



# Les différents algorithmes de l'IA

Suite de "Les différents types d'IA"

# Rubriques



## Introduction aux algorithmes d'IA

Introduction  
L'apprentissage Machine  
L'apprentissage profond  
L'apprentissage par renforcement



Apprentissage supervisé  
VS  
Apprentissage non  
supervisé



Pour aller plus loin



## Zoom sur certains algorithmes

Classification  
Clustering  
Association – Regression  
Réduction dimensionnelle

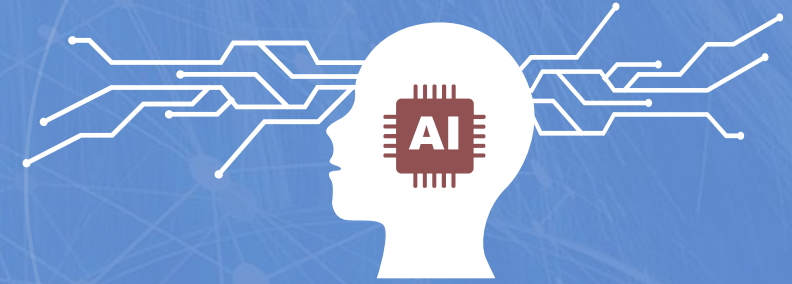


# Introduction aux algorithmes d'IA

Apprentissage machine, apprentissage profond ...



# Introduction



Ici l'objectif n'est pas de faire une classification des différents types d'IA (ou l'on retrouve des termes d'IA faible, forte, machines réactives, mémoire limitée, conscient de soi, théorie de l'esprit...) mais simplement d'avoir un aperçu des différentes approches et algorithmes utilisés en IA.

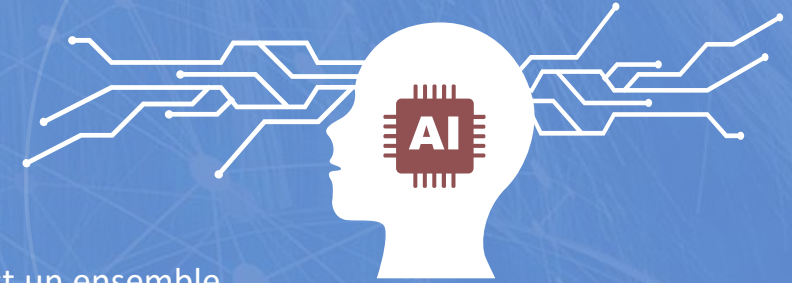
Tandis que l'apprentissage machine (**Machine Learning**) fonctionne à partir d'une base de données contrôlable, l'apprentissage profond (**Deep Learning**) a besoin d'un volume de données bien plus considérable. Le système doit disposer de plus de 100 millions d'entrées pour donner des résultats fiables.

Les **réseaux de neurones** d'apprentissage **profond** ont également évolué et contiennent désormais bien plus de couches différentes. Le **Deep Learning** de Google Photos comporte par exemple 30 couches. Une autre évolution massive est celle des **réseaux de neurones** convolutifs.

Le **machine learning**, ou apprentissage automatique, est un « concept qui tend à rendre une **machine** capable d'apprendre de ses expériences ». La **machine** récupère des quantités gigantesques d'informations, qu'elle réutilise pour s'adapter à de nouvelles situations pour les anticiper.



# Le Machine Learning – Apprentissage machine



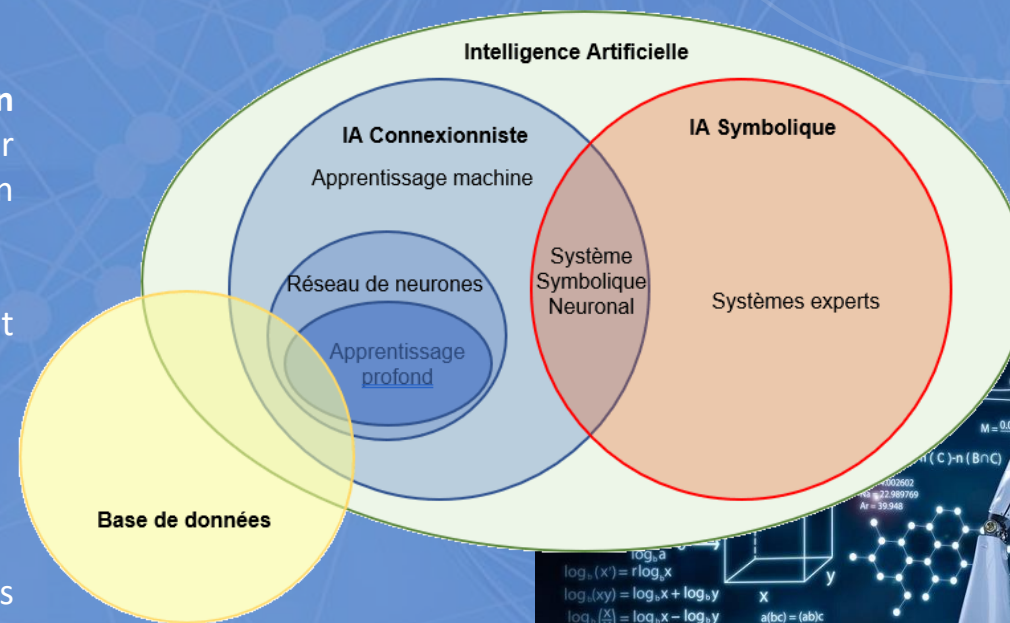
ML = Machine Learning (ou apprentissage automatique) : sous-branche de l'IA, le machine learning est un ensemble de méthodes statistiques appliquées à l'IA. Elles permettent en particulier à l'IA d'apprendre à partir de données d'exemple (apprentissage supervisé : on montre à l'IA plein d'appartements avec leur prix, puis elle peut deviner le prix d'un nouvel appartement). L'algorithme le plus simple et le plus connu du ML est la régression linéaire, qui consiste à trouver la droite approximant un nuage de points (grâce à ça, on peut facilement « classifier » des données en fonction de leur position par rapport à cette droite de régression)

Cette technologie permet de réaliser des **prédictions en consultant des données** (ou datas) en se basant sur des **statistiques**. Les premiers algorithmes ont été élaborés en 1950 : le plus célèbre d'entre eux est le Perceptron.

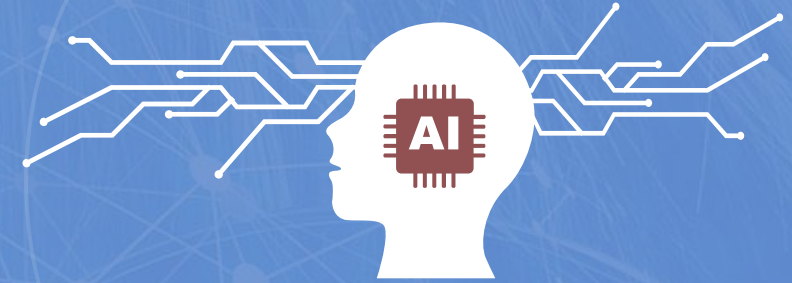
Toute démarche requérant l'usage du Machine Learning doit passer par 4 étapes :

- Le pré-traitement des données
- La modélisation
- Le déploiement
- La maintenance

Une fois les données filtrées, ordonnées et sauvegardées, elles peuvent être exposées aux algorithmes du Machine Learning.



# Le Machine Learning – Apprentissage machine



Partons d'un exemple simple : imaginons que vous vouliez créer une IA qui vous donne le prix d'un appartement à partir de sa superficie. Dans les années 1950, vous auriez fait un programme du type « si la superficie est inférieure à 20m<sup>2</sup>, le prix vaut 60 000€, si elle est entre 20m<sup>2</sup> et 30m<sup>2</sup>, le prix vaut 80 000€, etc... », ou peut-être « prix = superficie\*3 000 ».

Si vous avez un ami **statisticien**, il pourrait alors vous dire que ces approximations ne sont pas satisfaisantes, et qu'il suffirait de constater le prix de plein d'appartements dont on connaît la superficie pour estimer le prix d'un nouvel appartement de taille non-référencée ! Votre ami vient de donner naissance au machine learning (qui est donc un sous-domaine de l'intelligence artificielle).

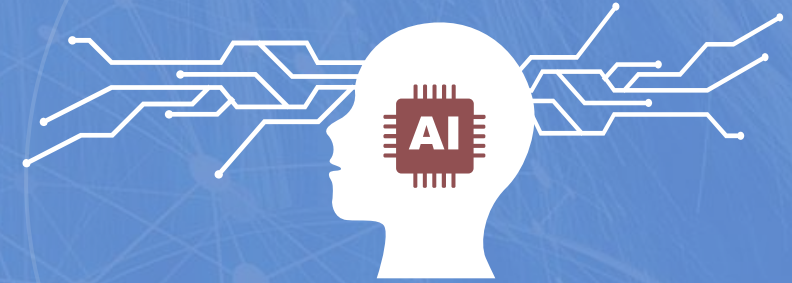
En effet, apparu dans les années 1980, le machine learning (ML) est l'application de méthodes statistiques aux algorithmes pour les rendre plus intelligents. L'enjeu du ML est bien de construire des **courbes qui approximent les données** et permettent de généraliser facilement. Il repose donc sur la capacité des algorithmes à recevoir beaucoup de données et à « apprendre » d'elles (corriger les courbes d'approximation) !

Le machine learning est un domaine large, qui comprend de très nombreux algorithmes. Parmi les plus célèbres, on retrouve :

- Les **régressions** (linéaires, multivariées, polynomiales, régularisées, logistiques...)
- L'algorithme de **Naïve Bayes** : l'algorithme donne la probabilité de la prédiction, sachant les événements antérieurs. Par exemple, quel est le prix le plus probable sachant que l'appartement fait 43.7m<sup>2</sup>.
- Le **clustering** : toujours grâce aux mathématiques, on va grouper les données en paquets de manière à ce que dans chaque paquet les données soient les plus proches possibles les unes des autres. C'est utilisé notamment pour des recommandations de films « proches » des films que vous avez déjà vus !
- Les **arbres de décision** : en répondant à un certain nombre de questions et en suivant les branches de l'arbre qui portent ces réponses, on arrive à un résultat (avec un score de probabilité)
- Ainsi que des algorithmes plus perfectionnés reposant sur plusieurs techniques de statistiques : **Random Forest** (une forêt d'arbres de décision qui votent), **Gradient Boosting**, **Support Vector Machine**...



# Le Deep Learning – Apprentissage profond



Le ML est « l'IA auto-apprenante » : l'algorithme détermine lui-même quelle est la meilleure approximation des données, et ce n'est pas un humain qui détermine « les coefficients de l'équation » (en général, c'est impossible pour un humain de trouver ces valeurs, car il y a trop de données à prendre en compte... tandis qu'une machine s'appuie sur des calculs).

En dépit de sa puissance, le ML pur a beaucoup de failles. La première est qu'un expert humain doit, au préalable, faire du tri dans les données. Par exemple, pour notre appartement, si vous pensez que l'âge du propriétaire n'a pas d'incidence sur le prix, il n'y a aucun intérêt à donner cette information à l'algorithme, car si vous lui en donnez trop, il pourrait voir des relations là où il n'y en a pas...

Ensuite, la seconde (qui découle de la première) : comment faire pour reconnaître un visage ? Vous pourriez donner à l'algorithme plein d'informations sur la personne (écart entre les yeux, hauteur du front, etc...), mais ce ne serait pas très adaptatif ni précis.

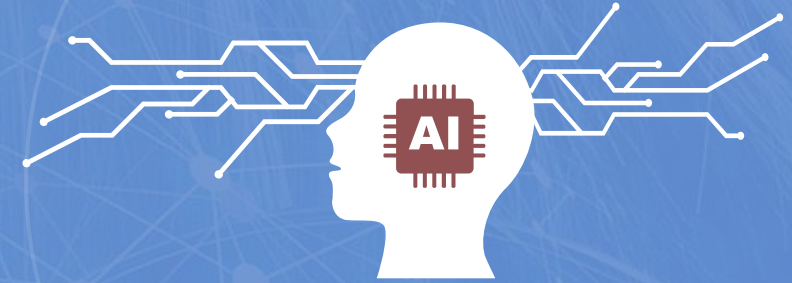
Ainsi est née l'idée du **Deep Learning** (DL – Apprentissage profond) vers 2010 : s'inspirer du fonctionnement de notre cerveau (avec des réseaux de neurones) pour pousser l'analyse plus loin et savoir extraire les données soi-même !

Le DL (qui est donc un sous-domaine du ML) repose donc sur ce qu'on appelle des réseaux de neurones artificiels (profonds), c'est-à-dire un ensemble de neurones sur plusieurs couches jusqu'à des neurones de sortie. Grâce à cette architecture, le DL est capable de reconnaître des visages, de synthétiser des textes ou encore de conduire une voiture autonome !

Vous avez du mal à voir où sont les statistiques dans tout ça ? En fait, l'algorithme va adapter les liaisons entre ses neurones (il va les renforcer ou les détruire), pour qu'en sortie on ait une bonne approximation des données d'entrée. Voir complément sur « Les réseaux de neurones »



# Le Deep Learning – Apprentissage profond

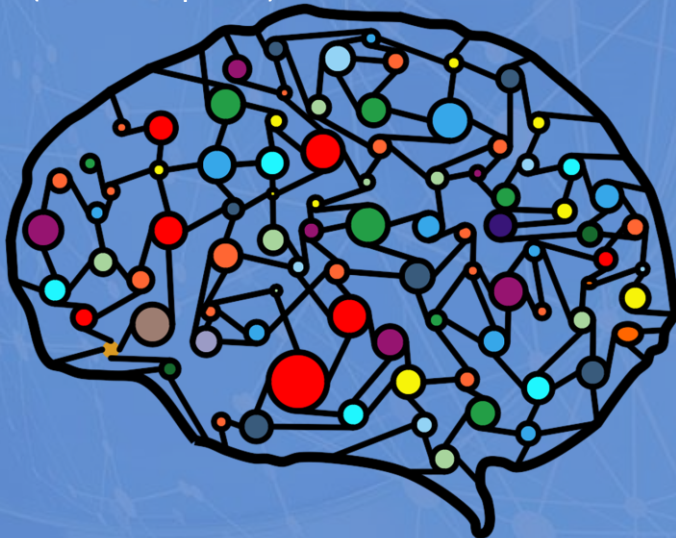


Voici quelques exemples d'algorithmes de Deep Learning :

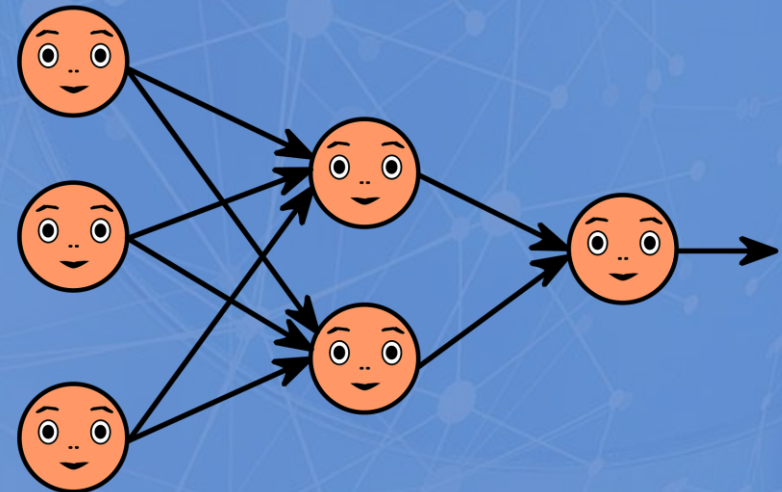
- Les réseaux de neurones artificiels (ANN) : ce sont les plus simples et sont souvent utilisés en complément car ils trient bien les informations
- Les réseaux de neurones convolutifs (CNN) : spécialisés dans le traitement de l'image, ils appliquent des filtres à des données pour en faire ressortir de nouvelles informations (par exemple, faire ressortir les contours dans une image peut aider à trouver où est le visage)
- Les réseaux de neurones récurrents (RNN) : les plus connus sont les LSTM (Long Short-Term Memory), qui ont pour faculté de retenir de l'information et de la réutiliser peu après. Ils servent pour l'analyse de texte (NLP), puisque chaque mot dépend des quelques mots précédents (pour que la grammaire soit correcte).

Ainsi que des versions plus avancées, comme les auto-encoders, les machines de Boltzmann, les self-organizing maps (SOM)...

En conclusion le Deep Learning permet de se passer d'un expert humain pour faire le tri dans les données, puisque l'algorithme trouvera de lui-même ses corrélations. Pour reprendre l'exemple de la reconnaissance faciale, l'algorithme de DL déterminera de lui-même s'il doit tenir compte de l'écart entre les yeux (entre les pixels) ou si cette information n'est pas assez déterminante comparée à d'autres (et c'est effectivement le cas).

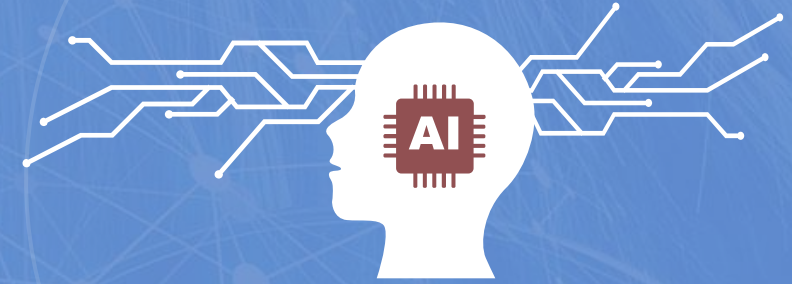


Images Pixabay





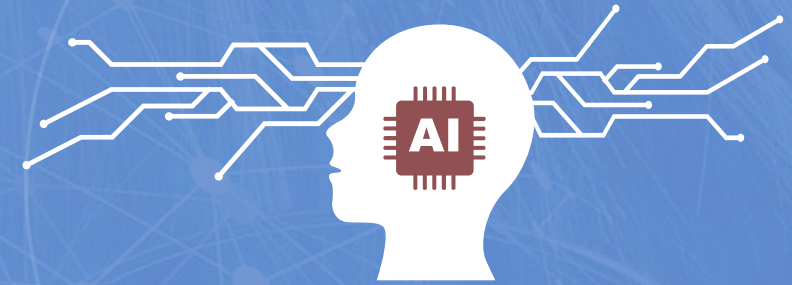
# Le Machine Learning



Une représentation non exhaustive des applications et des domaines du machine learning

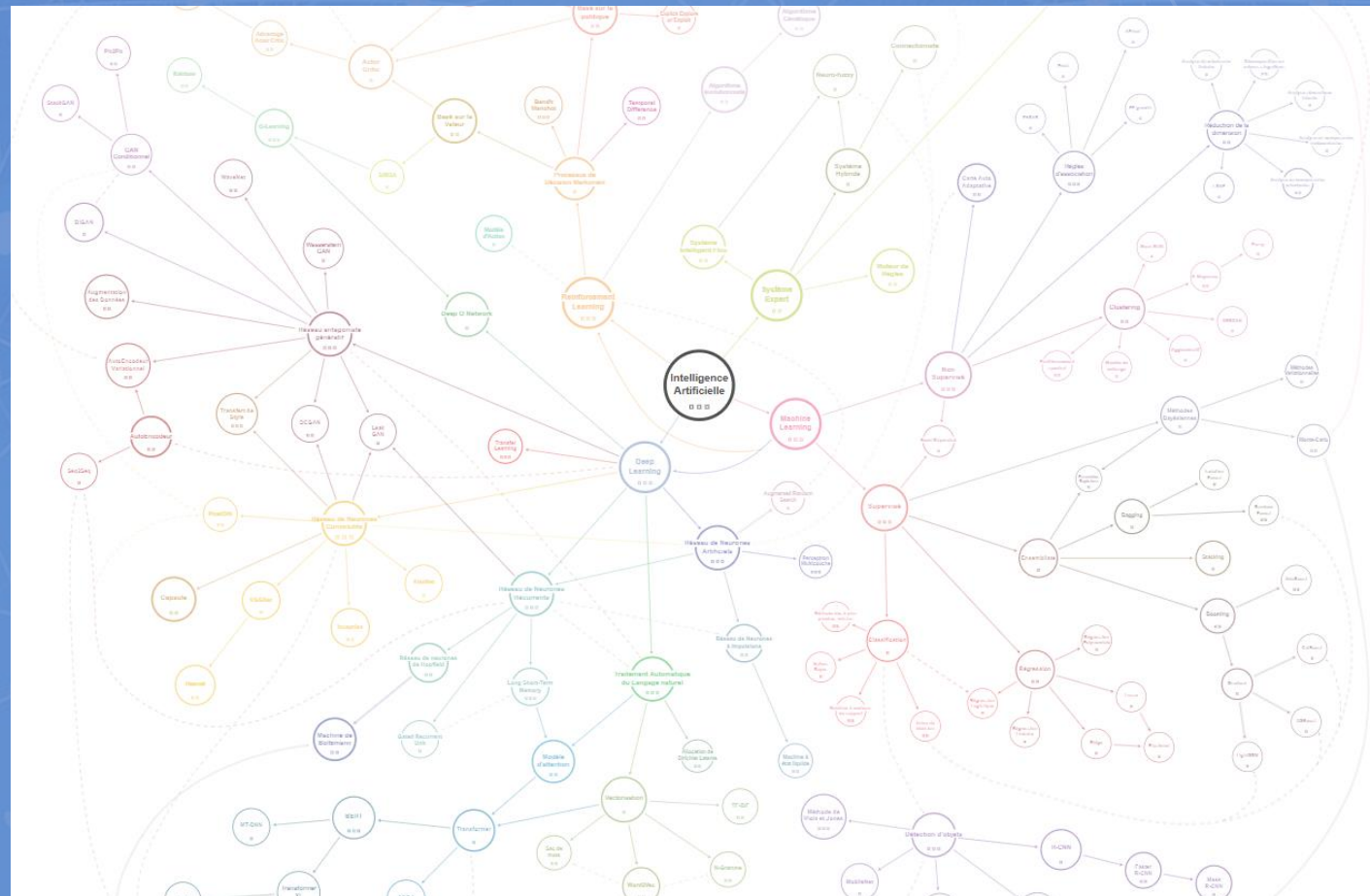


# Les domaines de l'Intelligence Artificielle

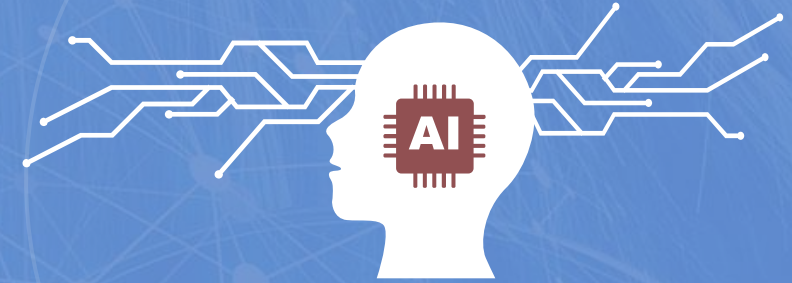


Vous retrouverez sur ce site tous les algorithmes et domaines de l'IA, avec à chaque fois une explication détaillée

– [Lien page Knowmap.org](https://www.knowmap.org)



# Apprentissage par renforcement



Et l'apprentissage par renforcement dans tout ça ?

C'est une méthode d'apprentissage dite « par renforcement » qui est utilisée sur certains algorithmes pour permettre, notamment, à une voiture d'apprendre à conduire toute seule par la pratique ou à un robot de sortir d'un labyrinthe. C'est ce type d'apprentissage qui a aussi permis à Google DeepMind de gagner aux échecs.

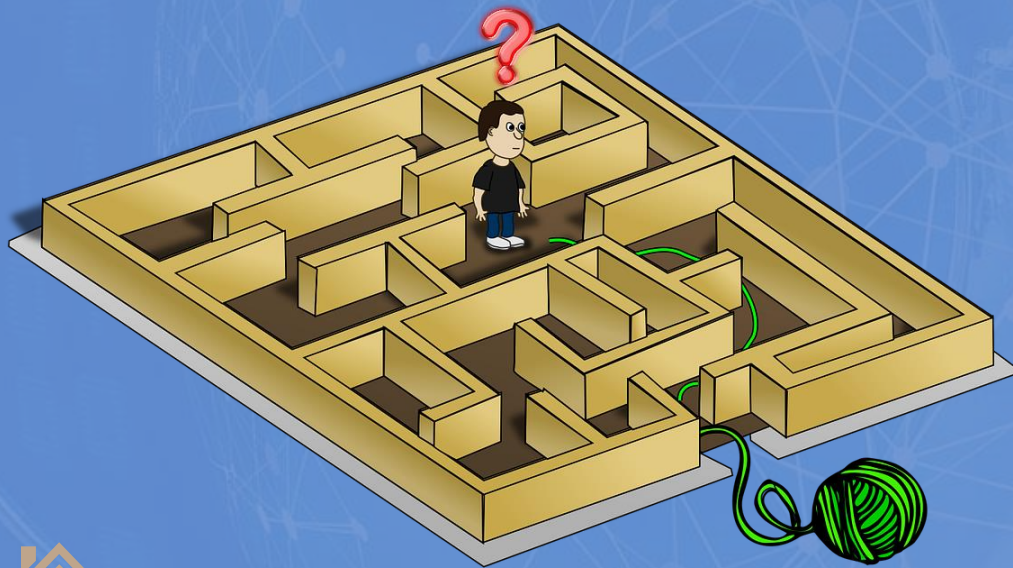
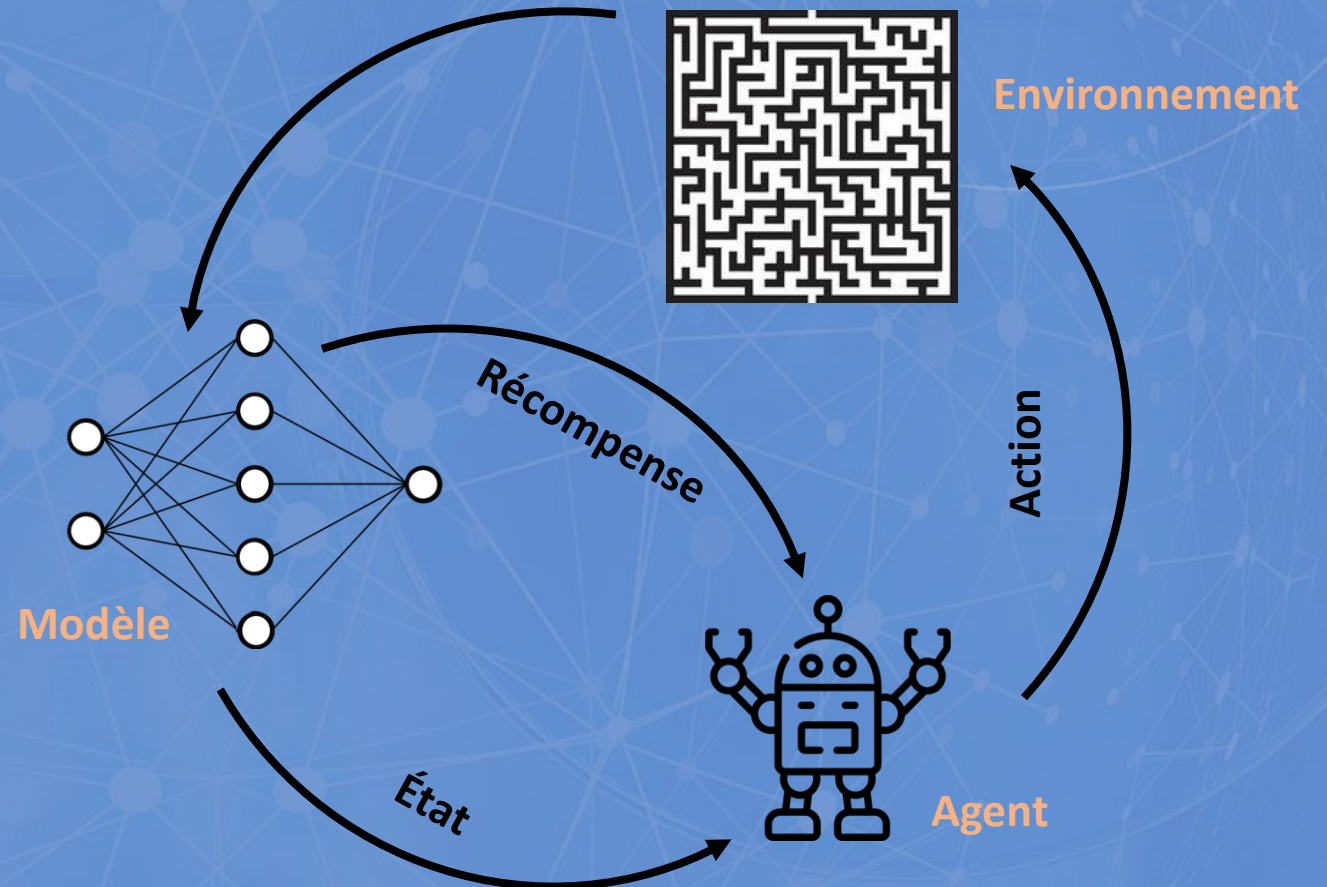


Image Pixabay





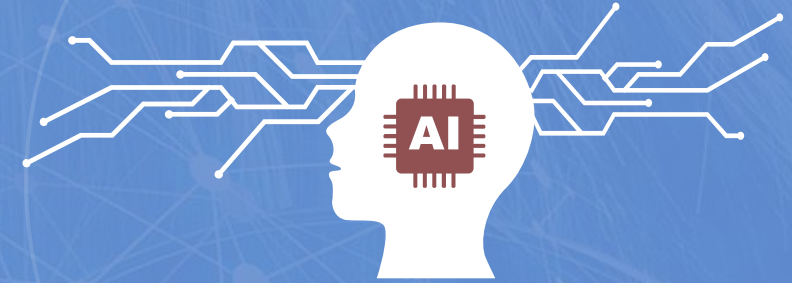
**Apprentissage  
supervisé**

**VS**

**Apprentissage  
non supervisé**



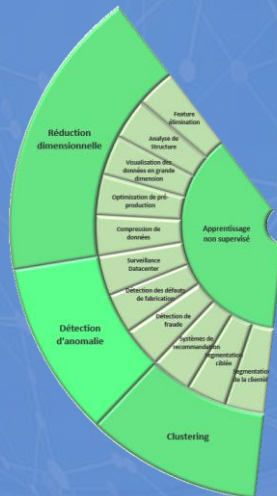
# Apprentissage supervisé et non supervisé



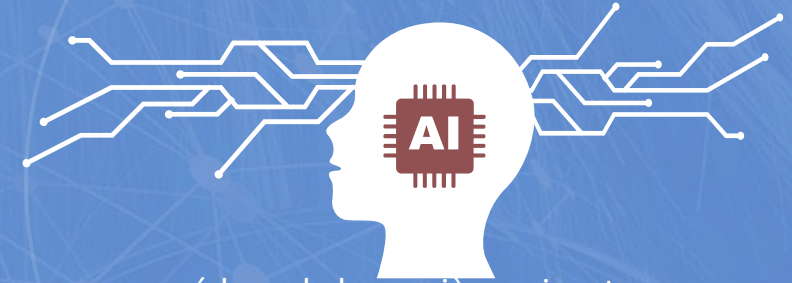
Pour illustrer la différence entre ces 2 modes d'apprentissage, prenons les 2 domaines de la classification (apprentissage supervisé) et du regroupement (clustering - apprentissage non supervisé) souvent confondus et quelques peu similaires en termes de finalité mais d'approche différente.

Le **clustering** est une méthode de regroupement d'objets de telle sorte que les objets ayant des caractéristiques similaires se rejoignent et les objets ayant des caractéristiques différentes se séparent.

La **classification** est un processus de catégorisation qui utilise un ensemble de données d'apprentissage pour reconnaître, différencier et comprendre des objets, on dit alors que les données d'entrées sont labellisées selon leurs paramètres descripteurs. Les usages de cette méthode sont par exemple liés à la détection de spams, ou à l'analyse du risque dans le domaine de la santé. Pour le premier, après avoir scanné le texte d'un email, et tagger certains mots et phrases, la « signature » du message peut être injectée dans un algorithme de classification pour déterminer si oui ou non il s'agit d'un spam. Dans le cas du second, les statistiques vitales d'un patient, son historique de santé, ces niveaux d'activités et les données démographiques peuvent être croisées pour attribuer une note (un niveau de risque) et évaluer la probabilité d'une maladie.

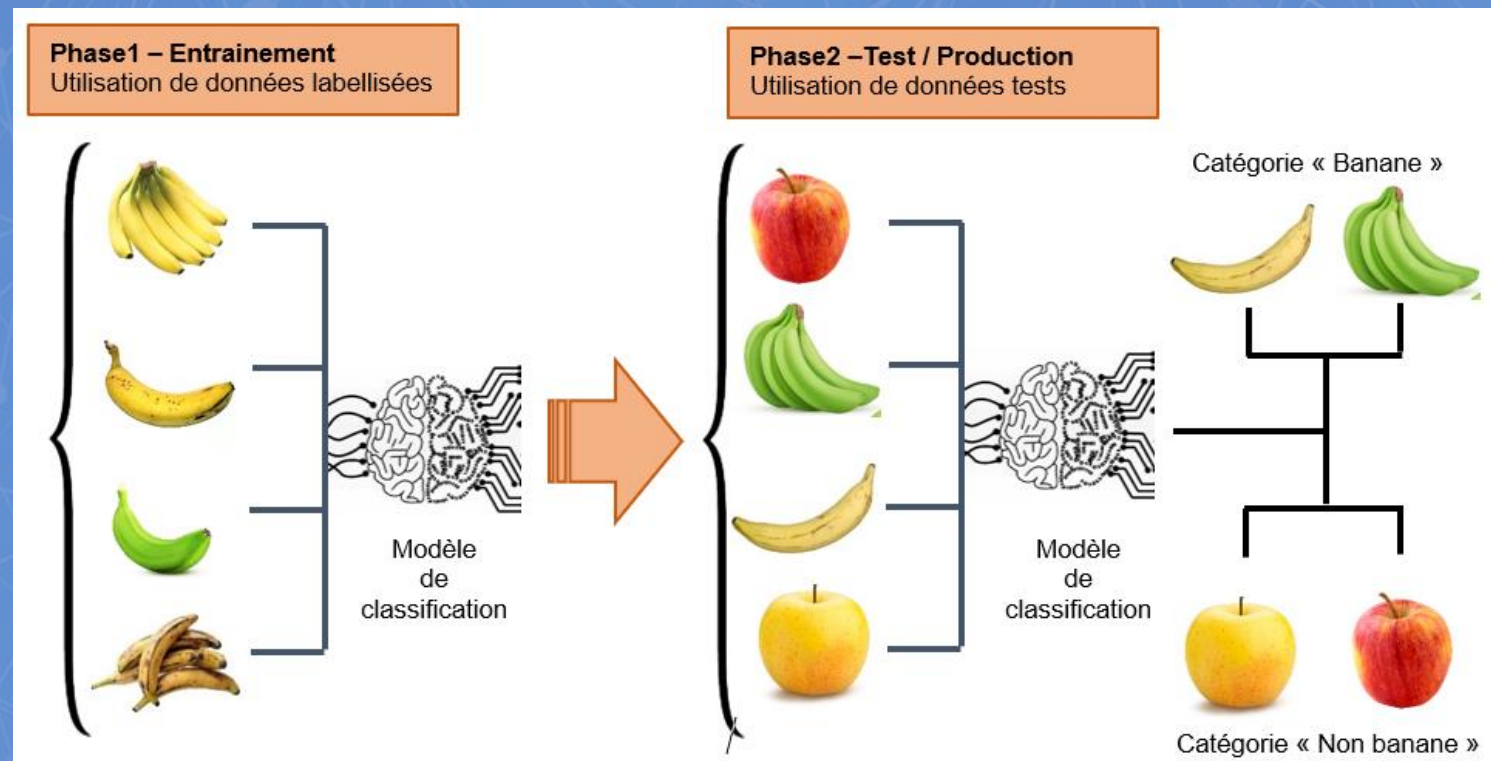


# Apprentissage supervisé et non supervisé

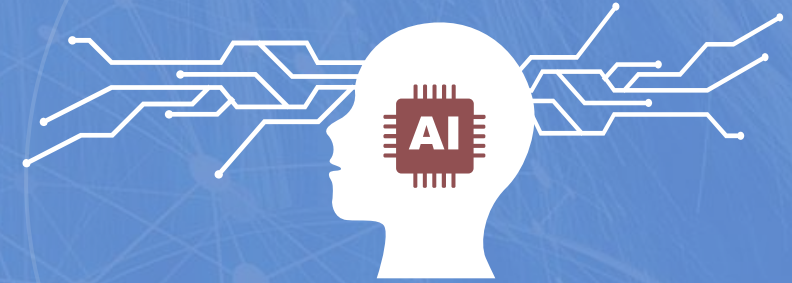


Prenons le cas de nos bananes sur le schéma ci-dessous. Si on désire reconnaître la forme d'une banane, nous procédons de la manière suivante en apprentissage supervisé.

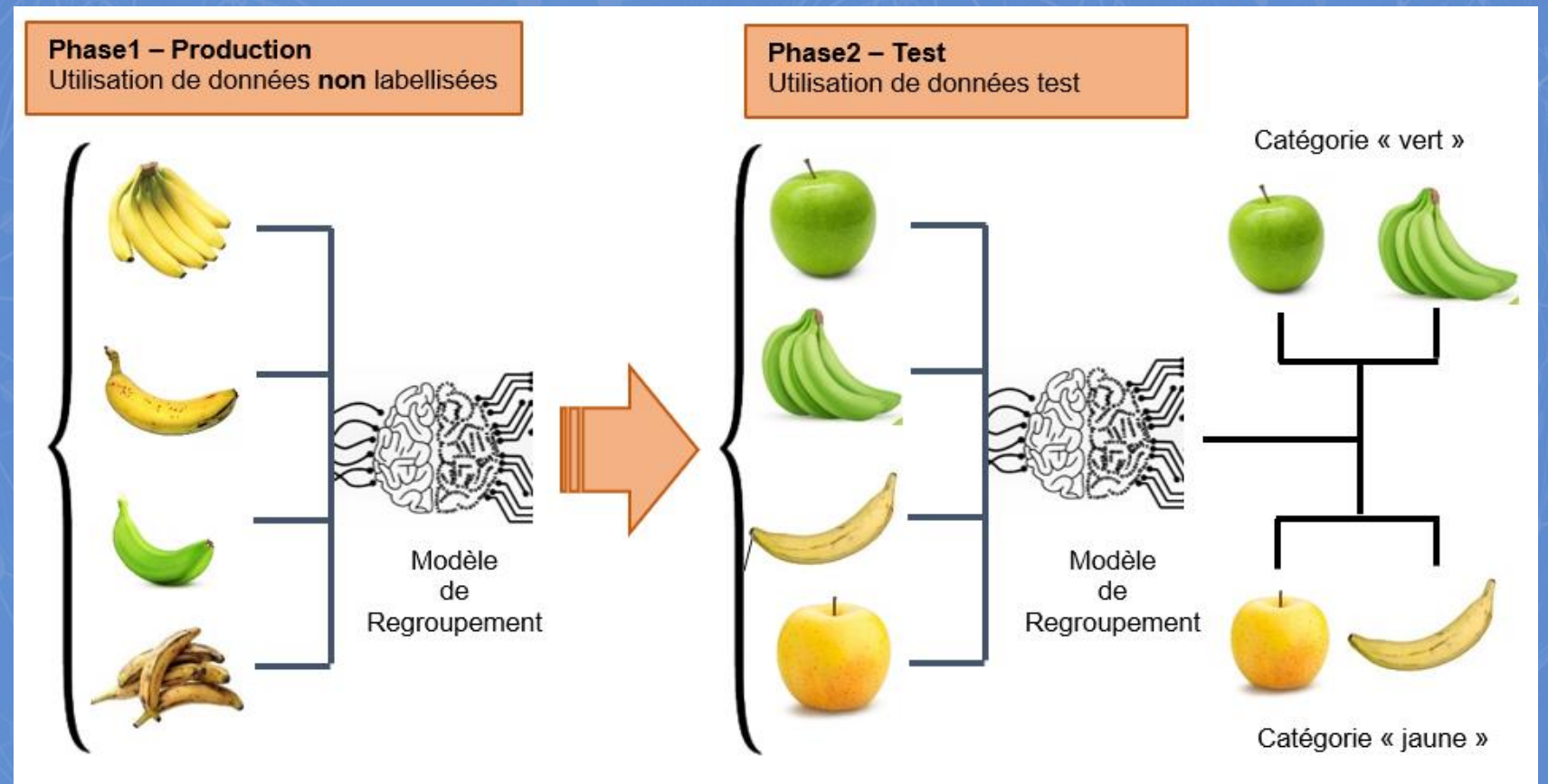
Nous disposons d'un échantillon de données labellisées « banane ». Nous prélevons 75 à 80% de l'échantillon pour entraîner notre modèle. Par la suite, les 25 ou 20% restants permettent de vérifier et valider la « robustesse » de notre modèle de classification. En y associant des pommes le modèle doit être capable de distinguer ce qui est une banane ou pas.



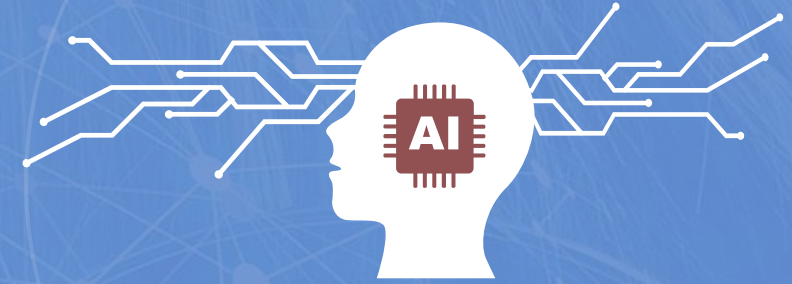
# Apprentissage supervisé et non supervisé



Dans le cas du regroupement pas de label, notre algorithme sera capable de trier nos données sans savoir s'il s'agit de pommes ou de bananes. Par contre nous pourrions obtenir un regroupement par couleur et non par forme.



# Apprentissage supervisé et non supervisé



Dans le domaine industriel ou scientifique, l'usage de la classification n'est pas limité à de la reconnaissance d'objets, elle permet par exemple d'optimiser des processus de fabrication à partir des données prélevées sur la machine de production (ici des mesures de température et de pression prélevées par des capteurs sur une machine d'injection) et des caractéristiques produits (dimension, couleur, résistance mécanique...).

La base de données se présente alors sous forme d'un fichier tableur csv (figure ci-contre). Elle se justifie lorsque que le nombre de paramètres et leur valeurs associées fournissent une base de données très importante. Dans des cas plus restreints la méthode des Plans d'Expérience permet d'obtenir des résultats probants.

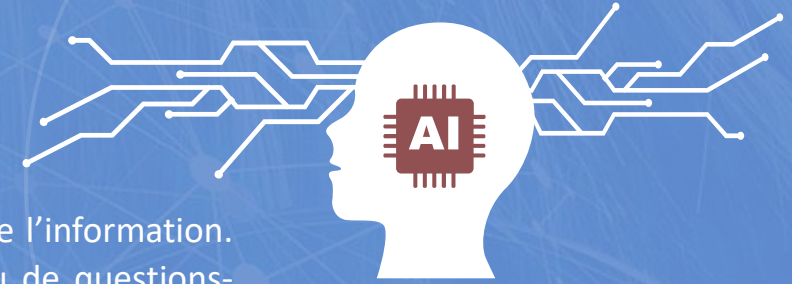
Les algorithmes de classification s'appuient sur la méthode ACP (Analyse en Composantes Principales) qui permet d'explorer des jeux de données multidimensionnels constitués de variables quantitatives) à laquelle sont adossés des algorithmes de tri (arbre de décision) et de régression (linéaire ou non) pour l'interprétation des résultats.

Humidity	Temperature	M1RMP5	M1Z1T	M1Z2T	M1MAmp	M1MPressur	M1MTemp	M1ExitTemp	M2RMP1
17.24	23.53	1207.77	72.3	72.7	49.37	431.12	77.3	75.1	12.59
17.24	23.53	1134.89	72.3	72.7	48.84	427.93	77.3	75.1	12.59
17.24	23.53	1246.27	72.3	72.7	48.97	433.71	77.3	75.1	12.59
17.24	23.53	1246.37	72.3	72.7	48.7	430.82	77.4	75.1	12.59
17.24	23.53	1247.26	72.3	72.7	49.1	439.28	77.4	75.1	12.59
17.24	23.53	1217.6	72.3	72.7	49.1	428.76	77.4	75.1	12.59
17.24	23.53	1222.42	72.3	72.7	49.1	421.37	77.4	75.1	12.59
17.24	23.53	1034.48	72.3	72.7	48.7	430.59	77.4	75.1	12.59
17.24	23.53	1213.23	72.3	72.7	48.43	434.32	77.5	75.1	12.59
17.24	23.53	1224.42	72.3	72.6	48.3	433.69	77.5	75.1	12.59
17.24	23.53	1246.69	72.3	72.6	48.57	420.96	77.5	75.1	12.59
17.24	23.53	1236.19	72.3	72.6	48.7	434.7	77.5	75.1	12.59
17.24	23.53	1241.22	72.4	72.6	47.62	423.68	77.5	75.1	12.59
17.24	23.53	1231.75	72.4	72.6	47.76	432.14	77.6	75.1	12.59
17.24	23.53	1242.61	72.4	72.6	47.89	431.37	77.6	75.1	12.59
17.24	23.53	1247.13	72.4	72.6	48.16	436.75	77.6	75.1	12.59
17.24	23.53	1242.94	72.4	72.6	48.7	420.54	77.6	75.1	12.59
17.24	23.53	1234.27	72.4	72.6	47.76	425.71	77.6	75.1	12.59
17.24	23.53	1140.22	72.4	72.6	47.36	425.94	77.7	75.1	12.59
17.24	23.53	1244.96	72.4	72.6	47.22	431.37	77.7	75.1	12.59
17.24	23.53	1245.84	72.4	72.5	48.03	427.87	77.7	75.1	12.59
17.24	23.53	1247.38	72.4	72.5	48.16	435.19	77.7	75.1	12.59
17.24	23.53	1247.74	72.4	72.5	48.43	427.19	77.7	75.1	12.59
17.24	23.53	1246.36	72.4	72.5	47.22	417.46	77.8	75.1	12.59
17.24	23.53	1171.1	72.4	72.5	46.95	428.26	77.8	75.1	12.59
17.24	23.53	1249.77	72.4	72.5	47.49	430.57	77.8	75.1	12.59
17.24	23.53	1245.59	72.4	72.5	47.76	417.46	77.8	75.1	12.59
17.24	23.53	1242.46	72.4	72.5	48.03	431.73	77.8	75.1	12.59
17.24	23.53	1237.5	72.4	72.5	47.36	421.98	77.8	75.1	12.59
17.24	23.53	1237.31	72.4	72.4	47.22	429.03	77.9	75.1	12.59
17.24	23.53	1237.76	72.4	72.4	47.36	428.39	77.9	75.1	12.59
17.24	23.53	1247.49	72.4	72.4	47.62	434.66	77.9	75.1	12.59
17.24	23.53	1241.72	72.4	72.4	48.03	418.96	77.9	75.1	12.59

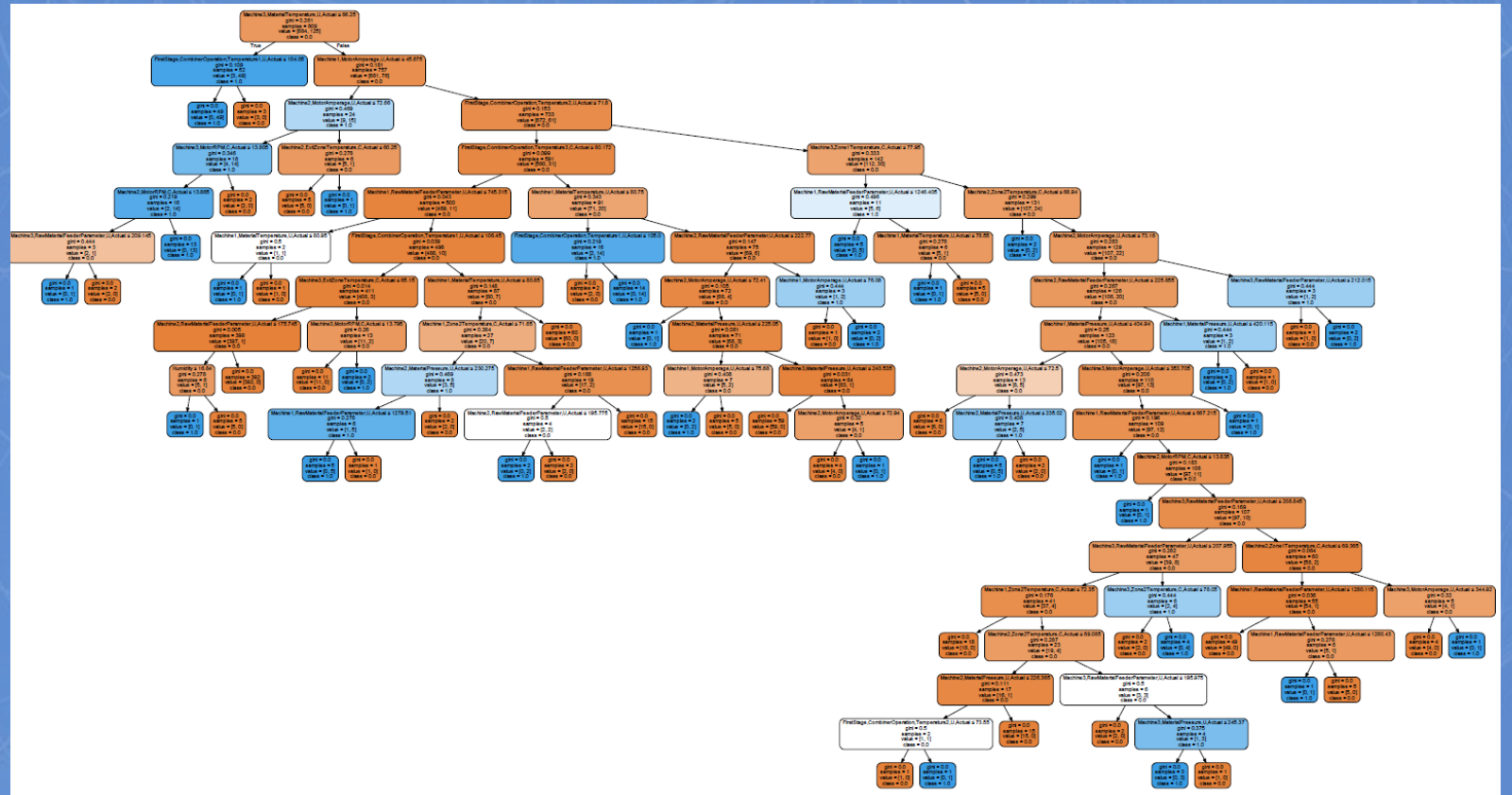




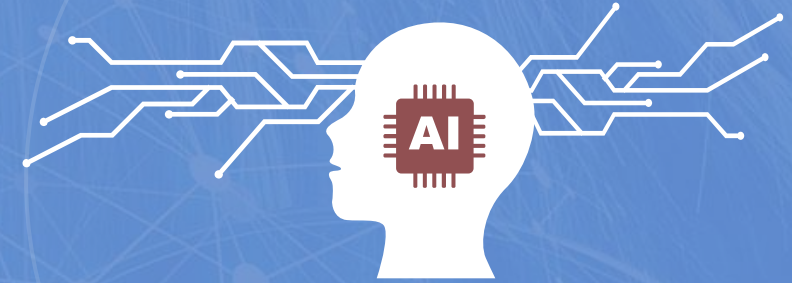
# Apprentissage supervisé et non supervisé



Le tri par arbre de décision s'appuie sur un organigramme, qui fournit une séquence hiérarchisée de l'information. Des tests sont effectués sur plusieurs paramètres d'entités classées. Ces tests peuvent être un jeu de questions-réponses (oui-non) ou inclure un ensemble plus étendu de variables distinctes. A chaque étape d'un arbre de décision, ils sont appliqués aux données pour affiner la classification et ce, jusqu'à la racine de l'arbre, les entités étant enfin séparées dans différentes classes.



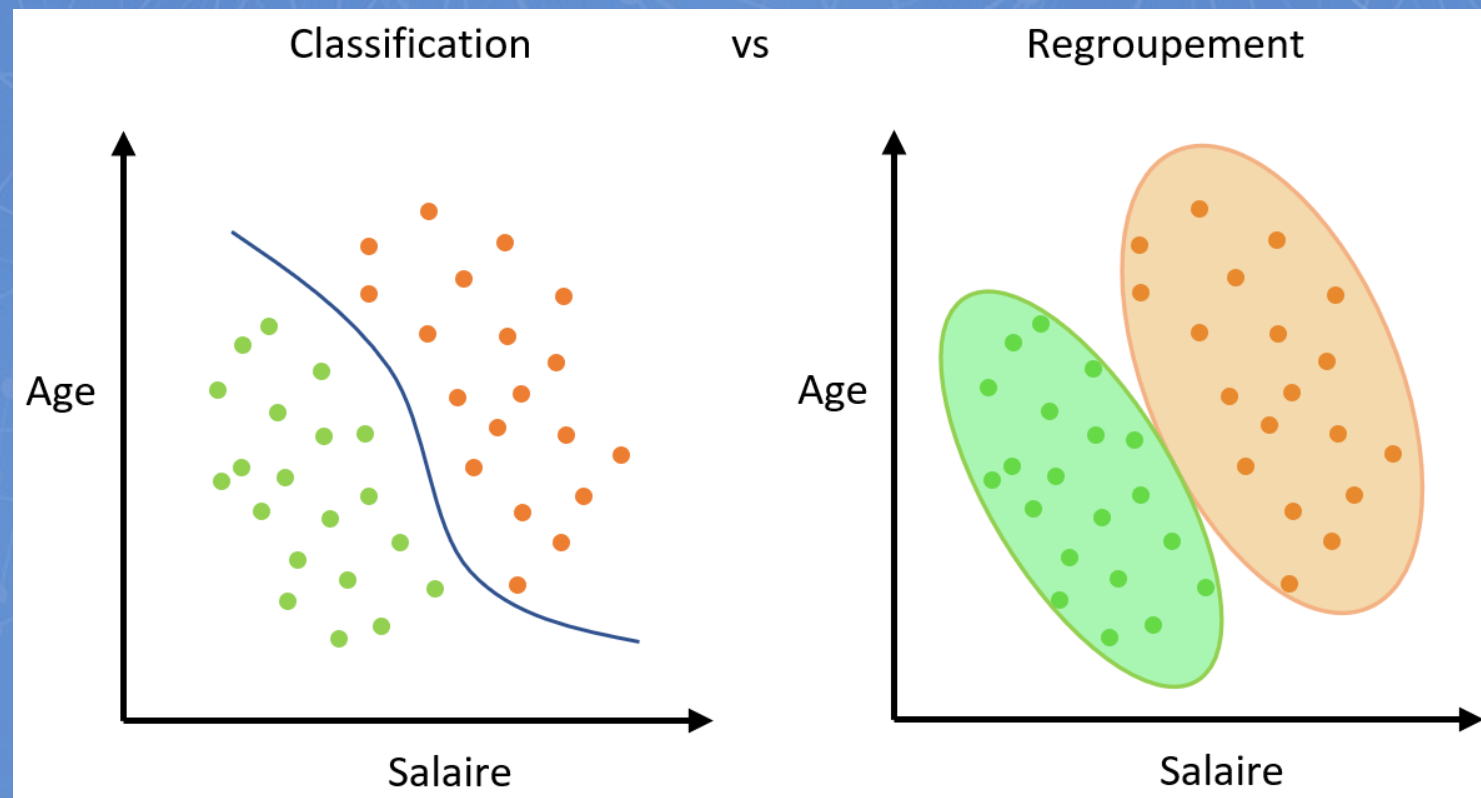
# Apprentissage supervisé et non supervisé



## En résumé :

Le choix de l'une ou l'autre technique va dépendre du type de données à disposition, si vous voulez utiliser la classification il vous faudra trier vos données et créer les labels, ce qui représente un travail substantiel. Deuxième critère, le type de sortie, dans le cas du regroupement, difficile d'avoir un regard objectif sur les résultats obtenus, notamment dans le cas de biais.

On a vu aussi que les algorithmes d'IA sont souvent associés à des méthodes statistiques et d'analyse de données classique et que plusieurs algorithmes peuvent être associés pour peaufiner les résultats.

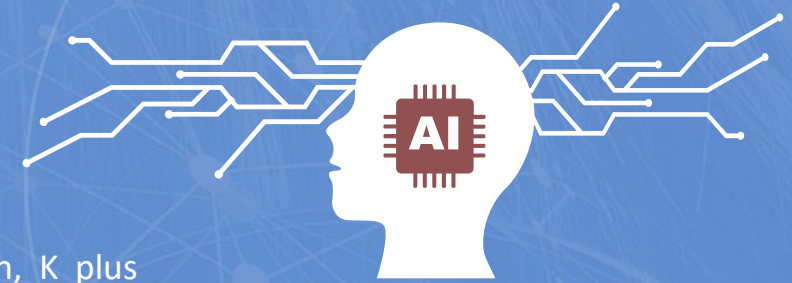




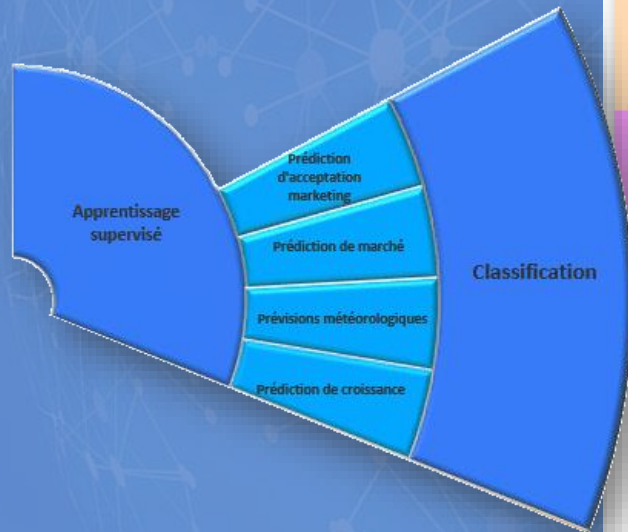
# Zoom sur certains algorithmes



# Apprentissage supervisé - Classification



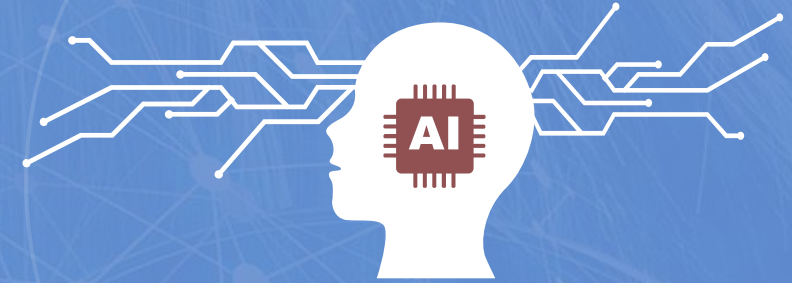
Un exemple pratique de classification utilisant différents algorithmes de classification (Perceptron, K plus proches voisins, réseau de neurones, arbre de décision...).



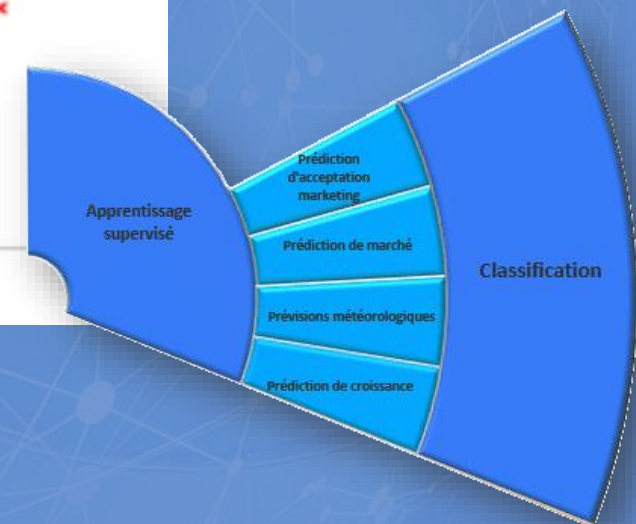
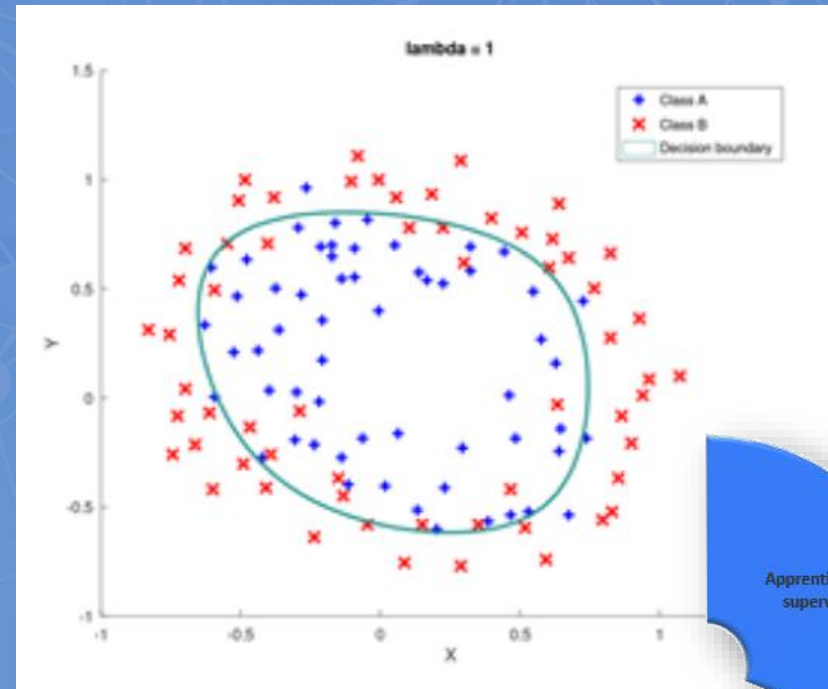
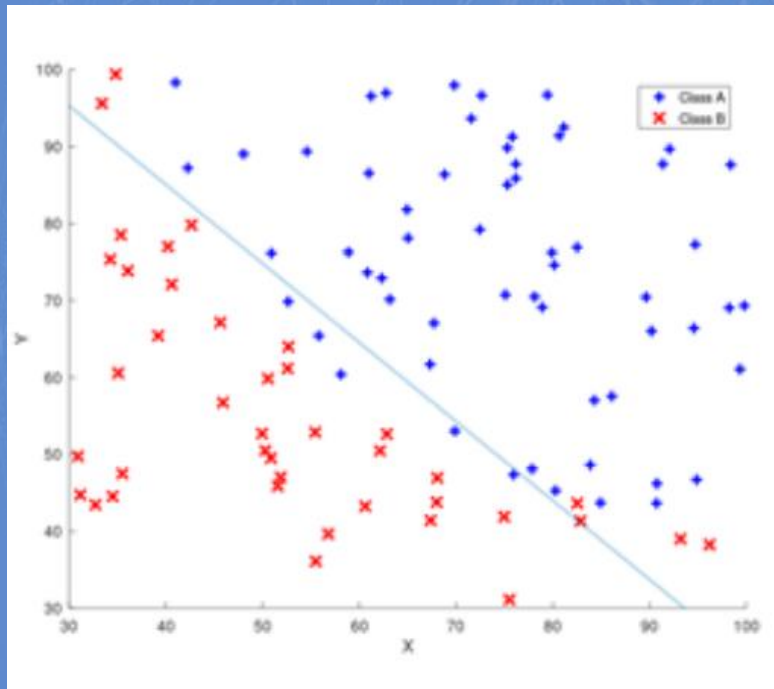
A screenshot of a web-based classification tool interface. On the left is a 2D scatter plot with orange and purple data points and a diagonal decision boundary. The plot is divided into an orange upper region and a purple lower region. On the right is a control panel with buttons for 'Upload Data', 'Save Data', and 'Clear all'. Below these are icons for 'K Nearest Neighbors', 'Perceptron', 'Support Vector Machine', 'Artificial Neural Network', and 'Decision Tree'. The 'Artificial Neural Network' icon is highlighted with a black border. Underneath is a 'Parameters' section with input fields for 'Learning rate: 0.02', 'Max Epochs: 500', and 'Max error %: 2'. Below the parameters is a neural network diagram with two layers of nodes, labeled 'X' and 'Y', and a 'Train' button at the bottom.



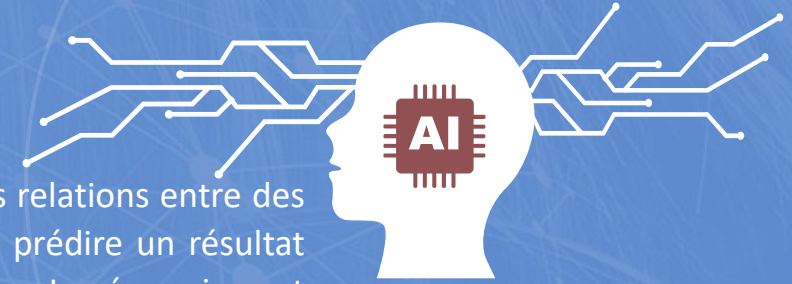
# Apprentissage supervisé - Classification



Dans la méthode de classification, nous pouvons avoir une classification binaire linéaire (figure de gauche) ou non linéaire (figure de droite).  
On peut aussi observer la présence de points mal classés, on parle alors de biais.

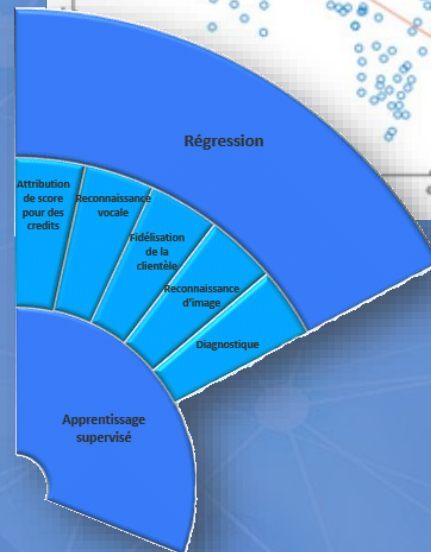
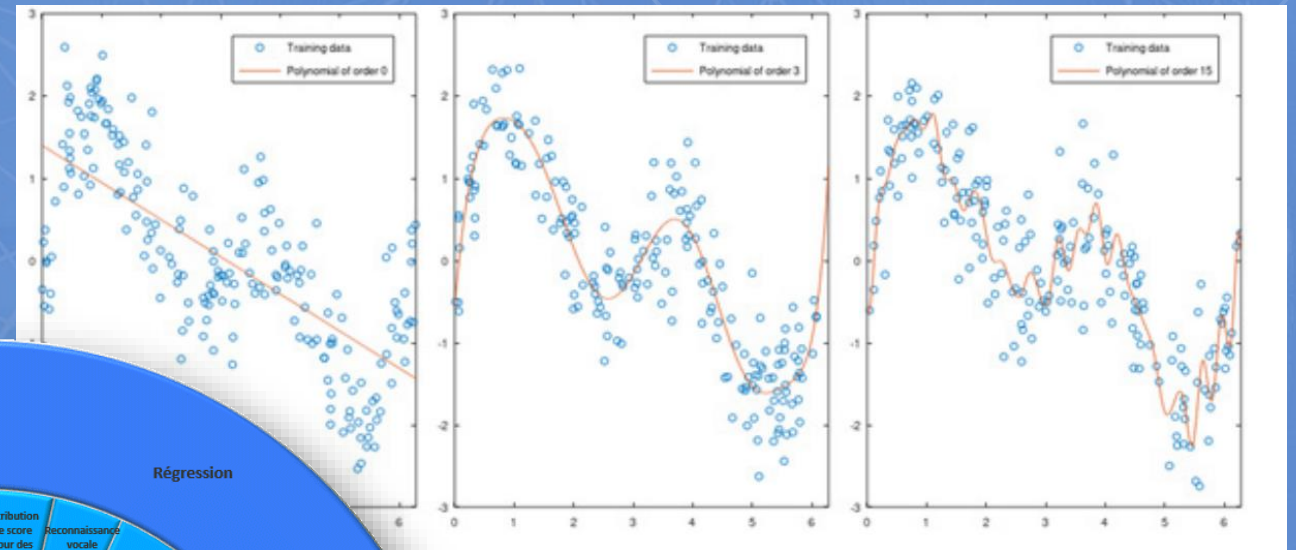


# Apprentissage supervisé - Regression

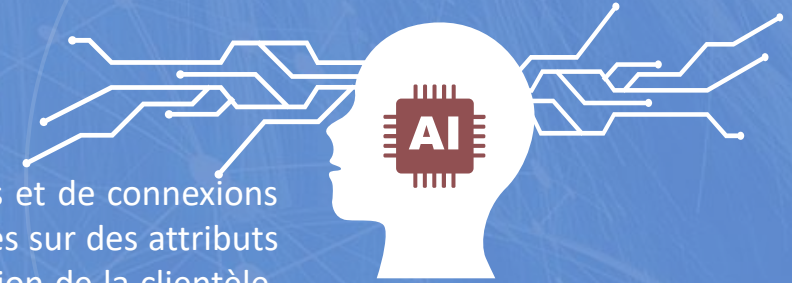


Une tâche de régression est une tâche qui analyse des données continues dans le but de trouver des relations entre des variables (généralement entre une variable dépendante et plusieurs variables indépendantes) pour prédire un résultat théorique pour lequel il n'y a pas de mesures de données disponibles, telles que des prévisions futures. La régression est principalement utilisée dans les modèles de prédiction et de prévision ; un algorithme de prédiction apprend et crée ses modèles sur les caractéristiques des états actuels ou historiques des variables pour créer et prédire une sortie de valeur continue. Cela peut conduire à une relation de régression linéaire simple ou à une relation logarithmique, exponentielle ou polynomiale plus complexe de différents degrés.

Le polynôme degré 0 (fonction linéaire – figure de gauche) montre un « sous-entraînement » (underfitting) du modèle de prédiction aux données d'apprentissage, manquant ainsi la plupart des points de données. le polynôme du 15e degré de l'ensemble d'apprentissage montre un « sur-entraînement » (overfitting) du modèle de prédiction qui s'aligne trop parfaitement sur les données d'apprentissage, mais qui fonctionnerait mal sur des données en dehors de l'ensemble d'apprentissage, La prédiction la plus juste correspond à celle du milieu,

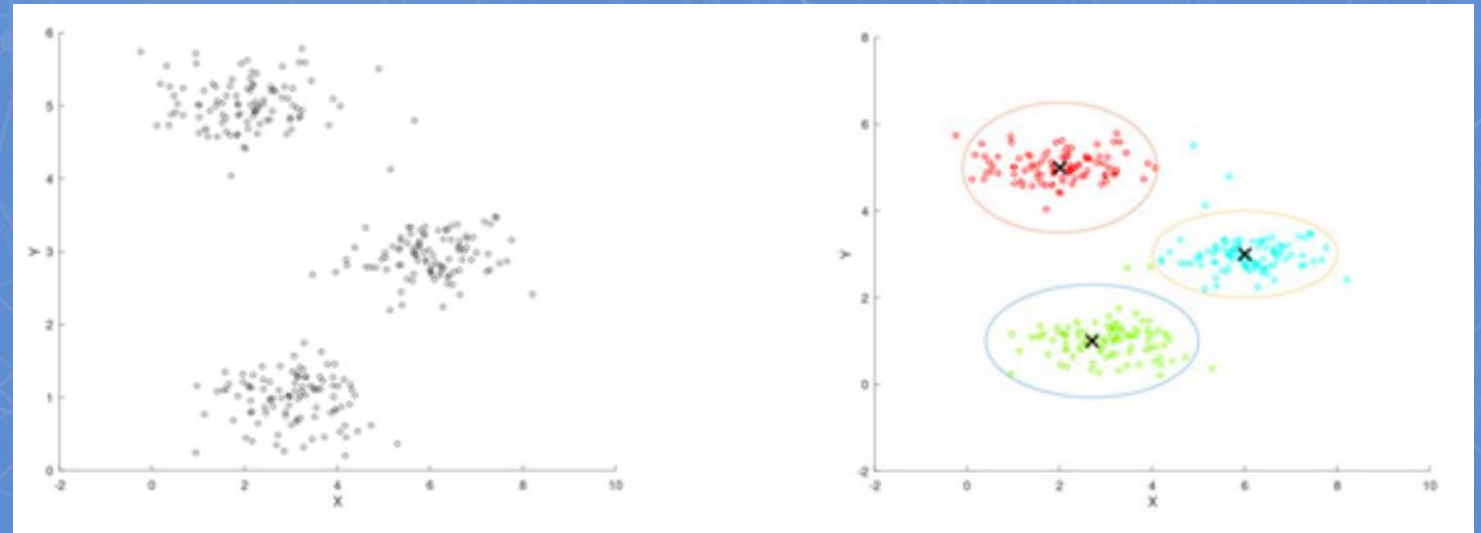


# Apprentissage supervisé - Clustering



Le clustering (regroupement), dans l'apprentissage non supervisé, implique la recherche de modèles et de connexions communes dans des ensembles de données non étiquetés et la création de groupes de données basés sur des attributs communs. Le clustering peut être vu dans le secteur du marketing où il est utilisé pour la segmentation de la clientèle, c'est-à-dire pour créer des « groupes » de clients qui partagent certains attributs communs.

En général, les algorithmes de clustering examinent un nombre défini de caractéristiques de données et mappent chaque entité de données à un point correspondant dans un graphique dimensionnel. Les algorithmes cherchent ensuite à regrouper les éléments en fonction de leur proximité relative les uns par rapport aux autres dans le graphique.



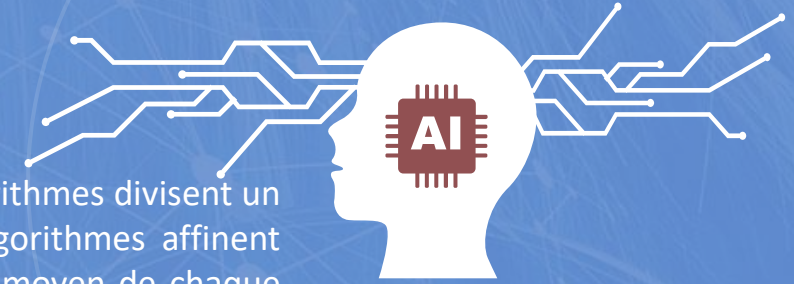
*Données brutes*

*Données regroupées*

Par exemple, une chaîne TV qui veut déterminer la répartition démographique des téléspectateurs par réseaux peut le faire en créant des clusters à partir des données disponibles sur les abonnés et de ce qu'ils regardent. Une chaîne de restaurants peut quant à elle regrouper sa clientèle en fonction des menus choisis par emplacement géographique, puis modifier ses menus en conséquence.

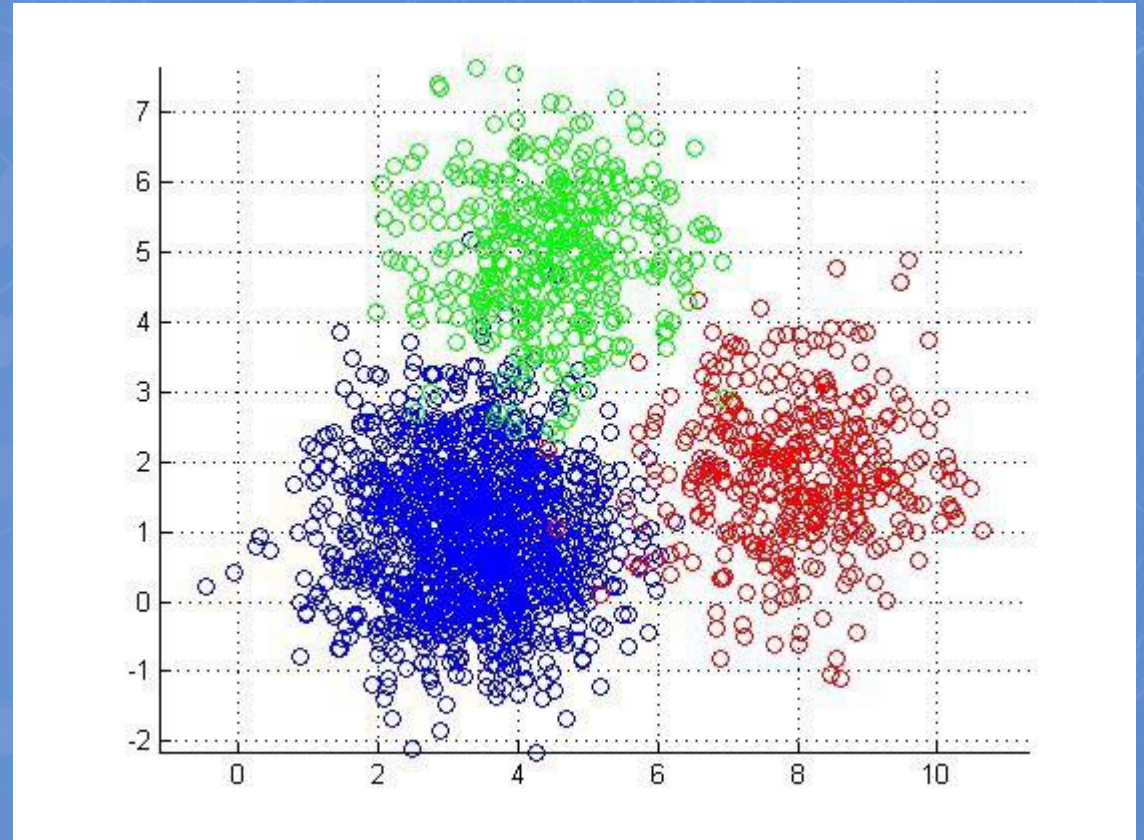


# Apprentissage supervisé - Clustering



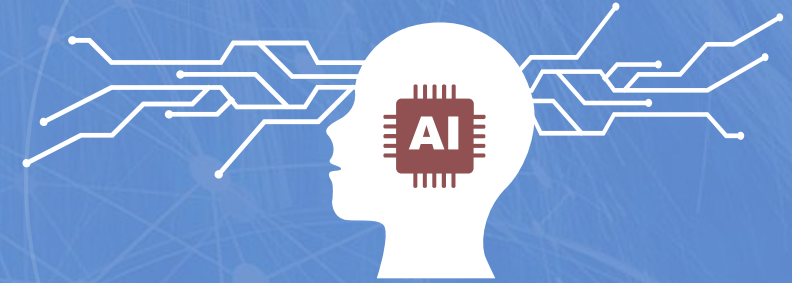
Un type couramment utilisé est l'algorithme de clustering k-moyennes (k-means en anglais). Ces algorithmes divisent un ensemble d'entités de données en groupes,  $k$  représentant le nombre de groupes créés. Les algorithmes affinent l'affectation d'entités à différents clusters en calculant itérativement le point médian ou centroïde moyen de chaque cluster (le point cité précédemment). Les centroïdes deviennent les points focaux des itérations, qui affinent leurs emplacements dans le tracé et réassignent les entités de données pour les adapter aux nouveaux emplacements. Un algorithme se répète jusqu'à ce que les regroupements soient optimisés et que les centroïdes ne bougent plus.

Source : [ledatascientist](#) – [Cliquer pour accéder](#)





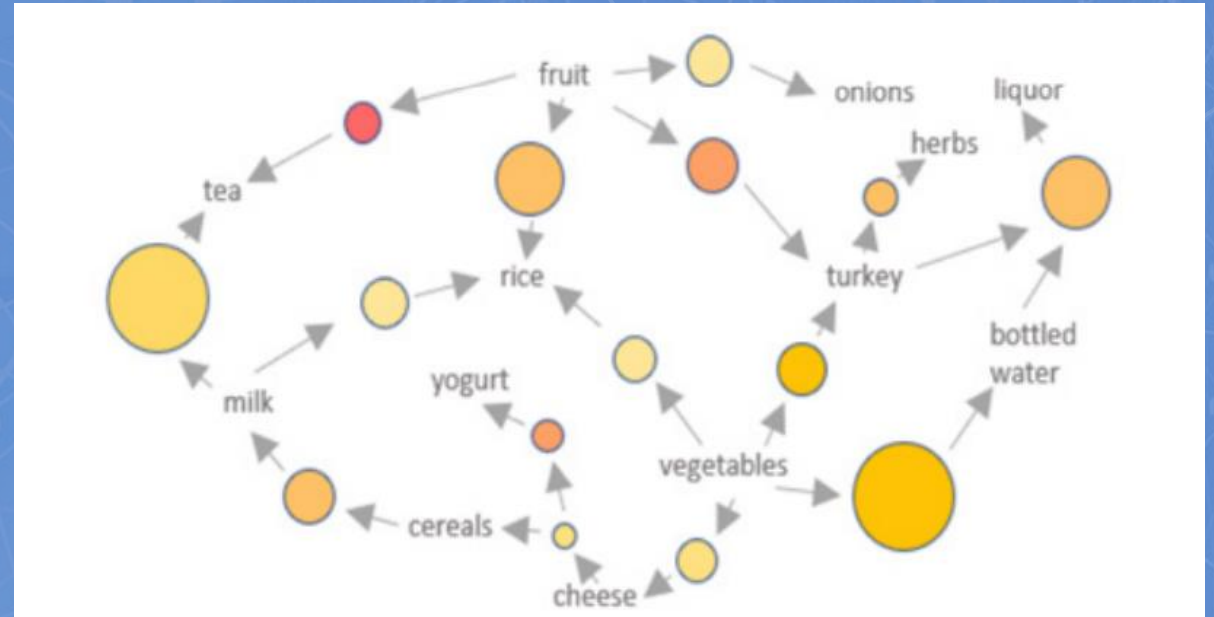
# Apprentissage supervisé - Association



La détection des règles d'association dans les données implique d'essayer de trouver et de décrire des relations auparavant inconnues et cachées dans des ensembles de données. Dans les données de transaction, un algorithme peut analyser les étapes et déterminer comment elles sont liées les unes aux autres. Il peut découvrir plus souvent quelles étapes précèdent ou succèdent à d'autres étapes et supposer des règles cachées de causalité.

*Visualisation des règles d'association à partir des transactions d'épicerie*

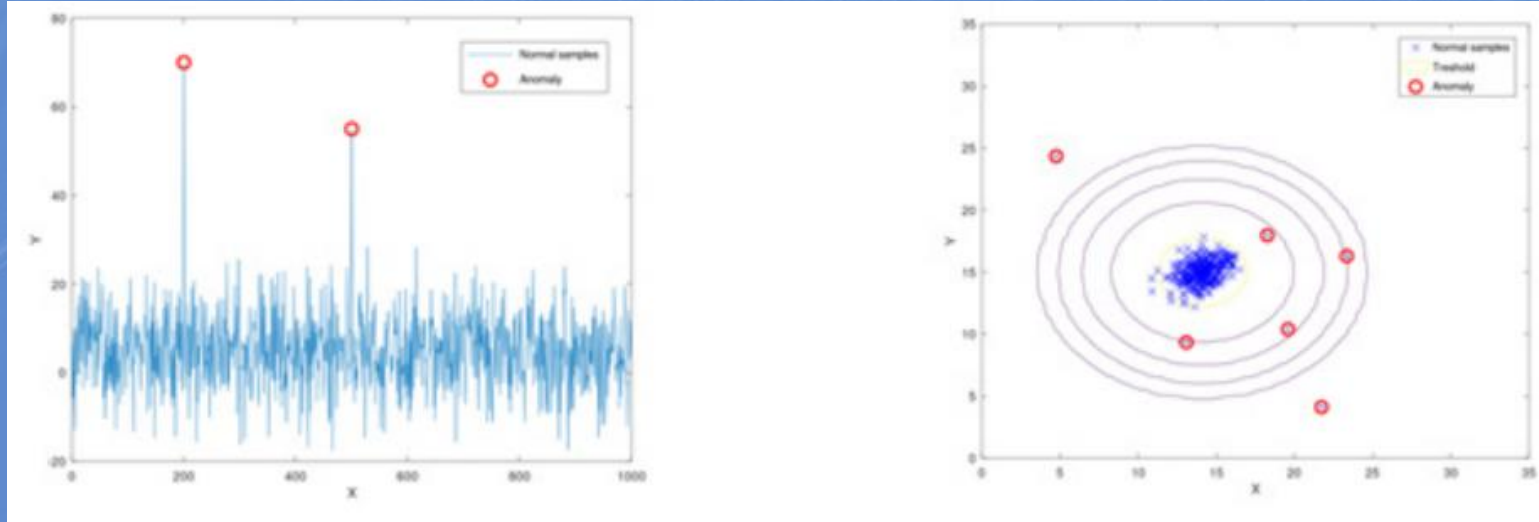
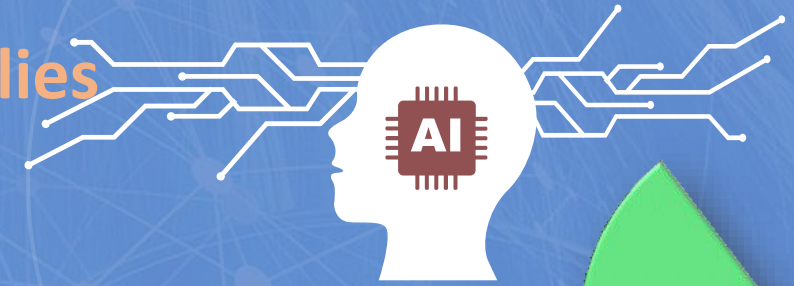
*(une taille de cercle plus grande implique une fréquence de la règle plus importante, la couleur du cercle implique la levée de la règle (force de la règle)).*



L'association est généralement utilisée dans les commerces de détail et en ligne pour analyser et comprendre les habitudes d'achat des clients et pour créer une stratégie marketing plus optimisée et ciblée. Les règles d'association peuvent également être utilisées dans la planification et la conception de l'aménagement, de l'aménagement commercial à l'urbanisme.



# Apprentissage non supervisé – Détection d'anomalies

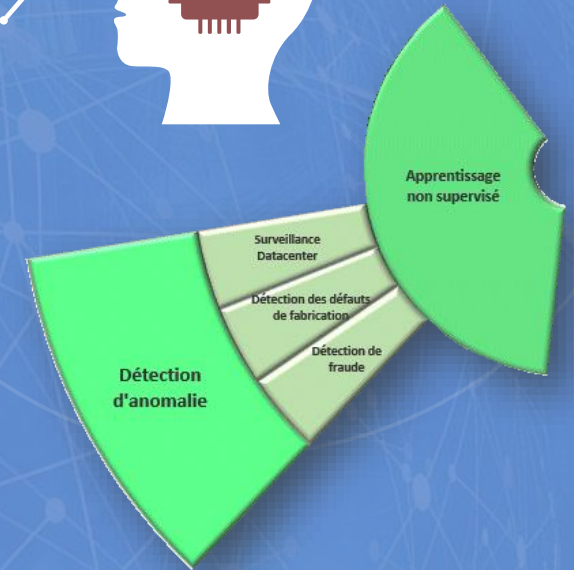


*Anomalies selon une distribution temporelle*

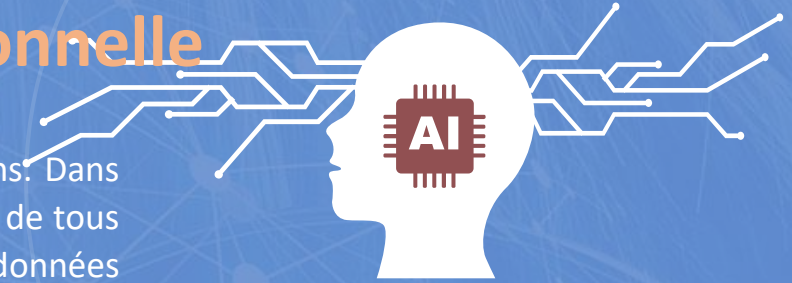
*Anomalies selon une distribution Gaussienne*

Reconnaître et identifier les anomalies dans les données peut présenter de nombreux avantages, en supposant que la plupart des données et des flux sont statistiquement cohérents. C'est-à-dire en suivant (par exemple) la distribution gaussienne, où tout écart significatif peut signifier des problèmes potentiels à étudier et à analyser plus avant. Dans le secteur des services financiers, la détection d'anomalies peut être utilisée pour identifier la fraude, car des habitudes de dépenses inhabituelles peuvent signifier un compte compromis ou un vol de carte de crédit.

Dans la fabrication, la détection des anomalies peut aider au contrôle qualité automatisé, car une anomalie peut signifier un produit défectueux. Dans la gestion des centres de données, il peut identifier une machine défectueuse ou un processus bloqué. Dans le domaine de la cybersécurité, un accès inhabituel aux ressources ou une activité réseau inhabituelle peut signifier un problème de sécurité potentiel ou une violation du réseau. Avec l'ampleur croissante de l'Internet des objets (IoT) et des concepts tels que les villes intelligentes qui dépendent fortement de l'augmentation des entrées de données provenant de divers capteurs, la détection d'anomalies pourrait être utilisée pour détecter des incidents réels (par exemple, liés au trafic, aux infrastructures urbaines ou à la navigation). Ces entrées pourraient éventuellement déclencher des systèmes d'alerte ou permettre des actions préventives et correctives avant que l'incident ne se produise.



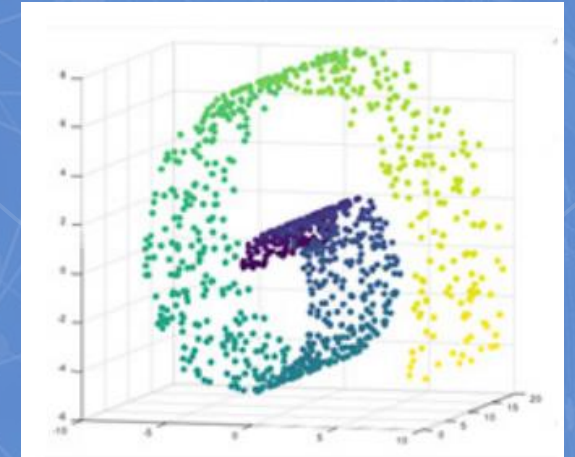
# Apprentissage non supervisé – Réduction dimensionnelle



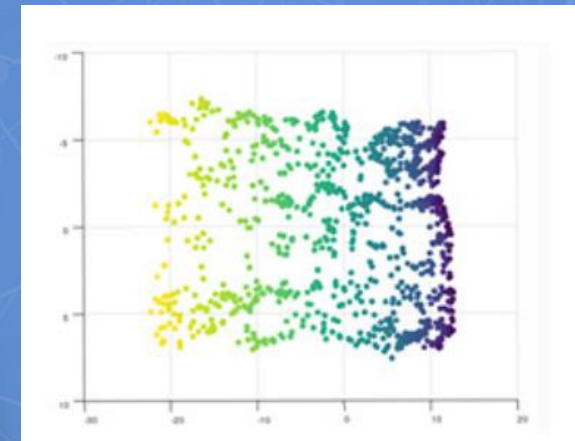
Lorsque nous analysons des données, nous pouvons constater qu'elles ont de nombreuses dimensions. Dans les données de santé, nous étudions généralement les antécédents et les paramètres physiométriques de tous les patients. L'industrie de l'assurance essaie de créer des modèles plus précis en prenant des données provenant de diverses sources pour des évaluations des risques plus précises. Les données de diffusion Web peuvent inclure des centaines ou des milliers de dimensions différentes avec des données fortement corrélées. Dans les ensembles de données très volumineux, généralement produits par le biais de mégadonnées, plusieurs dimensions peuvent contenir des redondances (par exemple, la hauteur en pieds, mètres et centimètres) ou des données non pertinentes pour un besoin spécifique. La réduction dimensionnelle simplifie l'analyse des données en créant un sous-ensemble d'entités de données ou en extrayant des ensembles spécifiques d'entités de données pour créer un nouvel ensemble de données.

Une IA traitant de la reconnaissance de formulaires pourrait, par exemple, convertir une image haute résolution colorée en une image noir et blanc, à une valeur d'intensité de couleur unique par pixel, à une résolution inférieure suffisante pour les tâches de reconnaissance ultérieures.

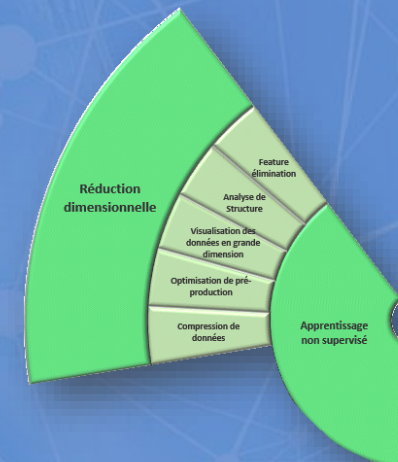
La réduction dimensionnelle est généralement destinée à améliorer l'efficacité de calcul pour les processus ultérieurs. Elle peut également être utilisée dans la visualisation de données, d'ensembles de données de grande dimension, réduites et affichées sous forme de visualisations 3D ou 2D.



Visualisation 3D



Données dépliées en 2D

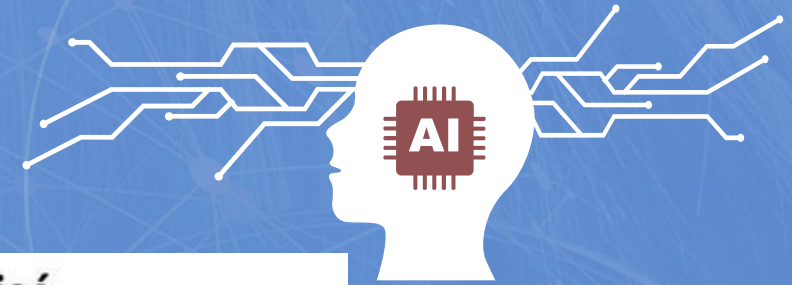




**Pour aller plus  
loin**



# Principales méthodes de machine learning



## supervisé

### classification

données avec label  
(pixels) -> (nombre)



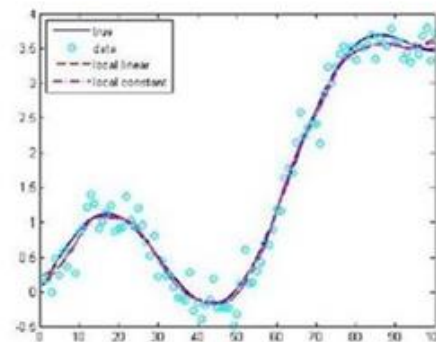
**regrouper** automatiquement les objets en classes, et prédire l'appartenance d'un nouvel objet à une classe identifiée

**type d'objet** : complexes (image, voix, ...)

**exemples** : prédire si une tumeur est cancéreuse en fonction de critères multiples, identifier un spam

### régression

données avec label  
prévoir (y) en fonction de (x)



**prédire** une valeur en fonction de données d'entraînement multidimensionnelles

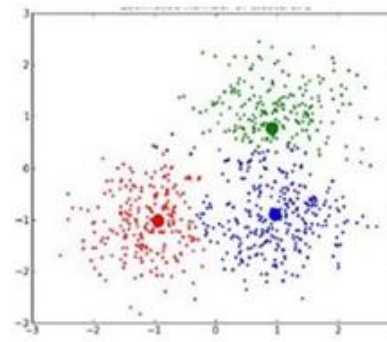
**type d'objet** : valeurs numériques continues

**exemples** : anticipation de churn client, de demande client, évaluation de pipe client, prévision de panne, prévision de récive

## non supervisé

### clustering

données sans label  
(x, y, z, ...)



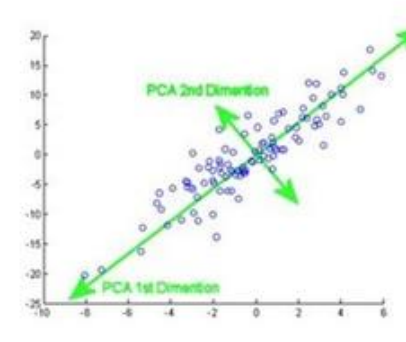
**segmenter** automatiquement un jeu de vecteurs (x, y, z, etc)

**type d'objet** : ensemble de n-uplets de valeurs numériques

**exemples** : détection de fraude, blanchiment d'argent sale, détection de faille de sécurité

### réduction dimensions

données sans label  
(x, y, z, ...)



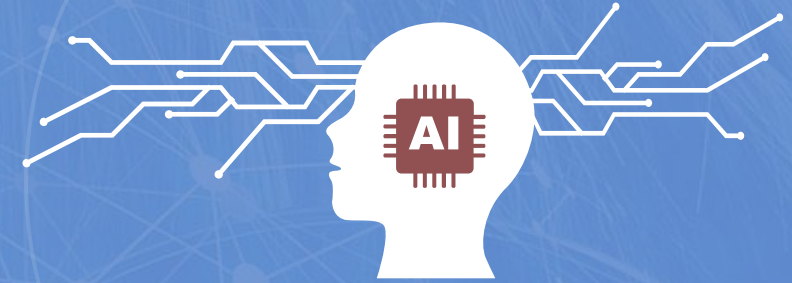
**déterminer** automatiquement les paramètres discriminants d'un jeu de données par rapport à une variable cible

**type d'objet** : ensemble de n-uplets de valeurs numériques

**exemples** : identifier les paramètres déterminants la corrélation entre des paramètres clients et leur comportement futur



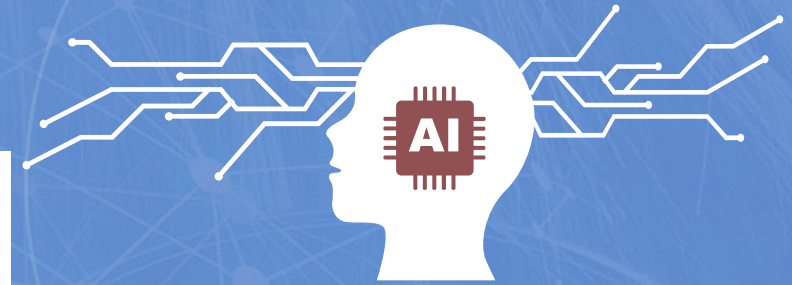
# Comparaison des 2 apprentissages



Classification	Clustering (Regroupement)
Utilisation de données labellisées en entrée	Utilisation de données non labellisées en entrée
2 phases	Seule phase
La valeur de sortie est connue	La valeur de sortie n'est pas maîtrisée
Domaine de l'apprentissage machine supervisé	Domaine de l'apprentissage machine non supervisé
Des données d'entraînement sont nécessaires	Des données d'entraînement ne sont pas nécessaires
Exemple d'algorithmes : Support vecteur machine (SVM), arbres de décision, méthode Bayésienne...	Exemple d'algorithmes : Regroupement hiérarchique, k-means, algorithme DBScan ...
Peut être plus complexe que le Clustering	Peut être moins complexe que le Clustering
Ne précise pas les axes d'amélioration	Précise les axes d'amélioration
Les conditions aux limites doivent être précisées	Les conditions aux limites ne sont pas toujours précisées



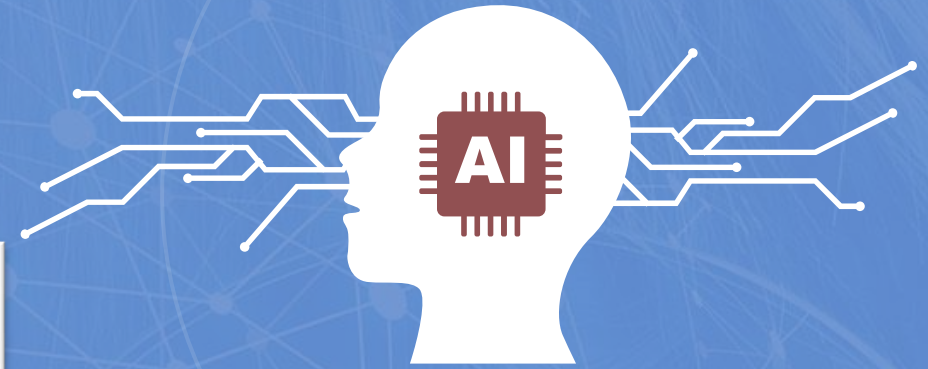
# Utilisation des algorithmes de l'IA



	EDUCATION	JUSTICE	SANTE	SECURITE	TRAVAIL	CULTURE	AUTRES
Générer du savoir	Mieux cerner les aptitudes d'apprentissage des élèves	Mettre en évidence des manières différentes de rendre la justice selon les régions	Tirer profit de la quantité immense de publications scientifiques	Repérer des liens insoupçonnés pour la résolution d'enquêtes par les services de gendarmerie	Comprendre les phénomènes sociaux en entreprise	Créer des oeuvres culturelles (peinture, musique...)	Affiner le profil de risque d'un client d'un assureur
Croiser des infos, «matching»	Répartir les candidats dans les formations d'enseignement (ex APB)		Répartir des patients pour participation à un essai clinique		Faire correspondre une liste de candidatures avec une offre d'emploi		Mettre en relation des profils « compatibles » sur des applications de rencontres
Prédire	Prévoir des décrochages scolaires	Prédire la chance de succès d'un procès et le montant potentiel de dommages intérêts	Prédire des épidémies Repérer des prédispositions à certaines pathologies afin d'en éviter le développement	Détecter les profils à risque dans la lutte antiterroriste Prédire l'occurrence future de crimes et délits	Détecter les collaborateurs qui risquent de démissionner dans les prochains mois	Créer des oeuvres ayant un maximum de chance de plaire aux spectateurs (Netflix)	
Conseiller	Recommander des voies d'orientation personnalisées aux élèves	Recommander des solutions de médiation en fonction du profil des personnes et de cas similaires			Proposer des orientations de carrière adaptées aux profils	Recommander des livres (Amazon), des séries télévisées (Netflix), etc.	Individualiser des messages politiques sur les réseaux sociaux
Aider à la décision		Suggérer au juge la solution jurisprudentielle la plus adéquate pour un cas donné	Suggérer au médecin des solutions thérapeutiques adaptées	Proposer aux forces de police les zones prioritaires dans lesquelles patrouiller			Aider à trouver le chemin le plus court (GPS)



# Approfondissements



A voir sur le site de la Dane

« Les réseaux de neurones »



A venir

« Apprentissage par renforcement »

« Reconnaissance d'images, d'objet ou de visage »

The 7 steps of machine learning :

