



---

# PROSIT 1

---

Intelligence Artificielle



JEUDI 2 JUIN 2022

CESI

Mathilde Ballouhey

# Table des matières

<b>TABLE DES MATIERES .....</b>	<b>1</b>
<b>1. CONTEXTE.....</b>	<b>3</b>
<b>2. MOTS CLES.....</b>	<b>3</b>
2.1. INTELLIGENCE ARTIFICIELLE .....	3
2.2. DENDROGRAMME.....	3
2.3. K-MEANS .....	3
2.4. SECTEUR DE RECENSEMENT .....	3
2.5. ANALYSE INTELLIGENTE.....	3
2.6. ALGORITHME .....	3
2.7. CENTROÏDE .....	4
2.8. ALGORITHME DE CLASSIFICATION ASCENDANTE HIERARCHIQUE .....	4
2.9. ALGORITHME SUPERVISE / NON-SUPERVISE.....	4
<b>3. BESOINS / PROBLEMATIQUES.....</b>	<b>4</b>
<b>4. CONTRAINTES.....</b>	<b>4</b>
<b>5. LIVRABLES.....</b>	<b>4</b>
<b>6. GENERALISATION .....</b>	<b>4</b>
<b>7. HYPOTHESES.....</b>	<b>4</b>
<b>8. PLAN D'ACTION .....</b>	<b>5</b>
<b>COURS.....</b>	<b>6</b>
<b>1. INTRODUCTION SUR L'INTELLIGENCE ARTIFICIELLE .....</b>	<b>6</b>
1.1. POINT D'HISTOIRE.....	6
1.2. DEFINITION .....	7
1.3. TECHNIQUE .....	8
<b>2. ALGORITHMES DE CLASSIFICATION .....</b>	<b>11</b>
2.1. CONCEPT DU D'APPRENTISSAGE AUTOMATIQUE.....	11
2.2. APPRENTISSAGE SUPERVISE .....	11
2.3. APPRENTISSAGE NON SUPERVISE.....	12
2.4. ALGORITHME DE CLASSIFICATION ASCENDANTE HIERARCHIQUE .....	12
2.5. ALGORITHME K-MEANS .....	13
<b>3. IA &amp; ETHIQUE.....</b>	<b>15</b>
3.1. BIAIS ALGORITHMIQUE : .....	15
3.2. CONFIDENTIALITE DES DONNEES : .....	15
3.3. RESPONSABILITE : .....	15
3.4. IMPACT SUR L'EMPLOI : .....	15
3.5. UTILISATION MILITAIRE : .....	15
<b>4. REGLEMENTATION .....</b>	<b>17</b>
<b>5. PYTHON.....</b>	<b>18</b>
5.1. TENSORFLOW .....	18
5.2. KERAS .....	18
5.3. PYTORCH .....	18
5.4. SCIKIT-LEARN .....	18
5.5. THEANO .....	18
5.6. CAFFE.....	18
5.7. NLTK .....	18
<b>6. PRE-TRAITEMENT D'INFORMATIONS .....</b>	<b>19</b>

6.1.	NETTOYAGES DES DONNEES .....	20
6.2.	PIPELINE.....	21
6.2.1.	<i>Détail au sein du workshop &amp; livrable.</i> .....	22
<b>I-</b>	<b>HYPOTHESES.....</b>	<b>23</b>
	<b>NOUS ALLONS UTILISER UNE ECHELLE POUR PRETRAITER NOS DONNEES .....</b>	<b>23</b>
<b>7.</b>	<b>LIVRABLE .....</b>	<b>23</b>
<b>8.</b>	<b>VALIDATION DES HYPOTHESES.....</b>	<b>23</b>
	<b>BIBLIOGRAPHIE .....</b>	<b>24</b>

# 1. Contexte

Une agence immobilière voudrait optimiser sa proposition de biens aux futurs acheteurs tout en minimisant le nombre de visites en utilisant l'intelligence artificielle

## 2. Mots clés

### 2.1. Intelligence artificielle

L'intelligence artificielle (IA) est un ensemble de techniques et de méthodes informatiques visant à créer des machines capables de simuler l'intelligence humaine. Elle comprend notamment des domaines tels que l'apprentissage automatique, la reconnaissance d'image, la compréhension du langage naturel, la robotique et la prise de décision autonome. L'objectif principal de l'IA est de permettre aux machines d'effectuer des tâches qui nécessitent normalement une intervention humaine, voire même de surpasser les capacités humaines dans certains domaines.

### 2.2. Dendrogramme

Un dendrogramme est un diagramme arborescent utilisé en analyse de données pour représenter graphiquement la similarité entre différents objets ou groupes d'objets. Les objets sont représentés par des feuilles et les groupes d'objets sont représentés par des noeuds internes, qui sont reliés entre eux par des branches qui reflètent la distance ou la similarité entre les objets ou groupes. Le dendrogramme est souvent utilisé pour visualiser les résultats d'une analyse de classification ou de regroupement de données.

### 2.3. K-means

K-means est un algorithme de regroupement non supervisé utilisé en analyse de données pour diviser un ensemble de données en K groupes (clusters) distincts. L'algorithme fonctionne en assignant chaque observation à l'un des K clusters de manière à minimiser la variance intra-cluster, c'est-à-dire la somme des distances entre chaque observation et le centre de son cluster. Les centres de chaque cluster sont mis à jour itérativement jusqu'à ce que la variance intra-cluster soit minimisée. K-means est largement utilisé dans de nombreux domaines, y compris la biologie, l'informatique, la finance et la recherche en marketing.

### 2.4. Secteur de recensement

Un secteur de recensement est une subdivision géographique utilisée par Statistique Canada pour la collecte de données lors des recensements de la population. Chaque secteur de recensement est défini comme une zone géographique homogène d'un point de vue socio-économique et démographique, et contient généralement entre 400 et 700 résidences. Les données collectées dans chaque secteur de recensement comprennent des informations sur la population, le logement, l'éducation, le revenu et l'emploi, entre autres. Les secteurs de recensement sont utilisés pour produire des statistiques précises et fiables sur la population et les ménages à travers le pays.

### 2.5. Analyse intelligente

L'analyse intelligente est l'utilisation de techniques d'intelligence artificielle pour extraire des connaissances à partir de grandes quantités de données, permettant de découvrir des tendances et des modèles cachés pour prendre des décisions éclairées dans de nombreux domaines.

### 2.6. Algorithme

Un algorithme est une suite d'instructions ordonnées et finies permettant de résoudre un problème ou accomplir une tâche spécifique. Il peut être utilisé dans de nombreux domaines, tels que l'informatique, les mathématiques ou l'ingénierie.

## 2.7. Centroïde

Le centroïde est le point de rencontre de toutes les médianes d'un triangle, c'est-à-dire le point où ces médianes se croisent. Il est également considéré comme le centre de gravité du triangle.

## 2.8. Algorithme de classification ascendante hiérarchique

L'algorithme de classification ascendante hiérarchique (CAH) est une méthode d'analyse de données qui permet de regrouper des individus ou des variables en fonction de leurs similarités. Cette méthode consiste à construire une hiérarchie de groupes en fusionnant itérativement les groupes les plus proches jusqu'à l'obtention d'un groupe unique contenant tous les individus ou toutes les variables. Le choix du critère de similarité et de la méthode de fusion dépendent des objectifs et des caractéristiques des données à analyser.

## 2.9. Algorithme supervisé / non-supervisé

Un algorithme supervisé est un algorithme qui utilise un ensemble de données étiquetées pour apprendre à prédire des étiquettes sur de nouvelles données. En revanche, un algorithme non-supervisé est un algorithme qui explore les données sans étiquettes pour en extraire des structures, des relations ou des regroupements.

# 3. Besoins / Problématiques

Quel type de prétraitement appliquer pour optimiser un jeu de données pour les résultats d'un algorithme de K-means ?

# 4. Contraintes

- Utilisation d'un algorithme de K-means
- Le jeu de donnée est imposé / donné
- La zone géographique est imposée

# 5. Livrables

Notebook Jupiter contenant les étapes suivantes :

- Prétraiter le jeu de donnée afin que les résultats de l'algorithme de K-means soient meilleurs
- Exécuter le K-means
- Comparer les résultats (test de performances)

# 6. Généralisation

Prétraitement des informations avant son exploitation dans un algorithme.

# 7. Hypothèses

- K-means est un algorithme de classification
- Un algorithme supervisé demande une activité humaine régulière
- Un algorithme non supervisé tourne et s'entraîne seul pendant X temps
- K-means ne va pas forcément affiner les résultats
- Le K-means est mieux que l'algorithme de classification ascendante hiérarchique
- Il existe une classification descendante hiérarchique
- Ce jeu de donnée pourrai être mieux traité par un autre algorithme
- La métrique utilisée serai la fonction perte / coût (loss function)
- Nous allons utiliser une échelle pour prétraiter nos données

## 8. Plan d'action

- Mots clés
- Réalisation des cours
  - Introduction sur L'intelligence Artificielle
  - Algorithme de classification
    - Supervisé / non supervisé
    - K-means + Algorithme de classification ascendante hiérarchique
  - L'intelligence artificielle et l'éthique (réglementation de la commission européenne / CNIL)
  - Les bibliothèques en python pour l'Intelligence Artificielle (scikit learn)
  - Prétraitement / préparation des données
- Workshop
- Réalisation du livrable
- Traitement des hypothèses

# Cours

---

## 1. Introduction sur L'Intelligence artificielle

### 1.1. Point d'histoire

L'histoire de l'intelligence artificielle (IA) s'étend sur plusieurs décennies, avec des avancées technologiques significatives et des périodes de développement intense. Voici un bref résumé des étapes clés de l'histoire de l'IA :

- 1943-1956 : Les prémices de l'IA : Pendant cette période, des chercheurs tels que Warren McCulloch et Walter Pitts ont créé les premiers modèles de réseaux de neurones artificiels. En 1956, John McCarthy, Marvin Minsky, Claude Shannon et Nathaniel Rochester ont organisé la conférence de Dartmouth, qui est considérée comme le début de l'IA en tant que domaine de recherche distinct.
- 1956-1974 : L'âge d'or de l'IA : Pendant cette période, les chercheurs ont développé des algorithmes de recherche et de planification, ainsi que des programmes de traduction automatique. Des systèmes comme le théorème de géométrie de Prover ont été développés et ont été considérés comme des percées importantes.
- 1980-2010 : L'IA symbolique et les réseaux de neurones : Pendant cette période, l'IA symbolique a été largement utilisée pour la résolution de problèmes et l'expertise en domaine. Les réseaux de neurones ont également connu un regain d'intérêt, avec des avancées telles que le réseau de neurones convolutionnels et le réseau de neurones récurrents.
- Depuis 2010 : L'IA profonde : Au cours de la dernière décennie, l'IA profonde a connu un énorme succès grâce à l'utilisation de réseaux de neurones profonds, de l'apprentissage par renforcement et de la capacité croissante de calcul. Des systèmes tels que AlphaGo de Google DeepMind ont montré la capacité de l'IA à surpasser les humains dans des tâches complexes.

Voici quelques exemples d'IA marquante de ces époques :

- 1951 : UNIVAC, le premier ordinateur commercial, a été lancé.
- 1952 : Arthur Samuel développe un programme de dames qui peut apprendre de ses propres erreurs.
- 1956 : Le langage de programmation LISP est développé pour la recherche en IA.
- 1967 : Le système DENDRAL est développé pour la reconnaissance de structures moléculaires.
- 1972 : Le système MYCIN est développé pour le diagnostic de maladies infectieuses.
- 1985 : Le système XCON est développé pour la configuration de produits personnalisés.
- 1997 : Deep Blue de IBM bat le champion du monde d'échecs Garry Kasparov.
- 2011 : Watson de IBM bat les meilleurs joueurs humains au jeu télévisé Jeopardy!
- 2016 : Alpha GO de Google DeepMind bat le champion du monde de Go Lee Sedol.
- 2017 : AlphaGo Zero de Google DeepMind bat la version précédente d'AlphaGo sans utiliser de données d'entraînement humaines, mais en utilisant uniquement l'apprentissage par renforcement.
- 2018 : GPT-2 de OpenAI est un modèle de langage naturel qui peut générer du texte humain-like de qualité élevée, ce qui soulève des préoccupations quant à l'utilisation abusive de cette technologie pour la désinformation ou les attaques de phishing.

- 2019 : OpenAI développe un système d'IA de traitement du langage naturel nommé GPT-3, qui est capable d'effectuer une variété de tâches telles que la traduction, la génération de texte, la réponse à des questions et la création de code informatique.
- 2020 : AlphaFold de DeepMind utilise l'IA pour prédire avec précision la structure tridimensionnelle des protéines, une percée qui pourrait avoir un impact considérable sur la recherche en médecine et la conception de médicaments.
- 2021 : DALL-E de OpenAI est un modèle de génération d'images qui peut créer des images à partir de descriptions textuelles. Cette technologie a le potentiel de révolutionner la conception graphique et la production d'images pour les applications de réalité virtuelle et augmentée.

Il convient de noter que ces exemples ne représentent qu'une petite partie des réalisations notables en IA au fil des ans, mais ils illustrent la variété des domaines d'application de l'IA et la rapidité avec laquelle la technologie a progressé au fil des ans.

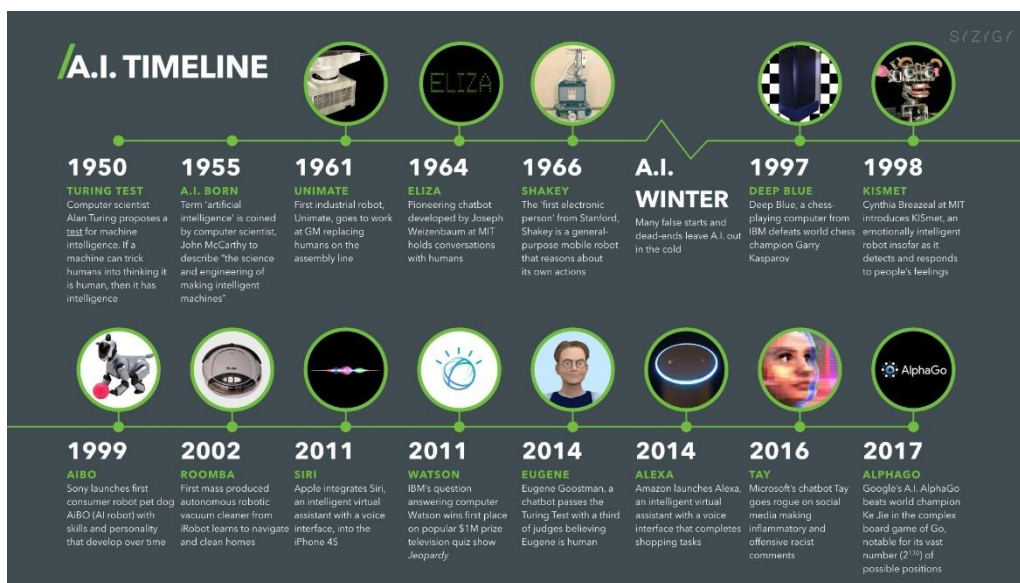


Figure 1 : Fiche chronologie des IA

## 1.2. Définition

L'**intelligence artificielle (IA)** est un domaine de l'informatique qui consiste à développer des systèmes informatiques capables d'accomplir des tâches qui nécessitent normalement l'intelligence humaine, telles que la **reconnaissance de la parole**, la **vision par ordinateur**, la **prise de décision** et le **raisonnement**. Les systèmes d'IA utilisent des **algorithmes** sophistiqués pour apprendre à partir de données, à travers des méthodes telles que l'**apprentissage automatique** (Machine Learning) et l'**apprentissage profond** (Deep Learning).

Au fil des années, l'IA a progressé de manière significative, permettant d'avancées dans de nombreux domaines tels que la médecine, l'industrie, les transports et la finance. Elle est maintenant utilisée pour résoudre des problèmes complexes tels que la détection de fraudes, la prédiction des maladies et la recommandation de produits personnalisés.

Cependant, malgré ces avancées, l'IA soulève également des questions et des défis **éthiques**, notamment en ce qui concerne la protection de la **vie privée**, la prise de décision **équitable** et **transparente**, et la **responsabilité** en cas d'erreur ou de biais.

Il est donc crucial de continuer à explorer et à approfondir l'IA, afin de maximiser ses avantages tout en minimisant ses risques. Les avancées futures pourraient inclure des systèmes encore plus sophistiqués, des améliorations de l'interaction homme-machine et des avancées dans la



compréhension des processus cognitifs humains. En somme, l'IA est un domaine passionnant et en constante évolution qui offre de nombreuses possibilités pour l'avenir.

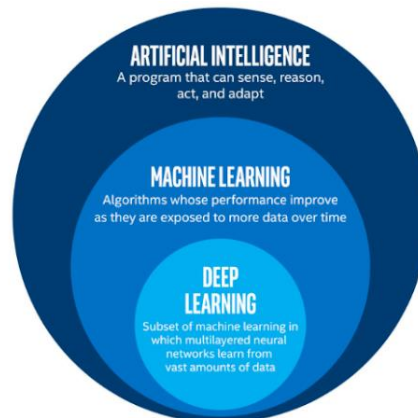


Figure 2 : Catégorisation d'apprentissage des intelligences artificielles

### 1.3. Technique

Pour comprendre comment fonctionne une IA, il est important de comprendre d'abord les différentes techniques utilisées en intelligence artificielle.

**L'apprentissage automatique** (Machine Learning en anglais) est l'une de ces techniques. Il consiste à entraîner un algorithme à partir d'un grand nombre de données pour lui permettre de généraliser et de faire des prédictions sur des données inconnues. Pour cela, on utilise souvent des modèles de réseaux de neurones artificiels, qui sont des structures informatiques composées de couches de neurones interconnectés.

**L'apprentissage profond** (Deep Learning en anglais) est une variante de l'apprentissage automatique qui utilise des réseaux de neurones profonds, c'est-à-dire avec de nombreuses couches intermédiaires, pour extraire des caractéristiques de plus en plus abstraites à partir des données d'entrée.

Les systèmes d'IA peuvent également utiliser des techniques de **traitement du langage naturel** (NLP en anglais) pour comprendre et générer des textes ou des conversations en langage naturel.

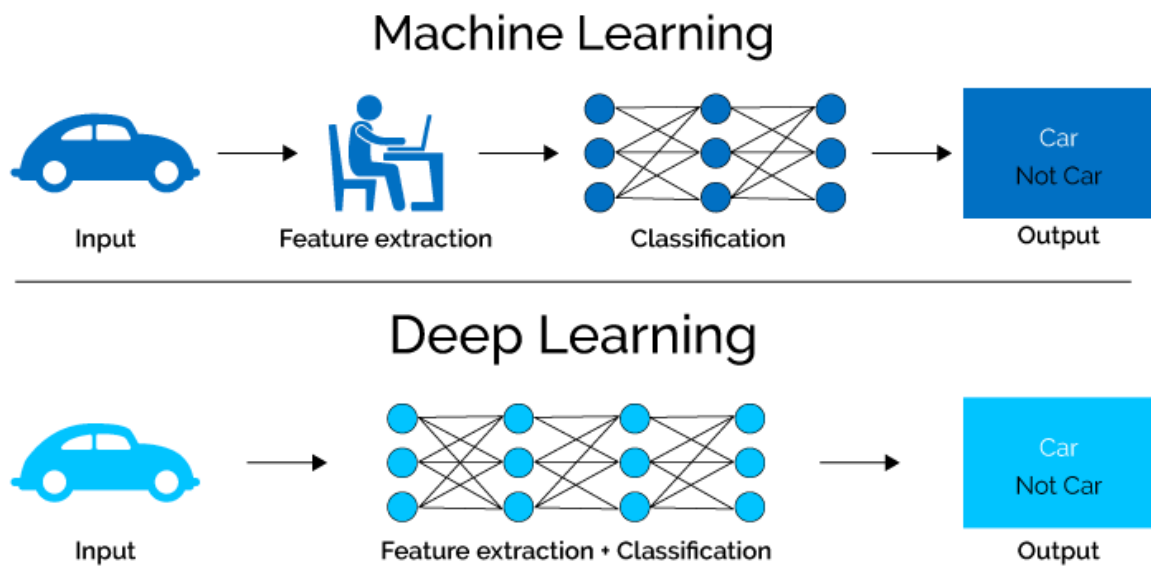
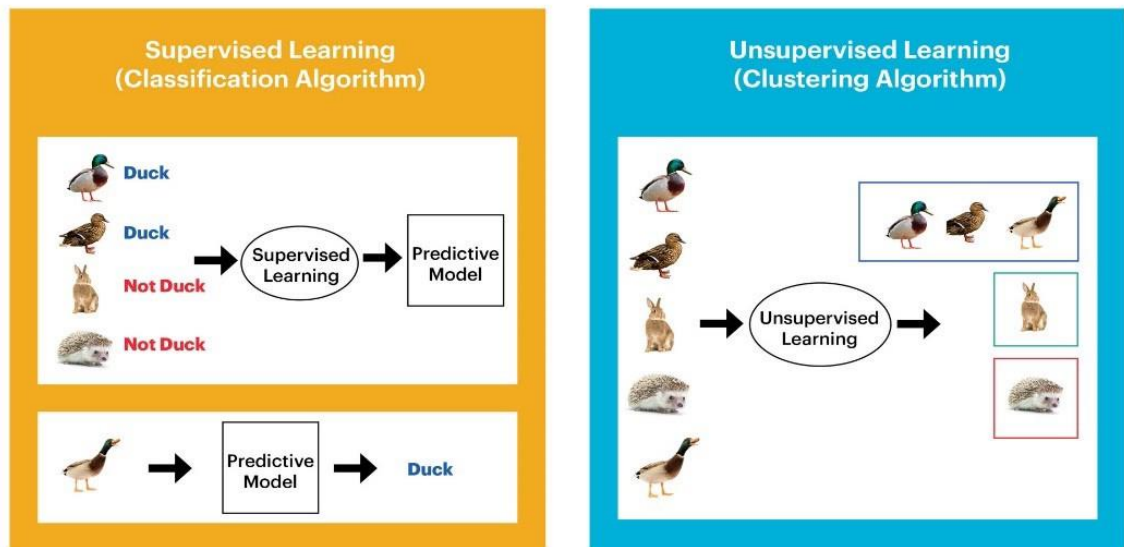


Figure 3 : schéma différence entre Machine Learning & Deep Learning

Enfin, il est important de noter que les systèmes d'IA peuvent être **supervisés** ou **non-supervisés**. Dans l'apprentissage supervisé, les données d'entrée sont étiquetées avec des réponses attendues, ce qui permet au modèle de s'entraîner à prédire ces réponses. Dans l'apprentissage non supervisé, les données ne sont pas étiquetées, et le modèle doit trouver des structures ou des similarités dans les données.

En résumé, les systèmes d'IA fonctionnent en utilisant des techniques d'apprentissage automatique, d'apprentissage profond et de traitement du langage naturel pour extraire des connaissances à partir des données. Ces systèmes peuvent être supervisés ou non supervisés, en fonction du type de données et des objectifs de l'application.



Western Digital.

Figure 4 : schéma différence entre apprentissage supervisés et apprentissage supervisés

## 2. Algorithmes de classification

Avant d'aborder le concept même d'algorithme de classification via deux exemples, il convient de rappeler le concept de Machine Learning et la différence entre les deux principaux types d'apprentissage :

- L'apprentissage supervisé
- L'apprentissage non supervisé

### 2.1. Concept du d'apprentissage automatique

L'apprentissage automatique (Machine Learning en anglais) est une catégorie d'algorithme qui permet aux applications de prédire de façon plus précise les résultats sans qu'ils soient programmées de façon explicite. Le principe de base est de créer des algorithmes capables de recevoir des données et d'utiliser une analyse statistique pour prédire un résultat de sortie tout en mettant à jour l'entrée à mesure que de nouvelles données deviennent disponibles.

### 2.2. Apprentissage supervisé

Une grande partie des apprentissages automatiques utilisent ce que l'on appelle l'apprentissage supervisé. Ce dernier consiste en divers variables d'entrée et une variable de sortie, le but étant d'appréhender la fonction de mapping de manière que, lorsque de nouvelles données d'entrée arrivent, il est possible de prédire les variables de sortie pour ces mêmes données.

On dit qu'il est supervisé car le processus ou training set peut être visualiser comme un enseignant qui supervise le processus d'apprentissage de l'algorithme. En résumé, les réponses correctes sont connues de l'enseignant, l'algorithme allant effectuer des prédictions sur les données qui sont corrigés par l'enseignant. L'apprentissage s'arrête lorsqu'un certain niveau de performance est atteint.

Il y a deux contextes qui peuvent permet d'effectuer un apprentissage supervisé :

- **La classification :** Un problème de classification survient lorsque la variable de sortie est une catégorie, telle que « rouge », « bleu » ou « maladie » et « pas de maladie ». Exemples :
  - En finance et dans le secteur bancaire pour la détection de la fraude par carte de crédit (fraude, pas fraude).
  - Détection de courrier électronique indésirable (spam, pas spam).
  - Dans le domaine du marketing utilisé pour l'analyse du sentiment de texte (heureux, pas heureux).
  - En médecine, pour prédire si un patient a une maladie particulière ou non.
- **La régression :** Un problème de régression se pose lorsque la variable de sortie est une valeur réelle, telle que « dollars » ou « poids ». Exemples :
  - Prédire le prix de l'immobilier
  - Prédire le cours de bourse
  - Certains types courants de problèmes fondés sur la classification et la régression incluent la prévision et la prévision de séries temporelles, respectivement.

Voici une liste de quelques algorithmes d'apprentissage automatique supervisé :

- Arbres de décision
- K Nearest Neighbours
- SVC linéaire (classificateur de vecteur de support)
- Régression logistique
- Naive Bayes

- Les réseaux de neurones
- Régression linéaire
- Régression vectorielle de support (SVR)
- Arbres de régression

### 2.3. Apprentissage non supervisé

Contrairement à son homologue supervisé, l'apprentissage dit non supervisé ne dispose pas de variables de sortie correspondantes. L'objectif de ce genre d'apprentissage est de modéliser la structure ou la distribution sous-jacente dans les données afin d'en apprendre davantage sur ces dernières.

Il est appelé apprentissage non supervisé car il n'y a pas de réponse correcte connue mais aussi pas d'enseignant. Les algorithmes de ce style d'apprentissage sont laissés à leurs mécanismes pour découvrir et présenter la structure intéressante des données.

Il existe, dans l'apprentissage non supervisé, deux catégories d'algorithmes :

- Les algorithmes de regroupement (clustering) : L'objectif est de séparer les groupes qui ont des traits similaires et de les assigner en grappes.
- Les algorithmes d'associations : L'association consiste à découvrir des relations entre les attributs des points de données.

Parmi les divers algorithmes d'apprentissage automatique non supervisé, deux vont nous intéresser :

- La classification ascendante hiérarchique
- Le K-means clustering

### 2.4. Algorithme de classification ascendante hiérarchique

Étant donnés des points et un entier  $k$ , l'algorithme CAH vise à diviser les points en  $k$  groupes, appelés clusters, homogènes et compacts. Concrètement comment s'y prend-on ?

- L'idée de départ est de considérer que chacun des points de votre jeu de données est un centroïde. Cela revient à considérer qu'à chaque point correspond une unique étiquette (0,1,2,3, 4...).
- Ensuite on regroupe chaque centroïde avec son centroïde voisin le plus proche. Ce dernier prend l'étiquette du centroïde qui l'a « absorbé ».
- On calcule alors les nouveaux centroïdes qui seront les centres de gravité des clusters nouvellement créés.
- On réitère l'opération jusqu'à obtenir un unique cluster ou bien un nombre de clusters préalablement défini.

Il y a trois points clés dans cet algorithme :

- Quelle est la métrique utilisée pour évaluer la distance entre les centroïdes ?
- Quel est le nombre de clusters à choisir ?
- Sur quel critère décide-t-on de regrouper les centroïdes entre eux ?

Dans la CAH, on utilise la distance euclidienne pour pouvoir évaluer la distance entre les centroïdes :

$$d(p, q) = \sqrt{\sum_{i=1}^n (p_i - q_i)^2}$$

A chaque étape de regroupement de deux centroïdes on obtiendra un nouveau cluster et un nouveau centroïde qui n'est autre que le centre de gravité du nuage de points.

Une notion qui est intéressante à noter dans cet algorithme est l'inertie intra classe. Cela correspond à la somme des distances euclidiennes entre chaque point associé à un regroupement et également au centre de gravité nouvellement calculé. L'idée principale est de minimiser l'augmentation évidente de cette inertie en regroupant des clusters entre eux en un nouveau cluster.

Pour un exemple avec un code python : [datascientest.com](https://datascientest.com)

L'algorithme CAH a quelques limites cependant. Tout d'abord, comme on demande en amont un nombre de partitions défini, cela demande d'avoir une idée précise sur ce nombre et même s'il existe une méthode pour optimiser ce dernier sur certains jeux de donnée, elle est loin d'être infaillible.

De plus, la Classification Ascendante Hiérarchique peut être très coûteuse en temps et en ressources mais il existe un moyen de contourner le problème en renseignant une matrice de connectivité (matrice creuse indiquant quelles paires d'observation sont voisines).

Enfin, suivant les données qui sont entrées, le CAH ne sera pas aussi simple car il faudra retravailler ces dernières pour avoir un résultat acceptable.

## 2.5. Algorithme K-means

Si on a divers points et un entier noté k, l'algorithme des k-means vise, comme le CAH, à diviser ces points en k groupes homogènes et compacts. Pour réaliser cela, il y a trois étapes à avoir :

- Définition aléatoire de trois centroïdes auxquelles on va associer des étiquettes ;
- Observations des distances entre chaque point et les trois centroïdes et associations de chacun d'eux avec l'étiquette du centroïde le plus proche ;
- Calcul de trois nouveaux centroïdes qui deviendront les centres de gravité de chaque nuage de points obtenu et étiqueté

Ces trois étapes sont répétées jusqu'à ce que les nouveaux centroïdes calculés ne varient plus.

Ainsi, il y a deux points clés dans cet algorithme :

- Quelle est la métrique utilisée pour évaluer la distance entre les points et les centroïdes ?
- Quel est le nombre de clusters à choisir ?

Tout comme dans le CAH, l'algorithme des k-means utilise la distance euclidienne. Pour rappel, elle permet d'évaluer la distance entre chaque point et les centroïdes. Pour chaque point on calcule la distance euclidienne entre ce point et chacun des centroïdes puis on l'associe au centroïde le plus proche c'est-à-dire celui avec la plus petite distance.

$$d(p, q) = \sqrt{\sum_{i=1}^n (p_i - q_i)^2}$$

Pour que l'algorithme puisse donner des résultats satisfaisants, il faut que l'on puisse déterminer le nombre de clusters idéal. Ceci est possible via diverses méthodes. Celle qui peut être intéressante

est la méthode du coude qui s'appuie sur la notion d'inertie (somme des distances euclidiennes entre chaque point et son centroïde associé). Logiquement, plus le nombre initial de clusters est élevé, plus l'inertie va se réduire et donc plus les points auront la chance d'être proche d'un centroïde.

Une autre méthode serait d'utiliser le coefficient de silhouette :

$$s = \frac{b - a}{\max(a, b)}$$

- $a$  = la moyenne des distances aux autres observations du même cluster (distance intra-cluster).
- $b$  = la distance moyenne au cluster le plus proche.

Le coefficient varie entre -1 et +1. Plus il est proche de la valeur positive, plus l'observation est située à l'intérieur de son propre cluster. Au contraire, si la valeur est proche de -1, alors cela veut dire que l'observation est dans le mauvais cluster. Dans le cas où la valeur est 0, l'observation est située près d'une frontière.

En plus de ces deux méthodes, une analyse descriptive est également utile pour permettre de déterminer les caractéristiques communes de chaque cluster et ainsi de comprendre les profils types de chaque regroupement.

Tout comme le CAH, l'algorithme k-means a ses limites. Tout d'abord, comme pour la classification ascendante hiérarchique qui nécessite de définir le nombre de partitions, il faut avoir une idée précise de ce dernier et bien qu'il existe des méthodes pour déterminer ce nombre, elles sont loin d'être infallibles. De plus, il faudra souvent retravailler les données qui sont utilisées pour pouvoir obtenir des résultats satisfaisants.

*Exemple avec code python : [datascientest.com](https://datascientest.com)*

### 3. IA & éthique

La technologie de l'IA est en constante évolution et suscite à la fois des espoirs et des craintes. Bien que cette technologie puisse apporter de nombreux avantages à la société, elle peut également poser des problèmes éthiques.

Dans cette section, nous allons examiner certains des défis éthiques les plus importants liés à l'IA, ainsi que les approches possibles pour y faire face.

#### 3.1. Biais algorithmique :

L'un des problèmes les plus importants de l'IA est le biais algorithmique. Les algorithmes peuvent être biaisés lorsqu'ils sont formés sur des données qui reflètent les préjugés et les inégalités de la société. Si les données utilisées pour entraîner un algorithme contiennent des biais, l'algorithme peut reproduire ces biais dans ses résultats. Par exemple, un algorithme de recrutement formé sur des données historiques peut discriminer les candidats en fonction de leur race, de leur sexe ou de leur orientation sexuelle. Il est important de prendre des mesures pour atténuer les biais algorithmiques, telles que l'utilisation de données plus diversifiées et l'audit régulier des résultats de l'algorithme.

#### 3.2. Confidentialité des données :

L'IA nécessite souvent l'accès à de grandes quantités de données personnelles pour fonctionner efficacement. Cela soulève des questions importantes sur la confidentialité des données. Les entreprises et les organisations doivent veiller à ce que les données des utilisateurs soient stockées et traitées de manière sûre et transparente, conformément aux réglementations en matière de confidentialité des données.

#### 3.3. Responsabilité :

L'IA soulève également des questions sur la responsabilité. Qui est responsable si une décision prise par un algorithme se révèle préjudiciable ou injuste ? Les entreprises et les organisations doivent être responsables de l'utilisation de l'IA et prendre des mesures pour minimiser les risques liés à cette technologie.

#### 3.4. Impact sur l'emploi :

L'IA peut automatiser de nombreux emplois, ce qui peut entraîner des pertes d'emplois et des changements majeurs dans la structure du marché du travail. Les gouvernements et les entreprises doivent anticiper ces changements et mettre en place des politiques pour aider les travailleurs à s'adapter à l'évolution du marché du travail.

#### 3.5. Utilisation militaire :

Enfin, l'IA peut également être utilisée à des fins militaires, notamment dans le développement d'armes autonomes. Il est important de mettre en place des réglementations internationales pour contrôler l'utilisation de l'IA dans le domaine militaire et éviter des conséquences catastrophiques pour la sécurité internationale.



En conclusion, l'IA est une technologie puissante qui peut apporter de nombreux avantages à la société, mais elle soulève également des questions éthiques importantes. Il est crucial que les entreprises, les organisations et les gouvernements travaillent ensemble pour mettre en place des réglementations et des pratiques qui garantissent que l'IA est utilisée de manière éthique et responsable.

## 4. Réglementation

La CNIL (Commission Nationale de l'Informatique et des Libertés) est l'autorité française en charge de la protection des données personnelles. Elle a mis en place des réglementations pour encadrer l'utilisation de l'IA et protéger les données personnelles des utilisateurs. En France, les entreprises qui utilisent l'IA doivent respecter le Règlement Général sur la Protection des Données (RGPD), qui est une réglementation européenne en matière de protection des données. Le RGPD impose des obligations strictes aux entreprises en matière de collecte, de traitement et de stockage de données personnelles.

En ce qui concerne l'IA, le RGPD exige que les entreprises prennent des mesures pour minimiser les risques de discrimination et de préjudice, ainsi que pour garantir la transparence et la responsabilité de l'algorithme. En outre, la CNIL a publié des lignes directrices spécifiques pour l'utilisation de l'IA en matière de recrutement et d'analyse de données de santé. Ces lignes directrices énoncent des principes éthiques clés pour l'utilisation de l'IA dans ces domaines et prodiguent des conseils sur la manière de respecter les obligations légales en matière de protection des données. Il est important que les entreprises qui utilisent l'IA en France se conforment aux réglementations de la CNIL et respectent les principes éthiques qui sous-tendent ces réglementations. Cela contribuera à garantir que l'IA est utilisée de manière éthique et responsable et protégera les droits des utilisateurs de l'IA.

## 5. Python

Voici quelques bibliothèques Python qui peuvent être utiles pour l'IA :

### 5.1. TensorFlow

TensorFlow est une bibliothèque open source développée par Google Brain Team. Elle est principalement utilisée pour la conception et l'entraînement de modèles d'apprentissage automatique, en particulier les réseaux de neurones. TensorFlow est également utilisé pour les tâches de traitement du langage naturel et les systèmes de recommandation.

### 5.2. Keras

Keras est une bibliothèque open source qui fonctionne au-dessus de TensorFlow. Elle fournit une interface simple pour la conception de modèles d'apprentissage automatique, en particulier les réseaux de neurones. Keras est populaire pour sa facilité d'utilisation et sa documentation détaillée.

### 5.3. PyTorch

PyTorch est une bibliothèque open source développée par Facebook AI Research. Elle est principalement utilisée pour la conception et l'entraînement de modèles d'apprentissage automatique, en particulier les réseaux de neurones. PyTorch est également utilisé pour les tâches de traitement du langage naturel et la vision par ordinateur.

### 5.4. Scikit-learn

Scikit-learn est une bibliothèque open source qui fournit des algorithmes pour l'apprentissage automatique supervisé et non supervisé. Elle est utilisée pour les tâches de classification, de régression, de clustering et de réduction de dimensionnalité.

### 5.5. Theano

Theano est une bibliothèque open source qui fonctionne au-dessus de TensorFlow. Elle est principalement utilisée pour la conception et l'entraînement de modèles d'apprentissage.

### 5.6. Caffe

Caffe est une bibliothèque open source développée par le groupe de recherche en vision par ordinateur de l'Université de Californie à Berkeley. Elle est principalement utilisée pour la conception et l'entraînement de modèles d'apprentissage automatique, en particulier les réseaux de neurones. Caffe est également utilisé pour les tâches de vision par ordinateur.

### 5.7. NLTK

NLTK (Natural Language Toolkit) est une bibliothèque open source utilisée pour le traitement du langage naturel. Elle fournit des outils pour le prétraitement des données textuelles, la tokenisation, la lemmatisation, l'étiquetage de parties du discours, la reconnaissance d'entités nommées, l'analyse syntaxique et la classification de texte.

## 6. Pré-traitement d'informations

Le prétraitement de l'information est une étape clé dans la conception d'un algorithme d'intelligence artificielle. Cela implique de nettoyer, normaliser et transformer les données en entrée de manière à les rendre exploitables par l'algorithme.

Voici quelques exemples de techniques de prétraitement de données couramment utilisées en intelligence artificielle :

1. Nettoyage des données : il est fréquent que les données d'entrée soient incomplètes, bruyantes ou erronées. Le nettoyage des données consiste à éliminer les données redondantes, à remplacer les valeurs manquantes et à corriger les erreurs.
2. Normalisation : les données peuvent être mesurées dans des unités différentes ou avoir des échelles différentes. La normalisation consiste à mettre les données sur une même échelle pour faciliter leur traitement.
3. Réduction de la dimension : les données peuvent contenir des variables qui ne sont pas pertinentes pour l'analyse ou qui sont corrélées avec d'autres variables. La réduction de la dimension consiste à éliminer ces variables pour améliorer la performance de l'algorithme.
4. Transformation des données : les données peuvent être transformées pour mieux refléter les caractéristiques des données. Par exemple, des données non-linéaires peuvent être transformées en données linéaires pour permettre l'utilisation d'algorithmes linéaires.
5. Échantillonnage des données : les données peuvent être échantillonnées pour éviter le sur-apprentissage et améliorer la généralisation de l'algorithme.

Ces techniques de prétraitement de données sont souvent appliquées en combinaison les unes avec les autres pour obtenir les meilleurs résultats. L'objectif est d'obtenir des données de qualité pour que l'algorithme puisse apprendre à partir de ces données et produire des prévisions ou des recommandations précises et fiables.

## 6.1. Nettoyages des données

Le nettoyage des données est une étape cruciale dans la préparation des données pour l'analyse et la modélisation en intelligence artificielle. Les données brutes peuvent contenir des erreurs, des doublons, des valeurs manquantes ou des valeurs aberrantes qui peuvent affecter négativement la qualité de l'analyse. Le nettoyage des données est le processus de détection et de correction de ces erreurs pour s'assurer que les données sont de qualité suffisante pour être utilisées par l'algorithme.

Voici quelques exemples de techniques de nettoyage de données couramment utilisées en intelligence artificielle :

1. **Suppression des doublons** : les données peuvent contenir des enregistrements en double qui peuvent fausser les résultats de l'analyse. La suppression des doublons permet de s'assurer que chaque enregistrement est unique.
2. **Traitement des valeurs manquantes** : les données peuvent contenir des valeurs manquantes, qui peuvent être traitées en supprimant l'enregistrement correspondant, en remplissant la valeur manquante par une valeur moyenne, ou en utilisant une méthode d'imputation plus avancée.
3. **Traitement des valeurs aberrantes** : les valeurs aberrantes sont des valeurs qui se situent très en dehors de la plage de valeurs normales. Ces valeurs peuvent être supprimées ou remplacées par une valeur moyenne.
4. **Détection des erreurs** : les données peuvent contenir des erreurs telles que des fautes de frappe ou des incohérences dans les enregistrements. Ces erreurs peuvent être détectées et corrigées à l'aide d'algorithmes de détection d'anomalies.
5. **Normalisation des données** : la normalisation des données consiste à mettre les données sur une même échelle pour faciliter leur traitement. Par exemple, les données peuvent être normalisées en les divisant par leur plage de valeurs.

Il est important de noter que le nettoyage des données peut être un processus long et fastidieux, en particulier lorsque les données sont volumineuses et complexes. Cependant, c'est une étape essentielle pour s'assurer que les données sont de qualité suffisante pour être utilisées dans l'analyse et la modélisation en intelligence artificielle.

Il n'y a pas de réponse unique à cette question, car la façon dont vous traitez les valeurs manquantes dépend du type de données et de l'analyse que vous effectuez. Dans certains cas, remplacer les valeurs manquantes par la médiane peut être une meilleure option que de les supprimer. Voici quelques points à considérer pour prendre cette décision :

1. La fréquence des valeurs manquantes : Si le nombre de valeurs manquantes est important dans votre jeu de données, la suppression de ces valeurs peut réduire considérablement la taille de votre jeu de données. Dans ce cas, il peut être préférable de remplacer les valeurs manquantes par la médiane plutôt que de les supprimer.
2. Le type de données : Le remplacement des valeurs manquantes par la médiane est particulièrement approprié pour les données numériques, car la médiane est une mesure de tendance centrale robuste qui n'est pas affectée par les valeurs aberrantes. Pour les données catégorielles, la suppression des valeurs manquantes peut être préférable.
3. L'impact sur l'analyse : Le remplacement des valeurs manquantes par la médiane peut modifier la distribution des données, ce qui peut avoir un impact sur l'analyse. Dans certains cas, cela peut être souhaitable, mais dans d'autres cas, cela peut affecter la validité des résultats.

En résumé, remplacer les valeurs manquantes par la médiane peut être une bonne option dans certains cas, mais pas dans tous les cas. Il est important de considérer la fréquence des valeurs manquantes, le type de données et l'impact sur l'analyse avant de prendre une décision.

## 6.2. Pipeline

Le principe de pipeline en IA consiste à diviser le processus de développement d'un système d'IA en plusieurs étapes distinctes, appelées "étapes de traitement". Chaque étape est conçue pour accomplir une tâche spécifique et elle est généralement reliée à l'étape suivante dans un flux séquentiel ou "pipeline".

Le pipeline d'IA typique comprend les étapes de prétraitement, de modélisation et d'évaluation. La première étape, le prétraitement, implique la collecte, la sélection et la préparation des données à utiliser pour entraîner le modèle d'IA. Cette étape comprend souvent des tâches telles que la normalisation des données, la suppression des valeurs aberrantes et la transformation des données brutes en un format utilisable par le modèle.

La deuxième étape, la modélisation, implique la création et l'entraînement d'un modèle d'IA à l'aide des données préparées lors de la première étape. Cette étape comprend le choix de l'algorithme de modélisation, le réglage des paramètres de l'algorithme, l'entraînement du modèle et la vérification de sa précision.

La dernière étape, l'évaluation, consiste à tester le modèle sur des données qui n'ont pas été utilisées lors de la phase de formation. Cette étape permet de déterminer la précision du modèle et sa capacité à généraliser au-delà des données d'entraînement.

Le pipeline en IA est un processus itératif et continu qui peut être amélioré en ajoutant des étapes supplémentaires ou en ajustant les paramètres existants pour améliorer la performance globale du modèle. Le principe de pipeline permet de diviser le processus complexe de développement d'IA en tâches plus petites et plus gérables, ce qui facilite la gestion du projet et l'amélioration continue du modèle d'IA.



*Figure 5 : schématisation d'un pipeline*

#### *6.2.1. Détail au sein du workshop & livrable.*

## I- Hypothèses

- K-means est un algorithme de classification
- Un algorithme supervisé demande une activité humaine régulière
- Un algorithme non supervisé tourne et s'entraîne seul pendant X temps
- K-means ne va pas forcément affiner les résultats
- Le K-means est mieux que l'algorithme de classification ascendante hiérarchique
- Il existe une classification descendante hiérarchique
- Ce jeu de données pourra être mieux traité par un autre algorithme
- La métrique utilisée sera la fonction perte / coût (Loss function)

Nous allons utiliser une échelle pour prétraiter nos données

## 7. Livrable

## 8. Validation des Hypothèses



# Bibliographie

---

- AI Now 2019/2020 Report*. (s.d.). Récupéré sur AI Now Institute: <https://ainowinstitute.org/>
- Antonis Papapantoleon, P. Y. (2020, Février 14). *Detection of arbitrage opportunities in multi-asset derivatives markets*. Récupéré sur Arxiv: <https://arxiv.org/abs/2002.06227>
- BNF. (s.d.). *Les enjeux éthiques et sociaux de l'intelligence artificielle*. Récupéré sur BNF.fr: <https://www.bnf.fr/fr/agenda/les-enjeux-ethiques-et-sociaux-de-lintelligence-artificielle>
- Brownlee, J. (s.d.). *Data Preparation for machine learning*.
- CNIL. (s.d.). *Comment permettre à l'Homme de garder la main ? Rapport sur les enjeux éthiques des algorithmes et de l'intelligence artificielle*. Récupéré sur cnil.fr: <https://www.cnil.fr/fr/comment-permettre-lhomme-de-garder-la-main-rapport-sur-les-enjeux-ethiques-des-algorithmes-et-de>
- Forbes. (s.d.). *AI Ethics: What It Is And Why It Matters*. Récupéré sur Forbes: <https://www.forbes.com/sites/nishatalagala/2022/05/31/ai-ethics-what-it-is-and-why-it-matters/>
- Géron, A. (s.d.). *Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow*. O'Reilly.
- Montréal, U. d. (s.d.). *La Déclaration de Montréal pour un développement responsable de l'IA*. Récupéré sur Montréal Declaration Responsible IA: <https://www.montrealdeclaration-responsibleai.com/>