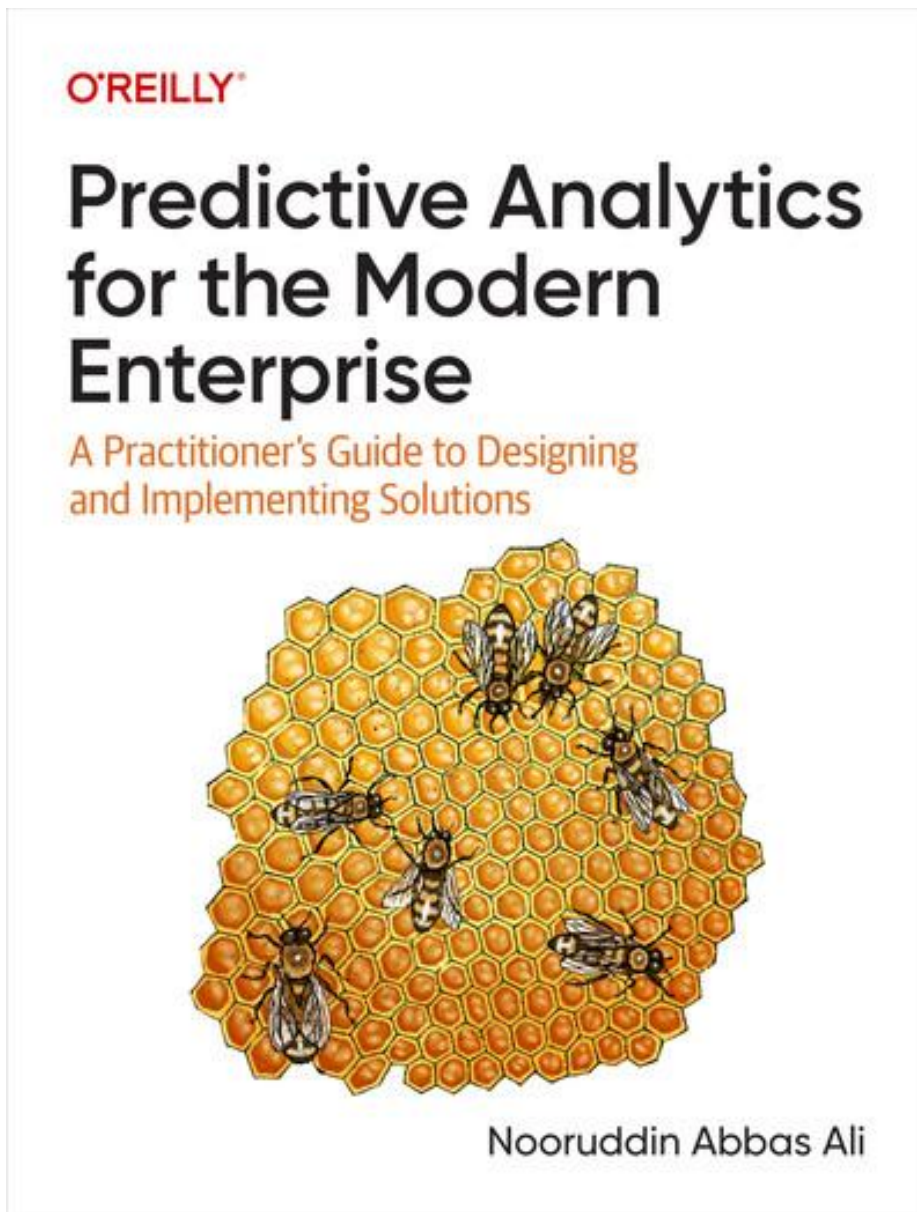


## Lecture Individuelle 2

### Les analyses prédictives



1

---

<sup>1</sup> <https://learning.oreilly.com/covers/urn:orm:book:9781098136857/400w/>

## Table des matières

1. Introduction.....	4
2. Différents types d'analyse de données.....	5
2.1. Analyse descriptive .....	5
2.2. Analyse diagnostique.....	5
2.3 Analyse prédictive.....	6
2.4 Analyse prescriptive.....	6
3. Acquisition de connaissances, apprentissage automatique et le rôle de l'analyse prédictive .....	8
3.1. Les bases de l'apprentissage automatique .....	8
3.2. Les types d'apprentissage automatique .....	9
3.3. Le rôle de l'analyse prédictive.....	9
3.4. Applications pratiques .....	9
3.5. Importance des prédictions précises.....	10
4. Outils, frameworks et plateformes dans le monde de l'analyse prédictive .....	11
4.1. Technologies et systèmes .....	11
4.2. Langages et bibliothèques .....	11
4.3. Services cloud pour l'IA et l'apprentissage automatique .....	11
4.3.1. Service AWS.....	11
4.3.2. Services GCP .....	12
5. Les défis liés à l'utilisation de l'analyse prédictive .....	13
5.1. Les personnes .....	13
5.2. Les données .....	14
5.3. La technologie.....	16
6. Les mathématiques et algorithmes derrière l'analyse prédictive .....	17
6.1. Les statistiques et algèbre linéaire .....	17

6.2. La régression .....	20
6.2.1. Analyse de la régression .....	20
6.2.2. Types de régression .....	21
6.3. Les arbres de décision .....	23
7. Traitement des données .....	26
7.1. Gestion des valeurs manquantes.....	26
7.2. Encodage des variables catégorielles .....	26
7.3. Transformation des données .....	26
7.4. Gestion des valeurs aberrantes (outliers).....	27
7.5. Gestion du déséquilibre des classes .....	27
7.6. Combinaison de données .....	27
7.7. Sélection des caractéristiques .....	28
7.8. Division des données .....	28
8. Python et scikit-learn pour les analyses prédictives .....	29
8.1. NumPy : manipulation efficace de données numériques .....	29
8.2. Pandas : gestion et préparation des données .....	29
8.3. Matplotlib : visualisation des données .....	30
8.4. Scikit-learn : modélisation et prédiction.....	31
8.5. Entraînement et prédiction avec un modèle de régression linéaire .....	32
8.6. Entraînement et prédiction avec un modèle d'arbre de décision .....	32
9. TensorFlow et Keras.....	35

## 1. Introduction

Le livre "Predictive Analytics for the Modern Enterprise" de Nooruddin Abbas Ali explore l'impact de l'analytique prédictive dans le monde des entreprises modernes. Dans un environnement où la prise de décisions basées sur des données devient incontournable, ce livre fournit un cadre précieux pour comprendre comment les outils d'analyse prédictive peuvent transformer les stratégies commerciales et opérationnelles. En abordant des concepts comme la modélisation statistique, les algorithmes d'apprentissage automatique et l'intégration des résultats dans les processus métier, l'auteur offre une perspective complète sur l'utilisation de l'analytique pour anticiper l'avenir et optimiser les performances.

Cette présentation vise à explorer les principales idées du livre, en mettant en lumière les méthodes abordées, les cas d'application réels, ainsi que les défis et les opportunités que présente l'utilisation de l'analytique prédictive dans le contexte actuel des entreprises. L'objectif est de montrer comment ces outils peuvent non seulement améliorer la prise de décision, mais aussi créer un avantage concurrentiel pour les organisations modernes.

Ce livre s'inscrit parfaitement dans le cadre de l'acquisition des compétences pour ce 5<sup>ème</sup> semestre, notamment la compétence B4 « Connaître les impacts et les apports du Machine Learning et de l'intelligence artificielle sur le système d'information de l'entreprise », ainsi que la compétence M7 « Connaître les principaux concepts mathématiques... et savoir les appliquer dans un cas d'utilisation du Machine Learning ».

## 2. Différents types d'analyse de données

L'analyse des données permet de mieux comprendre une entreprise en fournissant des informations sur ce qui s'est passé, ce qui pourrait se produire, et les actions à envisager. Cette pratique repose sur plusieurs types d'analyses, que l'on peut regrouper en différentes catégories.

### 2.1. Analyse descriptive

L'analyse descriptive répond à la question : "Que s'est-il passé ?" Elle se concentre sur l'examen des données historiques et actuelles pour identifier des tendances ou des comportements.

Exemple d'application :

- Domaine médical : les médecins collectent des informations sur les antécédents médicaux, les symptômes actuels et les constantes vitales d'un patient. Ces données permettent de dresser un portrait global de l'état de santé du patient.
- Plateformes de streaming : des entreprises comme Netflix analysent les contenus regardés pour identifier les genres ou titres populaires selon les périodes et les groupes démographiques.

Techniques utilisées :

- Agrégation : moyennes ou résumés à partir de larges ensembles de données.
- Segmentation : analyse de sous-ensembles spécifiques de données pour identifier des tendances ou des anomalies.

### 2.2. Analyse diagnostique

L'analyse diagnostique cherche à répondre à la question : "Pourquoi cela s'est-il produit ?" Elle vise à comprendre les relations de cause à effet entre différents événements ou variables.

Points clés :

- Différence entre corrélation et causalité : une corrélation indique une relation entre deux variables, mais pas nécessairement que l'une cause l'autre.
- Exemples d'utilisation :

- Un médecin investigate les causes d'une hypertension chez un patient (stress, alimentation, etc.).
- Une entreprise analyse les raisons pour lesquelles un produit s'est bien vendu à un moment précis.

Méthodes courantes :

- Analyse manuelle des tendances.
- Utilisation d'outils comme les modèles de régression ou Excel pour identifier des relations complexes entre des variables.

### 2.3 Analyse prédictive

L'analyse prédictive cherche à déterminer : "Que pourrait-il se passer ?" Elle utilise des données historiques et actuelles pour prédire des événements futurs avec une certaine probabilité.

Applications pratiques :

- Médecine : un médecin peut estimer les risques de maladies cardiovasculaires en fonction des antécédents médicaux et du mode de vie d'un patient.
- Commerce : les entreprises utilisent des modèles prédictifs pour personnaliser les recommandations de produits, comme le fait Netflix pour suggérer des contenus adaptés à chaque utilisateur.

Techniques :

- Modèles de données prédictifs et apprentissage automatique (Machine Learning).
- Analyse de scénarios simples (binaire) ou complexes (comme l'optimisation des prix).

### 2.4 Analyse prescriptive

L'analyse prescriptive répond à la question : "Que devrions-nous faire ?" Elle aide à identifier la meilleure action à entreprendre en tenant compte des contraintes et des objectifs d'une entreprise.

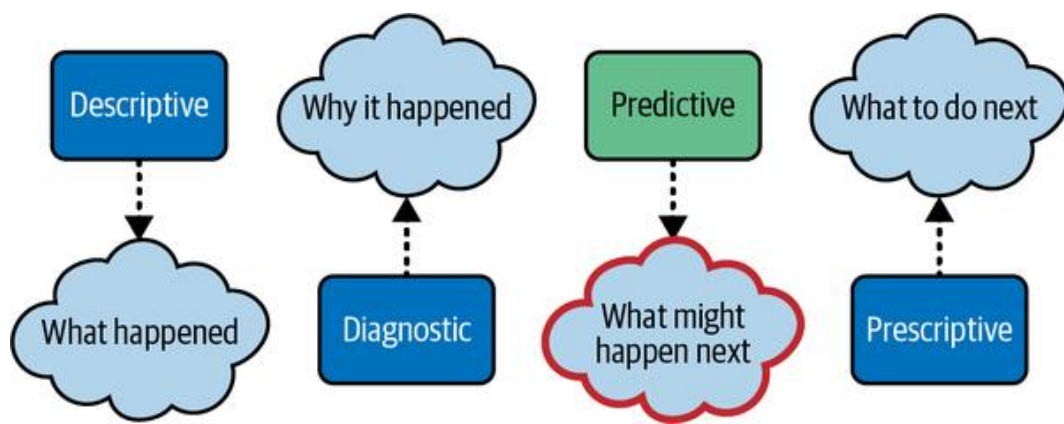
Caractéristiques :

- Intégration de plusieurs variables et scénarios pour déterminer l'action optimale.

- Utilisation de techniques d'optimisation comme la programmation linéaire ou les modèles de simulation.

Exemples d'application :

- Simulations pour prédire l'impact de certaines décisions sur un système complexe (par exemple, la propagation d'une maladie).
- Optimisation de processus logistiques pour réduire les coûts tout en respectant les délais.



### 3. Acquisition de connaissances, apprentissage automatique et le rôle de l'analyse prédictive

Lorsque vous conduisez, votre réaction naturelle face à un piéton est de ralentir. Décider de s'arrêter ou de continuer, ainsi que de maintenir une distance de sécurité, dépend de nombreux facteurs. Est-ce un enfant, un adulte ou une personne âgée ? Regarde-t-il le trottoir ou s'apprête-t-il à traverser ? Que révèle son langage corporel ? Y a-t-il un passage piéton à proximité ? Ces questions, parmi tant d'autres, sont évaluées mentalement pour prédire le comportement du piéton et adapter votre conduite en conséquence.

Ce processus repose sur l'acquisition de connaissances :

- Connaissances innées (instincts et comportements).
- Connaissances apprises (via des institutions, livres, échanges ou Internet).
- Connaissances acquises par l'expérience.

De manière similaire, les machines, omniprésentes dans nos vies, doivent également apprendre et s'adapter pour être efficaces. L'apprentissage automatique (machine learning) est la tentative humaine de reproduire ce processus d'apprentissage chez les machines.

#### 3.1. Les bases de l'apprentissage automatique

L'apprentissage automatique repose sur des algorithmes qui font des prédictions :

- Prédictions simples : classer un e-mail comme étant ou non un spam.
- Prédictions complexes : estimer la durée de vie d'un équipement industriel.

Trois composantes principales soutiennent ce processus :

- Le processus de décision : les données d'entrée sont analysées pour détecter des motifs.
- La fonction d'erreur : elle évalue l'exactitude des prédictions, en mesurant l'écart avec les données réelles.
- L'optimisation : elle ajuste les modèles pour minimiser l'erreur, permettant ainsi une amélioration continue et automatisée.



### 3.2. Les types d'apprentissage automatique

Apprentissage supervisé : les algorithmes utilisent des données étiquetées (avec des résultats connus) pour s'entraîner. Par exemple, prédire si une tumeur est bénigne ou maligne.

Apprentissage non supervisé : il identifie des modèles ou relations dans des données non étiquetées, sans aide externe.

Apprentissage semi-supervisé : combine des petits ensembles de données étiquetées avec de grandes quantités de données non étiquetées.

Apprentissage par renforcement : le modèle apprend par essais et erreurs, en recevant des récompenses pour les succès et des pénalités pour les échecs.

### 3.3. Le rôle de l'analyse prédictive

L'analyse prédictive permet aux entreprises de transformer des données brutes en connaissances, pour prévoir l'avenir et optimiser leurs décisions. Cela nécessite un cycle continu :

- Collecte et stockage des données issues de diverses sources.
- Prétraitement des données pour :
  - o Les adapter au contexte métier.
  - o Les nettoyer, classer et étiqueter si nécessaire.
- Analyse descriptive et diagnostique pour comprendre ce qui s'est passé et pourquoi.
- Prédiction des tendances futures.
- Décisions basées sur les prédictions.

### 3.4. Applications pratiques

Exemples de questions métiers et de prédictions associées :

- Hôtellerie : Combien de personnel déployer en fonction de la saison et de la météo ?
  - o Prédiction : nombre de clients attendus à chaque emplacement.
- Logistique : Quel stock maintenir pour la saison à venir ?
  - o Prédiction : quantité prévue de ventes pour chaque produit.
- Énergie : Combien d'énergie produire pour le prochain cycle ?
  - o Prédiction : consommation estimée des clients.

### 3.5. Importance des prédictions précises

Des prédictions erronées peuvent entraîner :

- Des coûts inutiles (surstockage, sureffectif).
- Des opportunités perdues (ruptures de stock, sous-effectif).
- Une insatisfaction client (manque de ressources, délais).

Ainsi, l'analyse prédictive est essentielle pour garantir l'efficacité opérationnelle et offrir une expérience client optimale. En exploitant efficacement leurs données, les entreprises peuvent s'adapter et prospérer dans un environnement de plus en plus compétitif.

## 4. Outils, frameworks et plateformes dans le monde de l'analyse prédictive

Les entreprises exploitent les analyses depuis longtemps, et leur usage ne cesse de croître. Dans un environnement concurrentiel et rapide, les organisations cherchent constamment à réduire leur délai de mise sur le marché et à améliorer l'expérience client pour maintenir leur avantage concurrentiel. L'analyse prédictive est un levier clé pour comprendre leur activité, anticiper les besoins des clients et réagir rapidement.

### 4.1. Technologies et systèmes

- Hadoop : un système de stockage distribué open source, basé sur un article technique de Google, qui permet de gérer de très grands ensembles de données. Associé à Hadoop, MapReduce facilite leur traitement.
- Spark et Flink : projets qui répondent aux besoins de traitement en temps réel. Des couches d'abstraction comme Pig et des bases de données non relationnelles comme HBase enrichissent l'écosystème Hadoop.

Les organisations choisissent de déployer ces plateformes sur site ou dans le cloud, en mode autogéré ou via des services SaaS proposés par des tiers.

### 4.2. Langages et bibliothèques

- Python : grâce à sa simplicité et à ses fonctionnalités, Python est devenu un pilier du développement d'applications analytiques. Il propose des bibliothèques spécialisées comme TensorFlow, conçue par Google pour effectuer des calculs numériques complexes et développer des applications d'apprentissage automatique.

TensorFlow est accompagné d'API comme Keras, qui simplifie l'entraînement et le test des modèles prédictifs.

### 4.3. Services cloud pour l'IA et l'apprentissage automatique

Les principaux fournisseurs de cloud, comme **AWS** et **GCP**, offrent des services spécifiques pour développer des fonctionnalités d'analyse prédictive.

#### 4.3.1. Service AWS

- SageMaker : un service complet pour le cycle de vie de l'apprentissage automatique, intégrant des outils comme Jupyter Notebook.

- Amazon Personalize : service prêt à l'emploi pour personnaliser l'expérience utilisateur via des recommandations produits.
- Amazon Forecast : aide à prévoir la demande grâce à des données historiques et météorologiques.
- Amazon Lookout for Equipment : permet la maintenance prédictive grâce à des données issues de capteurs industriels.

Ces services sont soutenus par des outils comme AWS Lambda (code sans serveur), AWS Kinesis (analyse en temps réel de données de flux) et AWS S3 (stockage d'objets).

#### 4.3.2. Services GCP

- Vertex AI : plateforme centralisée pour gérer le cycle complet des modèles d'apprentissage automatique, intégrant TensorFlow, scikit-learn et PyTorch.
- AutoML : automatise la création et le déploiement des modèles ML basés sur des données structurées.
- Timeseries Insights API : analyse les séries temporelles en temps réel pour détecter les anomalies et anticiper les tendances.
- BigQuery : entrepôt de données cloud avec des capacités ML intégrées, compatible avec des sources multcloud et des analyses géospatiales.

## 5. Les défis liés à l'utilisation de l'analyse prédictive

L'utilisation efficace de l'analytique prédictive au sein d'une organisation nécessite une approche globale du paradigme. Bien qu'il soit possible de commencer par un cas d'utilisation pour tester les eaux de l'analytique prédictive, le passage de la production de données à une organisation axée sur les données nécessite des changements à plusieurs niveaux au sein de l'organisation.

L'adoption de l'analytique prédictive implique des changements à trois niveaux :

- Les personnes
- Les données
- La technologie

### 5.1. Les personnes

Comme pour toute forme de changement, l'introduction de l'analytique prédictive entraîne son lot de défis. L'une des méthodes les plus créatives pour amener les employés à accepter ce changement est de leur poser la question suivante : « Selon vous, quel est le plus grand atout de votre organisation ? » Les réponses varieront, mais elles incluront généralement les employés, la réputation, le service, le produit et les clients de l'organisation. Si cela est également vrai pour votre organisation, il est essentiel de commencer à éduquer vos employés sur l'importance des données. Ils doivent comprendre que les données ne sont pas simplement une information produite par les applications de l'organisation ; elles constituent un atout majeur qui, lorsqu'elles sont utilisées correctement, peut avoir un impact significatif sur le résultat financier de l'entreprise.

Il est impératif que les managers, lorsqu'ils parlent de services et d'applications, pensent proactivement à la manière dont les données sont utilisées pour l'entreprise. Mais l'éducation ne doit pas s'arrêter là. Un ingénieur en informatique doit apprécier l'importance des données et prévoir des pipelines de données efficaces ainsi qu'une consolidation des données. Lors de la création d'applications, les développeurs doivent garder à l'esprit que les données produites par l'application doivent être facilement consommables par d'autres, et vice versa. Bien que j'évoque ici l'analytique prédictive, ce changement de mentalité s'applique également à l'ensemble du domaine de l'analyse des données.

L'analytique prédictive, comme beaucoup de bonnes choses dans la vie, demande du temps, des efforts et de l'argent pour être mise en œuvre. Cependant, contrairement à un abonnement à Netflix, où les raisons de payer sont immédiatement évidentes (du moins pour le premier mois), les bénéfices de l'analytique prédictive ne sont pas aussi évidents. Il est donc essentiel pour toute organisation de comprendre pourquoi elle entreprend cette démarche et ce qu'elle en retirera.

D'un point de vue technique, comprendre la douleur organisationnelle ou l'opportunité, ainsi que calculer le retour sur investissement (ROI) de cette initiative, garantira sa pérennité et son succès futur. Cela permettra de comprendre le problème à résoudre ou l'opportunité ainsi que les ressources nécessaires en temps et en argent. Il s'agit également de définir les critères nécessaires pour évaluer le succès ou l'échec de l'initiative. Compte tenu des nombreux changements requis au niveau organisationnel, un cas d'affaires bien défini, soutenu par un ROI positif, garantira le soutien nécessaire à différents niveaux.

Une autre considération à long terme (bien que non strictement liée aux personnes) est d'avoir en place des processus pour l'introduction, le développement, les changements et l'abandon des applications métier. À mesure que l'analytique prédictive se développe au sein de l'organisation, les changements apportés au paysage applicatif auront un impact direct sur le pipeline de données. Puisque ce pipeline est désormais lié, directement ou indirectement, au résultat financier de l'entreprise, ces changements doivent être régulés comme tout autre changement au niveau de l'entreprise.

## 5.2. Les données

Tout dans le domaine de l'analytique prédictive repose sur les données. Cependant, lorsqu'on se lance dans l'implémentation de l'analytique prédictive, plusieurs défis liés aux données doivent être pris en compte. Les données organisationnelles pertinentes pour l'entreprise ne se trouvent que rarement en un seul endroit. À mesure que les organisations évoluent de manière organique, leur paysage applicatif se développe également. Cela signifie que les données connexes se retrouvent souvent dans des silos. Ces silos résultent souvent d'une croissance organique (et parfois mal planifiée).

Cependant, des facteurs humains tels que les politiques internes ou des raisons bureaucratiques peuvent également empêcher les équipes de partager leurs données. Des

raisons techniques, comme l'absence d'une approche basée sur des microservices, de services d'intégration API, ou de contrôles d'authentification et d'autorisation granulaires, peuvent aussi nuire au partage de données saines entre les départements d'une organisation. L'un des aspects clés pour résoudre ce problème des données est de garantir que celles-ci puissent être collectées, stockées et partagées de manière évolutive et sécurisée. En résumé, les données doivent être gérées par une plateforme qui soit :

- Assez agile pour gérer différents types et structures de données
- Assez scalable pour traiter de très grands volumes de données provenant de toute l'organisation
- Assez sécurisée pour garantir que les données puissent être suivies et rendues accessibles aux utilisateurs et services autorisés

Bien entendu, ces problèmes sont d'autant plus complexes lorsque vous tentez de collecter des données au niveau d'un groupe à travers plusieurs sous-organisations. Dans ce cas, les silos de données entraînent plusieurs défis :

- Données incomplètes (par exemple, des données de transactions sans les détails du client)
- Données en double (par exemple, le même client existe dans plusieurs lignes de produits comme des clients distincts)
- Accès aux données pertinentes non disponible pour certains départements (par exemple, le marketing ne peut pas utiliser les nouvelles informations d'inventaire pour envoyer des offres personnalisées aux clients)

Cela conduit à des données sales. Il est erroné de penser que les données sales ne désignent que des données erronées ou incomplètes. Les données sales peuvent également être le résultat d'un manque de normalisation, de duplication, ou de l'absence de relation entre les données. Un aspect des données sales particulièrement pertinent pour l'analytique prédictive est leur tendancieuse. Cela peut être dû à un biais existant lors de la création des données. Par exemple, considérons une entreprise ayant des opérations principalement aux États-Unis et qui s'étend en Europe. Si la majorité de ses clients se trouvent aux États-Unis, un modèle d'analytique prédictive formé avec ces données pourrait en venir à prédire que les prospects américains sont plus susceptibles de devenir des clients que les prospects européens. Ce sont

là quelques-uns des défis que les organisations doivent surmonter avant de se lancer dans l'analytique prédictive.

### 5.3. La technologie

Concentrer une grande quantité de données provenant de toute l'organisation nécessite un investissement important dans la recherche, la conception, l'implémentation et l'exploitation des infrastructures informatiques et technologiques. Ajoutez à cela la complexité de l'entraînement (et parfois de la construction) des modèles d'analytique prédictive, et vous comprendrez que des investissements significatifs doivent être réalisés en matière d'infrastructure matérielle et logicielle pour concrétiser cette réalité.

Même si votre organisation a abordé les aspects culturels et technologiques de l'implémentation de l'analytique prédictive, il vous faut encore des personnes compétentes en modélisation des données et en apprentissage automatique. Vous avez également besoin de personnes capables de manipuler les données de manière créative pour les rendre accessibles à l'entreprise, et de personnes qui peuvent créer un lien efficace entre le côté technique et le côté business des données. Ces individus, tels que des data scientists, des data engineers, et des consultants business, possèdent des compétences spécialisées difficiles à trouver et encore plus difficiles à retenir.

C'est là qu'interviennent les fournisseurs de cloud, tels qu'AWS, Azure, GCP et Alibaba Cloud, ainsi que quelques entreprises spécialisées dans les données. Ils proposent des solutions tout en réduisant les exigences techniques pour permettre aux organisations de commencer à utiliser l'analytique prédictive. Bien que ces solutions évoluent rapidement, la quantité de connaissances nécessaires pour utiliser ces services et leur capacité à offrir un véritable avantage concurrentiel restent à débattre.



## 6. Les mathématiques et algorithmes derrière l'analyse prédictive

L'apprentissage est une caractéristique fondamentale de l'être humain. Nos expériences sont nos enseignements, et elles façonnent notre pensée et la manière dont nous vivons. Ce processus de maturation intellectuelle ne se limite pas à un individu mais s'étend à l'ensemble de notre espèce. Grâce à cet apprentissage collectif, l'Homme a su trouver des remèdes pour des maladies mortelles, construire des avions, explorer l'espace interstellaire, et observer des galaxies lointaines.

En tant qu'évolution intellectuelle, nous avons aussi créé des ordinateurs pour nous aider dans nos vies quotidiennes. Aujourd'hui, les ordinateurs sont présents dans presque tous les aspects de notre existence. Mais avant les ordinateurs, les mathématiques étaient déjà omniprésentes dans notre monde. Les mathématiques sont partout et, avec suffisamment de compréhension et d'effort, on peut représenter quasiment n'importe quelle situation sous forme de problème mathématique. C'est pourquoi on dit souvent que les mathématiques sont un langage universel.

Il n'est donc pas surprenant que les mathématiques soient à la base de l'informatique, depuis la représentation binaire des données au niveau du matériel jusqu'à la programmation logicielle et les algorithmes d'apprentissage automatique.

### 6.1. Les statistiques et algèbre linéaire

Les statistiques sont l'art (ou la science) de collecter des données numériques, de les comprendre et de les représenter, et d'en tirer des conclusions. Cela fait écho à la définition de l'analyse prédictive, où nous utilisons des données historiques et actuelles pour prédire des événements futurs. Cette similitude fondamentale n'est pas un hasard, car les mathématiques sous-jacentes de l'analyse prédictive reposent sur l'analyse statistique de grandes quantités de données, ainsi que sur d'autres formes de mathématiques appliquées, comme l'algèbre linéaire, la probabilité et le calcul.

Dans cette section, nous ne plongerons pas dans les statistiques spécifiques utilisées par les différents algorithmes d'analyse prédictive. Nous allons plutôt expliquer les concepts fondamentaux afin que, la prochaine fois que vous interagirez avec l'analyse prédictive, cela vous paraisse plus concret et moins mystérieux.

Imaginons un problème simple. Vous êtes propriétaire d'un magasin d'alimentation et vous avez 100 articles à vendre. Vous souhaitez répondre aux questions suivantes :

- Combien de chaque article je vends chaque mois ?
- Quels articles se vendent le plus (et le moins) dans un mois donné ?
- Quel est mon chiffre d'affaires moyen mensuel ?
- Comment mes ventes évoluent-elles d'un mois à l'autre ?

Toutes ces questions peuvent être résolues en collectant, traitant et analysant vos données de vente. Les statistiques permettent de traiter ces données pour en tirer des informations utiles qui aideront votre entreprise. Ces données peuvent être représentées visuellement ou sous forme de rapports pour une analyse plus poussée.

Voici quelques-unes des informations que vous pouvez obtenir de cette analyse :

- Quels articles génèrent le plus de revenus, et comment puis-je maintenir et améliorer leurs ventes ?
- Quels articles génèrent le moins de revenus, pourquoi, et que dois-je faire ensuite ?
- Quel est le bilan global de l'entreprise en termes de revenus mois par mois, et comment puis-je définir des objectifs commerciaux spécifiques et élaborer un plan pour les atteindre ?

Si l'entreprise se développe et devient une chaîne de magasins, la quantité et la complexité des données augmentent considérablement. Certaines questions restent les mêmes, comme « Comment faire croître mon entreprise ? », mais d'autres apparaissent, telles que « Pourquoi un magasin performe-t-il mieux qu'un autre ? » ou « Comment gérer efficacement l'inventaire pour répondre aux besoins du business ? » Pour répondre à ces questions, les statistiques nous offrent un ensemble d'outils.

### **Moyenne**

Elle donne la valeur moyenne d'une distribution de fréquences. Dans notre exemple, elle permet de répondre à la question du chiffre d'affaires moyen.

### **Médiane**

Elle représente la valeur qui se trouve au milieu d'une distribution de valeurs, de sorte qu'il y a une probabilité égale qu'une observation se situe au-dessus ou en dessous. La médiane nous aiderait à comprendre le prix de vente typique d'un article sans être influencée par des valeurs aberrantes.

### **Mode**

Elle représente l'observation la plus fréquemment observée. Cela serait utile pour répondre à la question de l'article le plus vendu.

### **Écart-type**

Il nous aide à comprendre à quel point notre distribution est dispersée par rapport à la moyenne. Cela nous permettrait d'étudier comment les revenus varient d'un magasin à l'autre.

Prenons un exemple simple. Si les ventes totales de mangues dans un mois donné sont représentées par  $Y$ , le prix par mangue est une constante  $c$ , et le nombre de mangues vendues est représenté par  $T$ , les revenus générés par la vente de  $Y$  mangues peuvent être exprimés par une équation simple. Cette dépendance directe entre le revenu et le nombre de mangues vendues montre que plus vous vendez de mangues, plus vous générez de revenus. Mais les questions commerciales sont souvent plus complexes. Par exemple, une entreprise pourrait demander à ses analystes de déterminer la quantité de chaque article à stocker dans chaque entrepôt.

En mathématiques, ces dépendances peuvent être représentées par différentes équations, certaines ayant des interdépendances. Bien que résoudre une seule équation avec une variable connue et une variable inconnue soit une tâche triviale, lorsqu'il y a plusieurs variables et équations, cela devient plus compliqué. C'est là que l'algèbre linéaire intervient.

### **Exemple d'application de l'algèbre linéaire :**

Dans notre exemple de mangues, la question de la rentabilité peut être analysée à l'aide de matrices représentant les équations liées à la vente des mangues. Manipuler ces matrices permet de résoudre des systèmes d'équations complexes, un processus essentiel dans l'analyse prédictive.

## 6.2. La régression

Jusqu'à présent, nous avons vu comment les problèmes réels peuvent être exprimés sous forme d'équations mathématiques et comment la compréhension des statistiques et de l'algèbre linéaire permet de résoudre ces équations et de tirer des conclusions sur le monde réel. Poussons un peu plus loin et essayons de comprendre comment prédire les valeurs futures à partir de données historiques.

Supposons que nous essayons de prédire la valeur de quelque chose que nous allons appeler  $Y$ , la variable dépendante.  $Y$  pourrait être n'importe quoi, du prix d'une action à un moment donné au nombre de verres de limonade vendus par un stand un jour donné. Introduisons également une autre variable, que nous appellerons  $X$ , la variable indépendante.  $X$  pourrait être quelque chose comme le jour de la semaine ou la température moyenne durant la journée.

Prenons un exemple où les variables sont représentées dans un tableau.

Nous pouvons tracer ces valeurs sous forme de graphique. À première vue, plusieurs observations peuvent être faites. La valeur de  $Y$  augmente avec la valeur de  $X$ , et l'augmentation de  $Y$  est proportionnelle à l'augmentation de  $X$ . Si quelqu'un nous demandait de prédire la valeur de  $Y$  lorsque  $X$  est égal à 40, nous devrions extrapoler les données pour deviner la valeur de  $Y$  à partir de cette extrapolation.

Pour rendre cela un peu plus scientifique, plutôt que de nous fier simplement à notre observation, nous pouvons tracer une ligne passant par nos observations historiques et prédire les valeurs de  $Y$  basées sur cette ligne. Cette ligne peut être représentée par une équation linéaire.

Cependant, si plusieurs lignes semblent correspondre aux observations, il devient nécessaire de choisir celle qui convient le mieux pour nos prédictions. Ce processus statistique d'analyse de la relation entre une variable dépendante ( $Y$ ) et une ou plusieurs variables indépendantes ( $X_1, X_2, \dots, X_n$ ) s'appelle la régression.

### 6.2.1. Analyse de la régression

L'utilisation de la régression pour la prévision (ou la prédiction) s'appelle l'analyse de régression. Bien que nous ayons utilisé un exemple simple, dans la réalité, les ensembles de

données sont bien plus grands et l'analyse de régression est souvent effectuée par un programme informatique.

La première question qui se pose est de savoir comment choisir la meilleure ligne pour notre analyse. Les scientifiques utilisent une fonction de perte pour évaluer l'exactitude du modèle de prédiction. Cette fonction de perte calcule la différence entre les valeurs observées ( $Y$ ) et les valeurs prédites ( $Y'$ ). Il existe plusieurs façons d'évaluer cette fonction de perte.

Une méthode courante est l'erreur quadratique moyenne (MSE), qui consiste à calculer la différence au carré entre les valeurs observées et les valeurs prédites, puis à les additionner, et enfin à en prendre la moyenne. L'idée est que plus la différence entre les valeurs observées et prédites est grande, plus la valeur moyenne sera élevée.

Une autre méthode est l'erreur absolue moyenne (MAE), où la différence entre les valeurs observées et prédites est prise en valeur absolue. La différence n'est pas pénalisée davantage si elle est grande, contrairement à MSE.

Les valeurs aberrantes (outliers), qui sont significativement différentes de la valeur prédite, peuvent poser problème car elles risquent d'augmenter la fonction de perte et de fausser les résultats.

### 6.2.2. Types de régression

Plusieurs types de régression sont utilisés pour définir des modèles de prédiction. Examinons ici la régression linéaire, polynomiale et logistique, ainsi que leurs limitations et leurs domaines d'application.

#### **Régression linéaire**

La régression linéaire est la technique de modélisation la plus courante. Elle est utilisée pour comprendre la relation entre une variable dépendante ( $Y$ ) et une ou plusieurs variables indépendantes ( $X_1, X_2$ , etc.). Une régression linéaire simple examine la relation entre  $Y$  et une seule variable indépendante.

Lorsque plusieurs variables indépendantes sont impliquées, cela devient une régression linéaire multiple. L'objectif est de trouver une ligne qui représente au mieux les données observées. Cette ligne peut être modélisée par une équation dans laquelle les valeurs prédites

sont exprimées comme une combinaison linéaire des variables indépendantes, avec des coefficients à déterminer.

Quelques points importants à retenir sur la régression linéaire :

- La variable dépendante peut être continue (comme la taille ou la température) ou discrète (par exemple, la présence ou l'absence d'un événement).
- Il doit y avoir une relation linéaire entre les variables dépendantes et indépendantes.
- La régression linéaire est sensible aux valeurs aberrantes qui peuvent influencer la ligne de prédiction.
- Il ne doit pas y avoir d'autocorrélation entre les observations, c'est-à-dire que les valeurs proches dans le temps doivent être plus similaires que celles éloignées dans le temps.

### **Régression polynomiale**

La régression polynomiale est utilisée lorsque la relation entre les variables dépendantes et indépendantes n'est pas linéaire. Elle représente la relation sous forme de courbe plutôt que de droite. Cette approche est utile lorsque les données montrent une tendance non linéaire, et elle utilise des puissances supérieures de la variable indépendante.

La régression polynomiale peut aussi être multivariée, c'est-à-dire impliquant plusieurs variables indépendantes. Cependant, il est important de noter qu'une régression polynomiale de trop haut degré peut entraîner un surajustement, où le modèle s'adapte trop aux données d'entraînement et devient moins performant sur de nouvelles données.

### **Régression logistique**

La régression logistique est couramment utilisée pour résoudre des problèmes de classification, en particulier lorsque la sortie est binaire (0 ou 1). Par exemple, elle peut être utilisée pour déterminer si un e-mail est un spam ou non, ou pour prédire si une tumeur est bénigne ou maligne. La régression logistique peut aussi être étendue à des cas où la sortie comporte plusieurs catégories, mais sans ordre particulier, comme dans le cas de la classification des préférences des électeurs dans une élection avec plusieurs candidats.

### 6.3. Les arbres de décision

Les arbres de décision sont une méthode d'analyse prédictive populaire qui permet de prendre des décisions sur la base de critères multiples. Ils représentent une **analyse multivariée** simple mais puissante, idéale pour des tâches de classification et de régression. L'avantage des arbres de décision réside dans leur capacité à modéliser des relations complexes tout en étant facilement interprétables, ce qui les rend très utiles dans la prise de décision basée sur des données.

Les arbres de décision sont particulièrement utiles dans des situations où il est nécessaire de prédire une **variable cible** (dépendante) en fonction de plusieurs **variables indépendantes**. Par exemple, ils peuvent être utilisés pour prédire si une personne est à risque élevé ou faible en fonction de variables telles que son salaire annuel et son solde moyen de compte bancaire.

#### Principe de fonctionnement

Un arbre de décision fonctionne en divisant un jeu de données en plusieurs sous-ensembles en utilisant des règles de décision basées sur des valeurs des variables indépendantes. Ces règles sont généralement des questions simples, comme « Est-ce que le salaire annuel est supérieur à 70 000 \$ ? ». Chaque "branche" de l'arbre représente une décision basée sur un critère spécifique, et chaque "feuille" (ou nœud terminal) représente un résultat final (par exemple, une prédiction de faible ou de haut risque).

Prenons l'exemple d'un problème de classification, où l'objectif est de prédire le risque d'un candidat à un prêt en fonction de son salaire annuel et de son solde moyen sur six mois. Le tableau suivant montre des données historiques de candidats à des prêts :

Salary (\$)	Average account balance (\$)	Risk classification
60,000	10,000	High risk
45,000	100,000	Low risk
80,000	2,000	Low risk
100,000	20,000	Low risk
35,000	4,000	High risk
75,000	30,000	Low risk
110,000	10,000	Low risk
68,000	3,000	High risk
12,000	40,000	High risk
40,000	10,000	High risk
42,000	20,000	High risk

Variables indépendantes :

- Le salaire annuel
- Le solde moyen du compte bancaire des six derniers mois

Variable dépendante :

- La classification du risque (Haut risque ou Faible risque)

L'objectif est de construire un arbre de décision qui nous permette de prédire si un candidat à un prêt représente un risque faible ou élevé en fonction de ces deux variables.

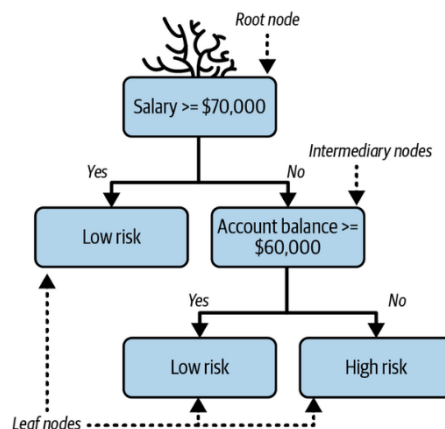
### Exemple de règles générées par l'arbre de décision

Voici un exemple de règles qu'un arbre de décision pourrait générer à partir de cet ensemble de données :

1. Si le salaire annuel est supérieur ou égal à 70 000 \$, le risque est faible.
2. Si le salaire annuel est inférieur à 70 000 \$, alors :
  - o Si le solde moyen du compte est supérieur ou égal à 60 000 \$, le risque est faible.
  - o Si le solde moyen du compte est inférieur à 60 000 \$, le risque est élevé.

Ces règles permettent de classer les nouveaux candidats à un prêt en fonction de leurs caractéristiques. Par exemple :

- Un candidat avec un salaire de 80 000 \$ et un solde moyen de compte de 2 000 \$ serait classé comme faible risque en raison de son salaire élevé.
- Un candidat avec un salaire de 45 000 \$ et un solde moyen de compte de 100 000 \$ serait également classé comme faible risque en raison de son solde élevé.





## **Avantages des arbres de décision**

Simplicité et interprétabilité : les arbres de décision sont facilement compréhensibles et peuvent être visualisés, ce qui permet aux utilisateurs de suivre les règles de décision et d'obtenir des insights explicites.

Flexibilité : ils peuvent être utilisés pour des tâches de classification (comme dans l'exemple ci-dessus) ou de régression (prédiction de valeurs continues).

Aucune nécessité de mise à l'échelle des données : contrairement à d'autres algorithmes de machine learning, les arbres de décision ne nécessitent pas de normalisation des données.

Gestion des données manquantes : les arbres de décision peuvent gérer les données manquantes de manière efficace en faisant une partition des données selon les valeurs disponibles.

## **Limites des arbres de décision**

- Sur-apprentissage (Overfitting) : un arbre de décision trop complexe risque de s'ajuster excessivement aux données d'entraînement, ce qui peut diminuer sa capacité à généraliser sur de nouvelles données.
- Instabilité : les arbres de décision peuvent être sensibles aux petites variations dans les données, ce qui peut entraîner des arbres très différents si les données sont légèrement modifiées.
- Biais vers des variables ayant plus de niveaux : l'algorithme peut être biaisé en faveur des variables ayant plus de catégories (par exemple, des variables avec de nombreuses valeurs uniques), ce qui peut nuire à la qualité du modèle.

## 7. Traitement des données

Avant de commencer toute analyse prédictive, il est essentiel de préparer les données. Cette étape est cruciale pour garantir des résultats fiables, car des données mal préparées peuvent mener à des modèles biaisés ou inefficaces. Le traitement des données inclut plusieurs sous-étapes permettant de nettoyer, transformer et organiser les données de manière appropriée avant d'entamer l'analyse.

### 7.1. Gestion des valeurs manquantes

Les données manquantes sont courantes dans de nombreux jeux de données. Il est important de décider comment les traiter avant de continuer l'analyse :

- Suppression des enregistrements incomplets : si la proportion de valeurs manquantes dans un enregistrement est trop importante, il peut être préférable de le supprimer.
- Remplacement des valeurs manquantes : si l'enregistrement reste utile malgré des valeurs manquantes, on peut les remplacer par des zéros, des moyennes ou des médianes, selon le type de donnée et la situation.
- Utilisation de modèles prédictifs : pour des données manquantes sur des variables spécifiques, des modèles peuvent être utilisés pour estimer et imputer les valeurs manquantes à partir des données disponibles.

### 7.2. Encodage des variables catégorielles

Les variables catégorielles, telles que les types d'habitation ("appartement", "villa"), doivent être transformées en données numériques pour être utilisées par les algorithmes d'apprentissage automatique. Cela peut se faire de différentes manières :

- Encodage par étiquette (Label Encoding) : chaque catégorie est associée à un nombre unique.
- Encodage One-Hot : chaque catégorie est transformée en une colonne binaire, où 1 indique la présence de la catégorie et 0 son absence.

### 7.3. Transformation des données

Certaines variables peuvent avoir des échelles de valeurs très larges, ce qui peut causer des problèmes dans certains modèles d'apprentissage automatique. Pour résoudre cela, il est nécessaire de normaliser ou de standardiser les données :

- Normalisation : mettre à l'échelle les données de manière qu'elles aient une plage de valeurs commune, souvent entre 0 et 1, ce qui est particulièrement utile lorsque les variables sont exprimées dans des unités différentes (par exemple, des millions pour une variable et des milliers pour une autre).
- Standardisation : transformer les données pour qu'elles aient une moyenne nulle et un écart-type de 1, utile pour les modèles qui reposent sur des distances (par exemple, les k plus proches voisins).

#### 7.4. Gestion des valeurs aberrantes (outliers)

Les valeurs aberrantes peuvent fausser les résultats des modèles. Il est crucial de les identifier et de les gérer correctement :

- Suppression des valeurs aberrantes : si une valeur est extrême et clairement incorrecte, elle peut être supprimée.
- Limites calculées : dans certains cas, les valeurs aberrantes peuvent être remplacées par une limite supérieure ou inférieure calculée sur la base des données existantes.
- Utilisation de méthodes robustes : certaines méthodes de modélisation sont plus robustes aux valeurs aberrantes (par exemple, les arbres de décision).

#### 7.5. Gestion du déséquilibre des classes

Un déséquilibre entre les classes dans les données peut entraîner des modèles biaisés, favorisant les classes majoritaires. Plusieurs approches peuvent être utilisées pour traiter ce déséquilibre :

- Pondération des classes minoritaires : accorder plus de poids aux classes sous-représentées lors de l'entraînement du modèle.
- Oversampling : générer des exemples synthétiques pour la classe minoritaire à l'aide de techniques comme SMOTE (Synthetic Minority Over-sampling Technique).
- Undersampling : réduire la taille des classes majoritaires pour équilibrer les proportions.

#### 7.6. Combinaison de données

Les données peuvent provenir de différentes sources, ce qui nécessite parfois des opérations complexes pour les combiner de manière cohérente :

- Joins dans des bases de données relationnelles : lorsque les données sont stockées dans différentes tables, des opérations de jointure permettent de les combiner en un seul jeu de données.
- Pipelines ETL (Extract, Transform, Load) : dans les systèmes de big data, des pipelines ETL permettent d'extraire les données de différentes sources, de les transformer selon les besoins (nettoyage, agrégation) et de les charger dans un entrepôt de données pour analyse.

### 7.7. Sélection des caractéristiques

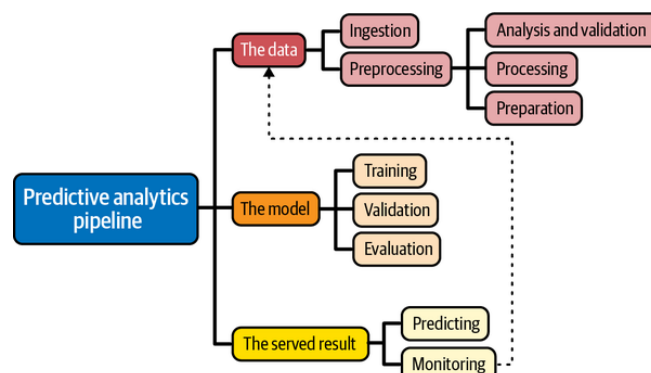
Il est important d'éliminer les variables inutiles ou redondantes pour réduire la complexité du modèle et améliorer sa performance. La sélection des caractéristiques permet de choisir les variables les plus pertinentes pour la prédiction.

- Méthodes statistiques : utilisation de tests de corrélation ou de techniques comme la régression pour sélectionner les variables significatives.
- Méthodes d'élimination automatique : certaines techniques comme le Lasso ou les arbres de décision permettent d'éliminer automatiquement les caractéristiques moins importantes.

### 7.8. Division des données

Avant d'entraîner un modèle, il est nécessaire de diviser les données en trois ensembles distincts : l'entraînement, la validation et le test.

- Ensemble d'entraînement : il est utilisé pour ajuster les paramètres du modèle.
- Ensemble de validation : il permet d'évaluer la performance du modèle pendant l'entraînement et d'ajuster les hyperparamètres.
- Ensemble de test : il est réservé pour évaluer la performance finale du modèle une fois l'entraînement terminé.



## 8. Python et scikit-learn pour les analyses prédictives

### 8.1. NumPy : manipulation efficace de données numériques

NumPy est une bibliothèque Python essentielle pour la manipulation des données numériques. Elle fournit des objets de type tableau multidimensionnel, appelés ndarray, qui permettent de stocker et de manipuler efficacement des données. Ces tableaux sont plus rapides que les listes Python classiques, notamment pour les grandes quantités de données.

Les principales fonctionnalités de NumPy incluent :

- Opérations mathématiques : permet d'effectuer des calculs vectoriels, des opérations arithmétiques, des transformations mathématiques sur des tableaux, etc.
- Opérations statistiques : calculs de moyennes, médianes, écart-types, corrélations, etc.
- Tri et sélection : offre des méthodes pour trier des tableaux, sélectionner des sous-ensembles spécifiques ou manipuler les indices.
- Gestion de matrices : facilitée par les opérations sur des matrices et les fonctions linéaires, utiles dans de nombreux algorithmes de machine learning.



### 8.2. Pandas : gestion et préparation des données

Pandas est une bibliothèque Python open-source qui permet de travailler de manière flexible avec des données structurées. Elle est particulièrement utile en data science et data analytics pour la manipulation de données tabulaires.

Les principales fonctionnalités de Pandas comprennent :

- Lecture et écriture de données : pandas permet de lire et d'écrire des données dans divers formats, tels que CSV, Excel, SQL, et bien d'autres, facilitant ainsi l'importation et l'exportation de jeux de données.

- Nettoyage des données : inclut des fonctions pour gérer les valeurs manquantes, supprimer les doublons, ou convertir des types de données. C'est une étape clé avant d'appliquer des algorithmes de machine learning.
- Manipulation de données : pandas propose des méthodes puissantes pour redimensionner les jeux de données, comme le "pivoting" ou le "melting", qui permettent de réorganiser les colonnes et lignes des DataFrames selon les besoins de l'analyse.
- Indexation et découpage : vous pouvez facilement sélectionner et manipuler des sous-ensembles de données avec des techniques d'indexation avancée, comme les filtres sur les colonnes et les lignes, et la gestion de données temporelles.
- Fusion de données : pandas permet de combiner plusieurs jeux de données à l'aide des fonctions merge et join, utiles pour l'intégration de données provenant de sources diverses.

Pandas est souvent utilisé en combinaison avec d'autres bibliothèques de visualisation de données, comme Matplotlib, pour une représentation graphique des résultats d'analyse.



### 8.3. Matplotlib : visualisation des données

Matplotlib est une bibliothèque Python utilisée pour créer des graphiques et des visualisations. Elle est particulièrement utile pour représenter visuellement les données contenues dans un DataFrame Pandas, rendant l'analyse plus accessible et intuitive.

Avec Matplotlib, vous pouvez créer une variété de graphiques :

- Graphiques linéaires : Pour visualiser des tendances ou des évolutions au fil du temps.
- Histogrammes : Pour examiner la distribution des données.
- Graphiques en barres : Utiles pour comparer des valeurs discrètes ou catégorielles.

Les graphiques générés par Matplotlib permettent de mieux comprendre les relations entre différentes variables et de détecter des motifs ou anomalies dans les données.

#### 8.4. Scikit-learn : modélisation et prédiction

Scikit-learn est l'une des bibliothèques les plus utilisées pour le machine learning en Python. Elle propose une large gamme de modèles d'apprentissage automatique, adaptés à des tâches de régression, classification, clustering, et réduction de dimensionnalité.

Les fonctionnalités clés de Scikit-learn incluent :

- Modèles d'apprentissage automatique : Scikit-learn offre des implémentations de nombreux algorithmes populaires, tels que les forêts aléatoires, les régressions linéaires et logistiques, les k-plus proches voisins (KNN), les machines à vecteurs de support (SVM), et les réseaux de neurones.
- Pipelines de machine learning : cette fonctionnalité permet de créer des chaînes d'opérations pour automatiser et structurer le flux de travail d'un projet de machine learning. Par exemple, les étapes de nettoyage, transformation, et modélisation peuvent être regroupées dans un pipeline pour éviter les erreurs de réutilisation des données.
- Évaluation des modèles : Scikit-learn offre de nombreuses métriques d'évaluation (précision, rappel, F1-score, etc.) et des outils pour la validation croisée, permettant de tester la robustesse des modèles avant leur déploiement.

Grâce à ses outils simples et efficaces, Scikit-learn est un choix privilégié pour ceux qui souhaitent créer des modèles de machine learning sans complexité excessive. Son intégration avec d'autres bibliothèques comme NumPy et Pandas facilite la mise en œuvre de solutions prédictives.



## 8.5. Entraînement et prédiction avec un modèle de régression linéaire

**Il y a un code qui est annexé à cette lecture individuelle, voici son explication :**

Le processus commence avec le jeu de données des manchots, qui présente des informations physiques sur trois espèces de manchots de l'Antarctique. Pour utiliser ce jeu de données, nous devons d'abord le nettoyer en supprimant les enregistrements comportant des champs vides, notamment pour la masse corporelle et le sexe des manchots.

Une fois le jeu de données nettoyé, il est divisé en deux sous-ensembles : un ensemble d'entraînement (75 % des données) et un ensemble de test (25 %). Cette division est réalisée à l'aide de la fonction `train_test_split` de Scikit-learn. Pour plus de propreté, des copies des sous-ensembles sont créées au lieu de simplement les référencer.

Un modèle de régression linéaire est initialisé, avec la longueur des nageoires comme variable indépendante et la masse corporelle comme variable dépendante. Le modèle est ensuite entraîné avec l'ensemble d'entraînement en utilisant la méthode `fit`.

Une fois le modèle formé, nous utilisons l'ensemble de test pour faire des prédictions sur la masse corporelle des manchots. Ces prédictions sont comparées aux valeurs réelles, et une évaluation visuelle est effectuée en traçant une ligne de régression sur un graphique de dispersion des données de test.

Le modèle est évalué à l'aide de la fonction `score`, qui calcule le coefficient de détermination  $R^2$  pour mesurer la performance du modèle. La précision sur les données d'entraînement est généralement meilleure que sur les données de test.

Pour évaluer davantage le modèle, nous utilisons la validation croisée, qui permet d'obtenir des scores pour plusieurs itérations de séparation des données, et ce, à l'aide de la méthode `cross_val_score` de Scikit-learn. Cela aide à vérifier la stabilité du modèle sur différentes partitions des données.

## 8.6. Entraînement et prédiction avec un modèle d'arbre de décision

**Il y a un code qui est annexé à cette lecture individuelle, voici son explication :**

Pour aborder ce problème de classification, nous avons opté pour un **arbre de décision**. Un arbre de décision est un modèle de machine learning qui divise de manière récursive les données en sous-groupes homogènes basés sur des critères de décision. Ce modèle est



particulièrement adapté aux problèmes de classification, car il permet de diviser les données en fonction de valeurs de caractéristiques, et de prédire une catégorie pour chaque observation.

### **Préparation des données**

Avant d'entraîner un arbre de décision, il est nécessaire de préparer les données. Nous avons utilisé un jeu de données contenant différentes caractéristiques des manchots, telles que l'espèce, la longueur et la profondeur du bec, et la masse corporelle. Pour que le modèle puisse utiliser ces informations, nous avons transformé les variables catégorielles, comme l'espèce et l'île, en variables numériques via un processus appelé **encodage**. Cela permet au modèle de comprendre et d'interpréter les valeurs textuelles de manière quantitative.

### **Séparation des données**

Une fois les données préparées, nous avons divisé le jeu de données en deux ensembles : un ensemble d'entraînement et un ensemble de test. L'ensemble d'entraînement est utilisé pour apprendre le modèle, tandis que l'ensemble de test sert à évaluer sa capacité à généraliser sur des données qu'il n'a pas vues auparavant. En règle générale, environ 80 % des données sont utilisées pour l'entraînement et 20 % pour le test.

### **Entraînement de l'arbre de décision**

L'étape suivante consiste à entraîner l'arbre de décision. En gros, l'arbre de décision apprend à prédire le sexe des manchots en fonction des caractéristiques disponibles. Le modèle commence par examiner la caractéristique qui sépare le mieux les données, puis continue à diviser les données en sous-groupes selon ce critère. Ce processus est itéré jusqu'à ce que le modèle ait créé un arbre de décision complet, où chaque feuille correspond à une prédiction (dans notre cas, le sexe d'un manchot).

### **Évaluation des performances**

Une fois l'arbre de décision entraîné, il est important d'évaluer sa performance. Nous utilisons l'ensemble de test pour voir à quel point le modèle peut prédire correctement le sexe des manchots en fonction des données qu'il n'a pas vues. L'une des métriques les plus simples pour évaluer un modèle de classification est le taux de précision, qui mesure la proportion de prédictions correctes parmi toutes les prédictions faites par le modèle.

## Visualisation de l'arbre de décision

Une des forces de l'arbre de décision est sa transparence. Contrairement à des modèles plus complexes comme les réseaux de neurones, l'arbre de décision peut être visualisé sous forme d'un graphique qui montre les règles de décision et les divisions effectuées par le modèle. Cette visualisation permet de comprendre comment le modèle prend ses décisions, ce qui peut être très utile pour l'interprétation et la validation des résultats.

## Optimisation du modèle

Bien que les arbres de décision soient relativement simples à comprendre, ils peuvent parfois être trop complexes et entraîner un **surapprentissage** (overfitting), où le modèle apprend trop de détails sur les données d'entraînement et échoue à généraliser correctement sur de nouvelles données. Pour éviter cela, il est possible de limiter la profondeur de l'arbre, ce qui le rend plus simple et plus général. En ajustant ce paramètre, nous avons pu améliorer la performance de notre modèle sur l'ensemble de test, tout en conservant une bonne précision.

## 9. TensorFlow et Keras

TensorFlow est une bibliothèque open-source développée par Google, principalement utilisée pour le machine learning (ML) et le deep learning. Initialement lancée en 2015, cette plateforme end-to-end permet de concevoir, entraîner et déployer des modèles de machine learning à grande échelle. TensorFlow est particulièrement adapté pour les modèles de deep learning, mais il peut également être utilisé pour des tâches plus classiques en machine learning. La plateforme est capable de travailler avec des réseaux neuronaux, des modèles de régression, des modèles de classification, et bien plus encore. Son principal avantage réside dans sa capacité à gérer des données massives et à effectuer des calculs sur des infrastructures distribuées, ce qui en fait un choix privilégié pour les chercheurs et les entreprises souhaitant développer des solutions d'intelligence artificielle de manière scalable.

TensorFlow supporte une exécution sur plusieurs plateformes, y compris sur des CPU, GPU et TPU (Tensor Processing Units, des unités spécialisées créées par Google), ce qui améliore considérablement les performances d'entraînement des modèles. La plateforme offre aussi un ensemble d'outils et de bibliothèques complémentaires permettant de réaliser des tâches variées, telles que la manipulation de données avec TensorFlow Data, la gestion des modèles avec TensorFlow Serving, ou encore la création d'interfaces d'applications avec TensorFlow.js.

### **Tenseurs dans TensorFlow**

Dans TensorFlow, les tenseurs sont des structures de données multidimensionnelles, qui représentent des tableaux de  $n$  dimensions. Les tenseurs sont essentiels car ils permettent de stocker et manipuler des données, qu'elles soient unidimensionnelles (vecteurs), bidimensionnelles (matrices) ou multidimensionnelles. En d'autres termes, un tenseur est une généralisation des matrices à des dimensions plus élevées. Par exemple :

- Un scalaire est un tenseur de zéro dimension.
- Un vecteur est un tenseur à une dimension.
- Une matrice est un tenseur à deux dimensions.

Et ainsi de suite pour des tenseurs à trois dimensions et plus.

Ces tenseurs permettent de faire des calculs mathématiques complexes tels que la multiplication scalaire, la multiplication matricielle, les produits de convolution, et bien d'autres opérations.

### **Keras : Une API simplifiée pour le Deep Learning**

Keras est une API de haut niveau pour le deep learning, initialement développée par François Chollet et ensuite intégrée à TensorFlow. Keras a été conçu pour être facile à utiliser, modulaire, et extensible, permettant ainsi aux chercheurs et aux développeurs de se concentrer sur les aspects innovants de leurs modèles, plutôt que sur la gestion des détails techniques. Keras est aujourd'hui un composant intégré de TensorFlow, ce qui signifie que tous les utilisateurs de TensorFlow peuvent utiliser Keras directement pour créer, entraîner et évaluer leurs modèles de deep learning.

Keras offre une interface simple pour créer des réseaux neuronaux, que ce soit pour des architectures de réseaux de neurones classiques (comme les perceptrons multicouches) ou pour des réseaux neuronaux plus complexes (comme les réseaux convolutifs - CNNs - ou les réseaux récurrents - RNNs). L'API Keras permet de :

- Construire des modèles de manière déclarative grâce à des blocs de construction simples et des fonctions très intuitives.
- Entraîner des modèles en utilisant des optimisateurs prédéfinis, des fonctions de perte et des métriques.
- Évaluer des performances sur les données d'entraînement et de test.
- Exporter les modèles pour la production en quelques lignes de code.

Keras est particulièrement apprécié pour sa simplicité d'utilisation, ce qui en fait une bibliothèque idéale pour les débutants en deep learning, mais aussi pour les experts qui souhaitent prototyper rapidement des idées de modèles.

### **Applications pratiques de TensorFlow et Keras**

TensorFlow et Keras sont utilisés dans une multitude d'applications dans le domaine de l'intelligence artificielle, telles que :

- Analyse d'image et de vidéo : Grâce aux réseaux neuronaux convolutifs (CNN), TensorFlow permet de résoudre des tâches de classification d'images, de détection d'objets, ou de segmentation d'images.
- Reconnaissance vocale et traitement du langage naturel (NLP) : Avec des architectures comme les réseaux neuronaux récurrents (RNN) ou les transformers, TensorFlow est largement utilisé pour la traduction automatique, la reconnaissance vocale, et l'analyse des sentiments dans des textes.
- Prédiction et recommandation : TensorFlow est aussi très utilisé dans des systèmes de recommandation, tels que ceux utilisés par Netflix ou Amazon, pour suggérer des produits ou des contenus aux utilisateurs.

