LEARNING



2.3. Source des données Publiques

5. **Data.world** : Data.world est une plateforme où vous pouvez trouver et partager des ensembles de données.

Il héberge une gamme diversifiée d'ensembles de données fournis par la communauté. Lien : <https://data.world/>



CHAPITRE 2 CONCEPTS CLES SUR MACHINE LEARNING



2.3. Source des données Publiques

6. Ensembles de données Reddit : le sous-reddit r/datasets est une communauté où les gens partagent et demandent des ensembles de données. Vous pourriez trouver des ensembles de données intéressants ici <https://www.reddit.com/r/datasets/>

The screenshot displays the Reddit interface for the r/datasets subreddit. The left sidebar shows navigation options like Home, Popular, and a list of recent topics including Gaming, Sports, Business, Crypto, and Television. The main content area features a post from user 'u/Swat_Sam2' asking for help with a data analysis project involving a MySQL server. The post text is as follows:

Help with data analysis project (mysql online server help)

I have to create a power BI project with a data which should be present in MySQL online hosted server But the problem is that the data which i have is 2 tables with 130k rows each (csv files), and i made a mysql server on freemysqlhosting.net but there are 2 problems, firstly it has a 5mb limit for the database Secondly each row takes about 4 seconds to upload And on this speed i think itll take 6 days to just upload 1 table Is there any other way to do this? Maybe something like, i could make the database in the local mysql server with the tables which doesn't take much time and then i could maybe set up this server to be accessible to public somehow Please help

The right sidebar contains community statistics: 190K Members, 35 Online, and Top 1% Rank by size. It also includes a 'COMMUNITY BOOKMARKS' section with a 'mod' button and a 'RULES' section.

CHAPITRE 2 CONCEPTS CLES SUR MACHINE LEARNING



2.3 Introduction aux algorithmes

- ❑ Les algorithmes d'apprentissage automatique sont des techniques utilisés pour créer des systèmes capables d'apprendre à partir de données et de faire des prédictions ou de prendre des décisions sans être explicitement programmés.
- ❑ Ces algorithmes se répartissent en différentes catégories en fonction du type d'apprentissage qu'ils prennent en charge : apprentissage supervisé, non supervisé, semi-supervisé ou par renforcement.

2.3.1. Algorithmes d'apprentissage supervisé

Dans l'apprentissage supervisé, l'algorithme est formé sur des données étiquetées (où l'entrée et la sortie correspondante sont connues). L'objectif est d'apprendre une correspondance entre les entrées et les sorties.

1. Régression linéaire [Linear Regression]

- **Objectif** : Prédire des valeurs continues (par exemple, les prix des maisons).
- **Description** : Modélise la relation entre les caractéristiques d'entrée (variables indépendantes) et une sortie continue (variable dépendante) en ajustant une ligne droite (ou hyperplan) aux données.

2. Régression logistique [Logistic Regression]

- **Objectif** : Problèmes de classification binaire (par exemple, détection de spam)
- **Description** : Similaire à la régression linéaire, mais utilisée pour prédire les résultats catégoriels.
Elle génère des probabilités à l'aide d'une fonction logistique.

CHAPITRE 2 CONCEPTS CLES SUR MACHINE LEARNING



2.3 Introduction aux algorithmes

2.3.1. Algorithmes d'apprentissage supervisé

3. Arbres de décision [Decision Trees]

- **Objectif** : Classification et régression.
- **Description** : Modèle de décisions de type arborescence, où les nœuds internes représentent des tests sur des caractéristiques, les branches représentent les résultats de ces tests et les nœuds feuilles représentent le résultat prédit.

4. Random Forest

- **Objectif** : Classification et régression.
- **Description** : Ensemble d'arbres de décision dans lesquels plusieurs arbres sont construits sur des sous-ensembles aléatoires de données et de caractéristiques. La prédiction finale est basée sur le vote majoritaire (classification) ou la moyenne (régression) de tous les arbres.

5. Support Vector Machines (SVM)

- **Objectif** : Classification et régression.
- **Description** : Recherche l'hyperplan optimal qui sépare au maximum les données en différentes classes.
Il fonctionne bien pour les ensembles de données de grande dimension.

CHAPITRE 2 CONCEPTS CLES SUR MACHINE LEARNING

2.3 Introduction aux algorithmes

2.3.1. Algorithmes d'apprentissage supervisé

6. K-Nearest Neighbors (k-NN)

- **Objectif** : Classification et régression.
- **Description** : Algorithme non paramétrique dans lequel la prédiction est faite sur la base de la classe majoritaire ou de la moyenne des « k » points les plus proches dans l'espace des caractéristiques.

7. Naive Bayes

- **Objectif** : Classification (par exemple, classification de texte, filtrage du spam)
- **Description** : Un classificateur probabiliste basé sur le théorème de Bayes, supposant que les caractéristiques sont indépendantes les unes des autres (ce qui est une hypothèse « naïve », d'où son nom).



CHAPITRE 2 CONCEPTS CLES SUR MACHINE LEARNING



2.3 Introduction aux algorithmes

2.3.2. Algorithmes d'apprentissage non supervisé

Dans l'apprentissage non supervisé, l'algorithme travaille avec des données non étiquetées et essaie de trouver des modèles ou des structures cachés en leur sein.

1. Clustering k-Means [K-Means Clustering]

- **Objectif** : Clustering (regroupement de points de données similaires).
- **Description** : Partitionne les données en « k » clusters où chaque point de données appartient au cluster avec le centroïde le plus proche (moyenne).

2. Clustering hiérarchique [Hierarchical Clustering]

- **Objectif** : Clustering.
- **Description** : Construit une hiérarchie de clusters en fusionnant ou en divisant de manière répétée des clusters.
Le résultat est un arbre de clusters.

3. Analyse en composantes principales [Principal Component Analysis (PCA)]

- **Objectif** : Réduction de la dimensionnalité.
- **Description** : Transforme les données de grande dimension en un espace de dimension inférieure tout en conservant autant de variance que possible. Souvent utilisé pour la visualisation des données ou la réduction du bruit.

CHAPITRE 2 CONCEPTS CLES SUR MACHINE LEARNING



2.3 Introduction aux algorithmes

2.3.3. Algorithmes d'apprentissage semi-supervisé

Dans l'apprentissage semi-supervisé, l'algorithme utilise une petite quantité de données étiquetées et une grande quantité de données non étiquetées.

1. Algorithmes d'auto-apprentissage [Self-training Algorithms]

Dans ces méthodes, le modèle est d'abord formé sur les données étiquetées, puis étiquette de manière itérative les données non étiquetées pour améliorer les performances du modèle.

2.3.4. Algorithmes d'apprentissage par renforcement

L'apprentissage par renforcement consiste à apprendre en interagissant avec un environnement, où l'agent apprend à prendre des mesures pour maximiser les récompenses cumulatives.

1. Q-Learning

- **Objectif** : Apprendre des politiques optimales pour les problèmes de prise de décision séquentielle.
- **Description** : Un algorithme basé sur la valeur où l'agent apprend la valeur de chaque action dans un état donné en maximisant les récompenses futures.

CHAPITRE 2 CONCEPTS CLES SUR MACHINE LEARNING

2.3 Introduction aux algorithmes



2.3.4. Algorithmes d'apprentissage par renforcement

2. Deep Q-Networks (DQN)

- **Objectif** : Version avancée de Q-Learning qui fonctionne bien pour les environnements complexes.
- **Description** : Combine le Q-learning avec l'apprentissage profond[Deep Learning] en utilisant des réseaux neuronaux pour estimer les valeurs Q, souvent utilisées dans les jeux vidéo ou la robotique.

3. Méthodes de gradient de politique [Policy Gradient Methods]

- **Objectif** : Apprendre directement la politique au lieu de la fonction de valeur.
- **Description** : L'agent apprend la distribution de probabilité des actions plutôt que d'estimer la valeur des actions. Couramment utilisé dans les espaces d'action continue.

CHAPITRE 2 CONCEPTS CLES SUR MACHINE LEARNING



2.3 Introduction aux algorithmes

2.3.5. Algorithmes d'apprentissage d'ensemble [Ensemble Learning Algorithms]

L'apprentissage d'ensemble consiste à combiner plusieurs modèles d'apprentissage automatique pour améliorer les performances globales.

1. Bagging

- **Objectif** : Réduire la variance et éviter le surajustement (Overfitting).
- **Description** : Plusieurs modèles (par exemple, des arbres de décision) sont formés sur différents sous-ensembles de données, et leurs prédictions sont combinées (par exemple, par vote majoritaire ou par moyenne) pour améliorer les performances.

2. Boosting

- **Objectif** : Réduire les biais et améliorer la précision.
- **Description** : Les apprenants faibles (modèles qui fonctionnent légèrement mieux que les devinettes aléatoires) sont formés séquentiellement, chaque nouveau modèle se concentrant sur la correction des erreurs des modèles précédents.
- **Exemples** : **AdaBoost**, **Gradient Boosting Machines** (GBM), **XGBoost**.

3. Stacking

- **Objectif** : Combiner les points forts de différents modèles.
- **Description** : Plusieurs modèles sont entraînés, et leurs prédictions sont utilisées comme entrées dans un méta-modèle de niveau supérieur, qui apprend à combiner les prédictions.

CHAPITRE 2 CONCEPTS CLES SUR MACHINE LEARNING

2.3 Introduction aux algorithmes

2.3.6. Résumé des algorithmes par catégorie



| Catégorie | Algorithmes |
|--------------------------------|-----------------------------------------------------------------------------------------------|
| Apprentissage supervisé | Linear Regression, Logistic Regression, SVM, k-NN, Random Forest, Decision Trees, Naïve Bayes |
| Apprentissage non supervisé | k-Means, Hierarchical Clustering, PCA |
| Apprentissage par renforcement | Q-Learning, DQN, Policy Gradients |
| Apprentissage d'ensemble | Bagging, AdaBoost, XGBoost, Stacking |

Regression Linear simple [Lab]

CHAPITRE 2 CONCEPTS CLES SUR MACHINE LEARNING



2.4 Exemple Pratique de la Prédiction

2.4.1. Régression linéaire simple [Simple Linear Regression]

- ❑ La prédiction d'une réponse à l'aide d'une ou de plusieurs caractéristiques est une méthode de prédiction de la variable dépendante (Y) en fonction des valeurs des variables indépendantes (X).
- ❑ On suppose que les deux variables sont linéairement liées. Par conséquent, nous essayons de trouver une fonction linéaire qui prédit la réponse en fonction de la caractéristique ou de la variable indépendante (x).

Les équations mathématiques qui décrivent une régression linéaire simple et régression linéaire multiple sont présentées dans les équations suivantes:

Simple
Linear
Regression

$$y = b_0 + b_1 * x_1$$

Multiple
Linear
Regression

$$y = b_0 + b_1 * x_1 + b_2 * x_2 + \dots + b_n * x_n$$

$$\begin{array}{c} \text{Constant/Intercept} \\ \downarrow \\ Y_i = \beta_0 + \beta_1 X_i \\ \uparrow \qquad \qquad \uparrow \\ \text{Dependent Variable} \quad \text{Slope/Coefficient} \end{array}$$

Independent Variable
↓

CHAPITRE 2 CONCEPTS CLES SUR MACHINE LEARNING



2.4 Exemple Pratique de la Prédiction

2.4.1. Régression linéaire multiple [Multiple Linear Regression]

□ L'équation mathématique qui décrit une régression linéaire multiple est présentée dans l'équation suivante:

$$y = \alpha + \beta_1(x_1) + \beta_2(x_2) + \dots + \beta_n(x_n)$$

Diagram illustrating the components of the Multiple Linear Regression equation:

- y : Predicted value
- α : Bias
- β_1 : Weight 1
- x_1 : Feature 1
- β_2 : Weight 2
- x_2 : Feature 2
- β_n : Weight n
- x_n : Feature n

Dans cette équation :

y est la valeur prédite ou variable dépendante

x est les caractéristiques ou variable indépendante

α est le biais

β est le poids de chaque caractéristique

CHAPITRE 2 CONCEPTS CLES SUR MACHINE LEARNING

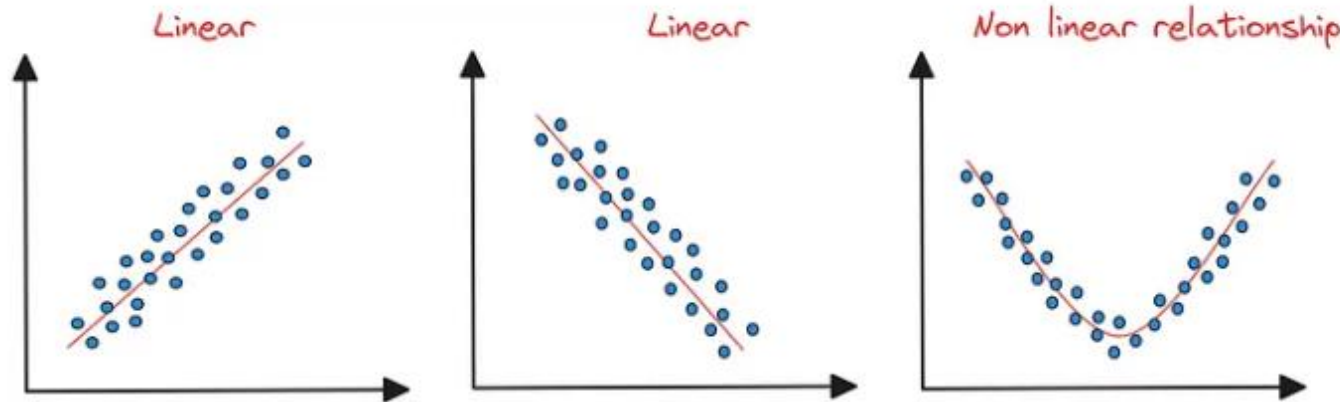


2.4 Exemple Pratique de la Prédiction

2.4.1. Régression linéaire multiple [Multiple Linear Regression]

❑ Remarque importante !

- ❑ Il existe plusieurs hypothèses importantes pour effectuer une analyse de régression.
- ❑ Certaines des hypothèses que vous devez confirmer sont les suivantes :
- ✓ **Linéarité** : la relation entre la variable dépendante et la variable indépendante doit être linéaire. Autrement dit, chaque changement d'unité dans la valeur de la variable indépendante entraîne le même changement dans la variable dépendante.



Relation entre les variables : linéaire (corrélation positive), linéaire (corrélation négative), relation non linéaire

CHAPITRE 2 CONCEPTS CLES SUR MACHINE LEARNING



2.4 Exemple Pratique de la Prédiction

2.4.1. Régression linéaire simple [Simple Linear Regression]

❑ L'équation mathématique qui décrit une régression linéaire simple est présentée dans l'équation (1).

Variable Dépendante

$$y = b_0 + b_1 x_1 \quad (1)$$

variables indépendantes

Exemple 1 : Prédiction du pourcentage de notes qu'un étudiant devrait obtenir en fonction du nombre d'heures qu'il a étudiées

Formulation du probleme :

Bias

Feature [variables indépendante]

$$Score = b_0 + b_1 hours$$

Valeur predate [variables dépendante]

Weight [Coefficient]

Dataset

| | Hours | Scores |
|---|-------|--------|
| 0 | 2.5 | 21 |
| 1 | 5.1 | 47 |
| 2 | 3.2 | 27 |
| 3 | 8.5 | 75 |
| 4 | 3.5 | 30 |

CHAPITRE 2 CONCEPTS CLES SUR MACHINE LEARNING



2.4 Exemple Pratique de la Prédiction

2.4.1. Régression linéaire simple [Simple Linear Regression]

Exemple 1 : Prédiction du pourcentage de notes qu'un étudiant devrait obtenir en fonction du nombre d'heures qu'il a étudiées

Nous allons suivre les étapes de prétraitement des données, puis construire le modèle de régression linéaire simple.
Les étapes sont les suivantes :

1. Importer les bibliothèques

```
1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
```

- ❖ Nous importons trois bibliothèques essentielles en Python qui sont couramment utilisées pour l'analyse, la manipulation et la visualisation des données.
- ❖ Nous connaissons déjà numpy et pandas. À la ligne 3, nous importons la bibliothèque matplotlib, qui est une bibliothèque de traçage pour le langage de programmation Python.

CHAPITRE 2 CONCEPTS CLES SUR MACHINE LEARNING



2.4 Exemple Pratique de la Prédiction

2.4.1. Régression linéaire simple [Simple Linear Regression]

Exemple 1 : Prédiction du pourcentage de notes qu'un étudiant devrait obtenir en fonction du nombre d'heures qu'il a étudiées

```
1 dataset = pd.read_csv("../Data/studentscores.csv")
2 dataset.head(5)
```

| | Hours | Scores |
|---|-------|--------|
| 0 | 2.5 | 21 |
| 1 | 5.1 | 47 |
| 2 | 3.2 | 27 |
| 3 | 8.5 | 75 |
| 4 | 3.5 | 30 |

- **pd.read_csv** : cette fonction lit un fichier CSV (comma-separated values) dans un DataFrame.
- **"../Data/studentscores.csv"** : le chemin d'accès au fichier CSV. Le chemin relatif ../Data/ signifie que le fichier se trouve dans le répertoire Data, un niveau au-dessus du répertoire de travail actuel.
- **dataset.head(5)** : cette méthode renvoie les cinq premières lignes du DataFrame. Si vous omettez l'argument, les cinq premières lignes sont affichées par défaut.

CHAPITRE 2 CONCEPTS CLES SUR MACHINE LEARNING



2.4 Exemple Pratique de la Prédiction

2.4.1. Régression linéaire simple [Simple Linear Regression]

Exemple 1 : Prédiction du pourcentage de notes qu'un étudiant devrait obtenir en fonction du nombre d'heures qu'il a étudiées

3. Vérifier les données manquantes

```
[10]: 1 dataset.isnull().sum()
```

```
[10]: Hours      0  
      Scores    0  
      dtype: int64
```

- La méthode `dataset.isnull().sum()` permet d'identifier et de compter le nombre de valeurs manquantes dans chaque colonne de notre DataFrame.
- `dataset.isnull()` : Cette méthode renvoie un DataFrame de la même forme qu'une dataset, mais des valeurs booléennes : True lorsque les éléments du DataFrame d'origine sont NaN (manquants) et False dans le cas contraire.
- `sum()` : Lorsqu'elle est appliquée au DataFrame renvoyé par `isnull()`, cette méthode compte le nombre de valeurs True dans chaque colonne. Essentiellement, elle additionne le nombre de valeurs manquantes pour chaque colonne du DataFrame.

CHAPITRE 2 CONCEPTS CLES SUR MACHINE LEARNING



2.4 Exemple Pratique de la Prédiction

2.4.1. Régression linéaire simple [Simple Linear Regression]

Exemple 1 : Prédiction du pourcentage de notes qu'un étudiant devrait obtenir en fonction du nombre d'heures qu'il a étudiées

4. Diviser l'ensemble de données

❑ Avant de diviser l'ensemble de données, nous devons séparer les variable indépendantes (Feature: X) et la variable dépendante (Target: y).

```
[3]: 1 X = dataset.iloc[ : , :1].values  
     2 y = dataset.iloc[ : , 1].values
```

```
[4]: 1 from sklearn.model_selection import train_test_split  
     2  
     3 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=0)
```

X : l'ensemble de fonctionnalités (données d'entrée).

y : les étiquettes cibles (données de sortie).

- **train_test_split** : une fonction qui divise les tableaux ou les matrices en sous-ensembles aléatoires d'entraînement et de test.
- **test_size=0.2** : cela signifie que 20 % des données seront utilisées pour les tests et 80 % pour l'entraînement.
- **random_state=0** : garantit que la division est reproductible. La même valeur `random_state` produira toujours la même division.
- **X_train, X_test** : il s'agit des sous-ensembles de l'ensemble de fonctionnalités pour l'entraînement et les tests, respectivement.
- **y_train, y_test** : il s'agit des sous-ensembles correspondants des étiquettes cibles pour l'entraînement et les tests, respectivement.

- **dataset.iloc[:, :-1]**: Ceci sélectionne toutes les colonnes de la dataset sauf la dernière (c'est-à-dire toutes les colonnes de Features : X).
- **.values**: Convertit les colonnes sélectionnées (généralement au format DataFrame) en un tableau NumPy.
- **dataset.iloc[:, -1]**: Ceci sélectionne la dernière colonne de la dataset (qui est souvent la colonne cible ou d'étiquette : y).

CHAPITRE 2 CONCEPTS CLES SUR MACHINE LEARNING



2.4 Exemple Pratique de la Prédiction

2.4.1. Régression linéaire simple [Simple Linear Regression]

Exemple 1 : Prédiction du pourcentage de notes qu'un étudiant devrait obtenir en fonction du nombre d'heures qu'il a étudiées

5. Ajustement du modèle de régression linéaire simple aux données d'entraînement

```
1 from sklearn.linear_model import LinearRegression
2
3 simple_linear_regression = LinearRegression()
4 simple_linear_regression = simple_linear_regression.fit(X_train, y_train)
```

- **from sklearn.linear_model import LinearRegression** : Ceci importe la classe **LinearRegression** du module **sklearn.linear_model**. Cette classe est utilisée pour effectuer une régression linéaire, qui est une méthode pour modéliser la relation entre une variable dépendante et une ou plusieurs variables indépendantes.
- **LinearRegression()** : Cela crée une instance de la classe **LinearRegression**, initialisant un nouveau modèle de régression linéaire.
- **simple_linear_regression** : Il s'agit de la variable qui stocke l'instance du modèle de régression linéaire.
- **simple_linear_regression.fit(X_train, y_train)** : Cette méthode entraîne le modèle de régression linéaire à l'aide des données d'entraînement. Elle trouve la ligne la mieux ajustée qui minimise l'erreur quadratique moyenne entre les valeurs prédites et les valeurs réelles dans les données d'entraînement.
- **X_train** : Les données d'entraînement pour les caractéristiques (variables indépendantes).
- **y_train** : Les données d'entraînement pour la variable cible (variable dépendante).

La méthode d'ajustement ajuste les paramètres du modèle (coefficients) en fonction des données d'entraînement.

CHAPITRE 2 CONCEPTS CLES SUR MACHINE LEARNING



2.4 Exemple Pratique de la Prédiction

2.4.1. Régression linéaire simple [Simple Linear Regression]

Exemple 1 : Prédiction du pourcentage de notes qu'un étudiant devrait obtenir en fonction du nombre d'heures qu'il a étudiées.

5. Prédire les résultats

```
1 y_pred = simple_linear_regression.predict(X_test)
```

Le code `y_pred = simple_linear_regression.predict(X_test)` est utilisé pour faire des prédictions sur les données de test en utilisant le modèle de régression linéaire entraîné.

- `simple_linear_regression.predict(X_test)` : Cette méthode utilise le modèle de régression linéaire formé (stocké dans `simple_linear_regression`) pour prédire les valeurs cibles des données de test (`X_test`).
- `y_pred` : Les valeurs prédites sont stockées dans la variable `y_pred`.

CHAPITRE 2 CONCEPTS CLES SUR MACHINE LEARNING



2.4 Exemple Pratique de la Prédiction

2.4.1. Régression linéaire simple [Simple Linear Regression]

Exemple 1 : Prédiction du pourcentage de notes qu'un étudiant devrait obtenir en fonction du nombre d'heures qu'il a étudiées.

5. Visualisation des résultats d'entraînement

```
1 plt.scatter(X_train, y_train, color='red') # Training Data
2 plt.plot(X_train, simple_linear_regression.predict(X_train), color='blue')
3 plt.title('Score Vs Hours (Trainig set)')
4 plt.xlabel("Score")
5 plt.ylabel('Hours')
6 plt.show()
```

Nous avons créé un scatter des points de données d'entraînement ainsi que la ligne de régression prédite par le modèle de régression linéaire entraîné. Expliquons maintenant chaque ligne de codes :

- **plt.scatter(X_train, y_train, color='red')** : cette ligne crée un nuage de points des données d'entraînement.
- **X_train** : les valeurs caractéristiques (par exemple, les heures d'étude).
- **y_train** : les valeurs cibles correspondantes (par exemple, les scores).
- **color='red'** : définit la couleur des points de données en rouge.
- **plt.plot(X_train, simple_linear_regression.predict(X_train), color='blue')** : cette ligne trace la ligne de régression prédite par le modèle de régression linéaire entraîné.
- **simple_linear_regression.predict(X_train)** : valeurs cibles prédites à l'aide du modèle de régression linéaire.

CHAPITRE 2 CONCEPTS CLES SUR MACHINE LEARNING



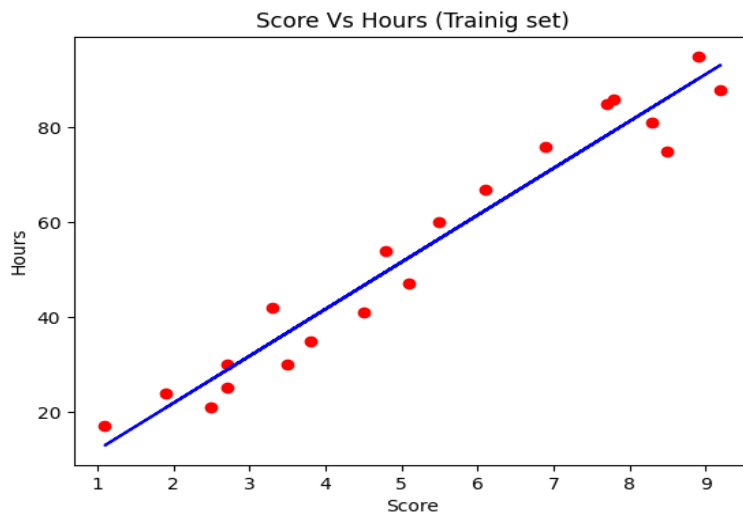
2.4 Exemple Pratique de la Prédiction

2.4.1. Régression linéaire simple [Simple Linear Regression]

Exemple 1 : Prédiction du pourcentage de notes qu'un étudiant devrait obtenir en fonction du nombre d'heures qu'il a étudiées.

5. Visualisation des résultats d'entraînement

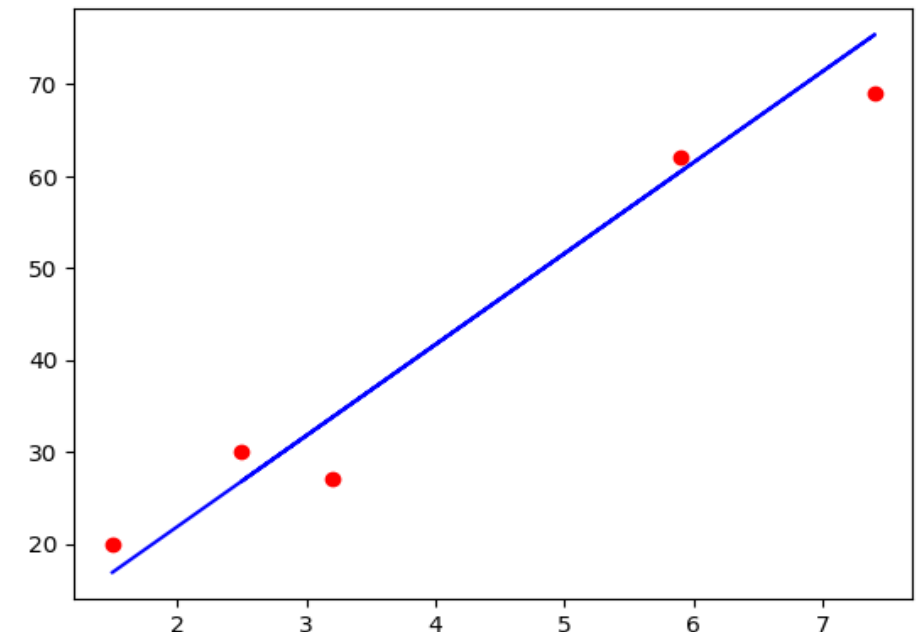
- `color='blue'` : Définit la couleur de la ligne de régression en bleu.
- `plt.title('Score Vs Hours (Training set)')` : Définit le titre du graphique.
- `plt.xlabel("Score")` : Définit l'étiquette de l'axe des x.
- `plt.ylabel('Hours')` : Définit l'étiquette de l'axe des y.
- `plt.show()` : Affiche le tracé avec les fonctionnalités spécifiées, la ligne de régression, le titre et les étiquettes.



6. Visualisation des résultats de test

```
1 plt.scatter(X_test, y_test, color='red') # Training Data
2 plt.plot(X_test, simple_linear_regression.predict(X_test), color='blue')
```

[<matplotlib.lines.Line2D at 0x2530e3d36d0>]



CHAPITRE 2 CONCEPTS CLES SUR MACHINE LEARNING



2.5 Exemple Pratique de la Prédiction

2.5.1 Régression linéaire Multiple [Multiple Linear Regression]

- ❑ Dans ce chapitre, nous allons découvrir la régression linéaire multiple à l'aide de scikit-learn dans le langage de programmation Python.
- ❑ La régression linéaire multiple, souvent appelée régression multiple, est une méthode statistique qui prédit le résultat d'une variable de réponse en combinant de nombreuses variables explicatives.
- ❑ La régression multiple est une variante de la régression linéaire dans laquelle une seule variable explicative est utilisée.

Formulation Mathématique:

$$y = \alpha + \beta_1(x_1) + \beta_2(x_2) + \dots + \beta_n(x_n)$$

Diagram illustrating the components of the multiple linear regression equation:

- y : Predicted value
- α : Bias
- β_1 : Weight 1
- x_1 : Feature 1
- β_2 : Weight 2
- x_2 : Feature 2
- β_n : Weight n
- x_n : Feature n

$$y = \beta_0 + \beta_1x_1 + \beta_2x_2 + \beta_3x_3 \dots$$

ici, y est la variable dépendante.

x_1, x_2, x_3, \dots sont des variables indépendantes.

β_0 = interception de la droite.

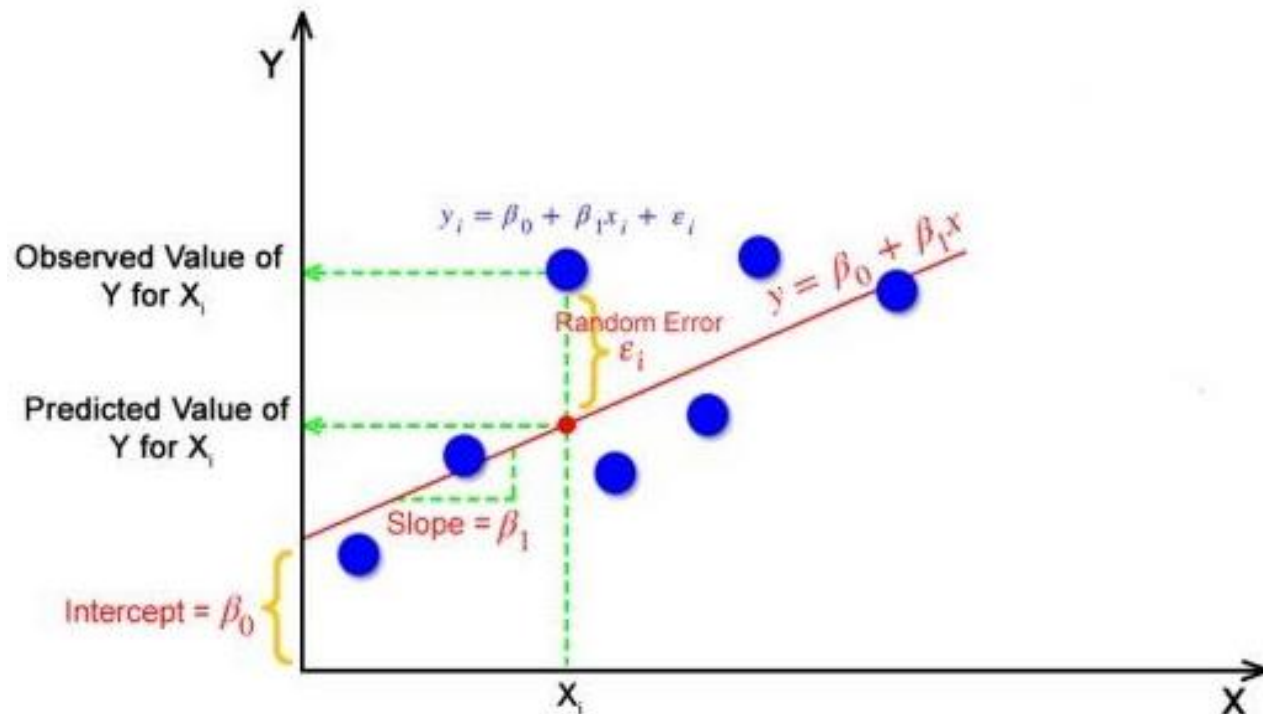
β_1, β_2, \dots sont des coefficients.

CHAPITRE 2 CONCEPTS CLES SUR MACHINE LEARNING



2.5 Example Pratique de la Prédiction

- ❑ La régression linéaire multiple, souvent appelée régression multiple, est une méthode statistique qui prédit le résultat d'une variable de réponse en combinant de nombreuses variables explicatives.



Formulation Mathématique:

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \dots$$

ici, **y** est la variable dépendante.

x₁, x₂, x₃,... sont des variables indépendantes.

b₀ = interception de la droite.

b₁, b₂, ... sont des coefficients.

Regression Linear Multiple [Lab]

Housing