

# Projet d'analyse de données

Dans cette partie, on présente les différentes étapes d'un projet d'analyse de données.

## 1 Projet d'analyse données

Un projet d'analyse de données peut se découper en cinq grandes étapes :

1. Définition des objectifs
2. Données
3. Élaboration et validation des modèles
4. Mise en oeuvre
5. Suivi de la performance et amélioration

### Remarque

Dans ce cours, on s'intéressera principalement à l'élaboration et à la validation de modèles.

Lors de la planification d'un projet, il faut prendre en compte que chaque étape à une importance différente, mais aussi que chacune ne prend pas le même temps d'exécution. Pyle (1999) donne une estimation du temps de chaque étape, ainsi que de leur importance dans la réussite du projet (donné en pourcentage du total).

| Étape                   | Temps | Importance |
|-------------------------|-------|------------|
| Comprendre le problème  | 10/   | 15/        |
| Explorer la solution    | 9/    | 14/        |
| Implementer la solution | 1/    | 51/        |
| Préparer les données    | 60/   | 15/        |
| Analyser les données    | 15/   | 3/         |
| Modéliser les données   | 5/    | 2/         |

On remarque deux faits importants : ce n'est pas parce qu'une étape est très importante qu'elle va prendre beaucoup de temps. L'implémentation de la solution est très importante (sinon il n'y a pas de résultat), mais ne sera généralement pas très longue à faire (possiblement en quelques lignes de code). À l'inverse, la préparation des données est un étape d'importance moyenne (encore que c'est discutable), mais elle prend la majeure partie du temps du projet. En effet, il faut, par exemple, gérer les données manquantes, les données aberrantes, les éventuels accents pour des données en français, etc.

## 1.1 Définition des objectifs

Est-ce que l'on veut : visualiser les données ? explorer et émettre des hypothèses ? tester ? regrouper ? comprendre ? prédire ?

Comment fait-on en pratique ? On pose des questions ! Tout d'abord, il faut clarifier les termes. Qui va utiliser le modèle et comment ? Quelle est la population cible ?

### Exemples

1. La Banque National du Canada veut lancer un nouveau produit d'épargne et souhaite mieux connaître ses clients pour prédire s'ils veulent l'acheter.
2. L'équipe de hockey des Canadiens de Montréal souhaite mieux connaître ses adversaires pour développer des nouvelles tactiques de jeu.
3. Pharmascience souhaite savoir si son nouveau médicament est efficace.

## 1.2 Données

- Inventaire et qualité
- Constitution de la base de données
- Exploration et traitement préliminaire

Qu'est-ce que l'on veut dire par qualité des données ?

- Est-ce que les données sont représentatives de la population cible ?
- Est-ce que les données permettront de tirer des conclusions de causalité ?
- Est-ce que les données sont fiables ?

Source de données :

Quelques liens pour récupérer des données.

Nettoyage de données : cf R (importation, nettoyage, tidyverse, types de variables, retirer les doublons, uniformiser les modalités, vérifier le format des valeurs spéciales, pivot, opérateur pipe, jointure).

Exploration des données : modalités rares, modalités trop nombreuses, asymétrie, déséquilibre des classes, valeurs extrêmes ou aberrantes, variables fortement corrélées, valeurs manquantes.

Statistiques descriptives

Pyle, Dorian. 1999. *Data Preparation for Data Mining*. Morgan Kaufmann.