

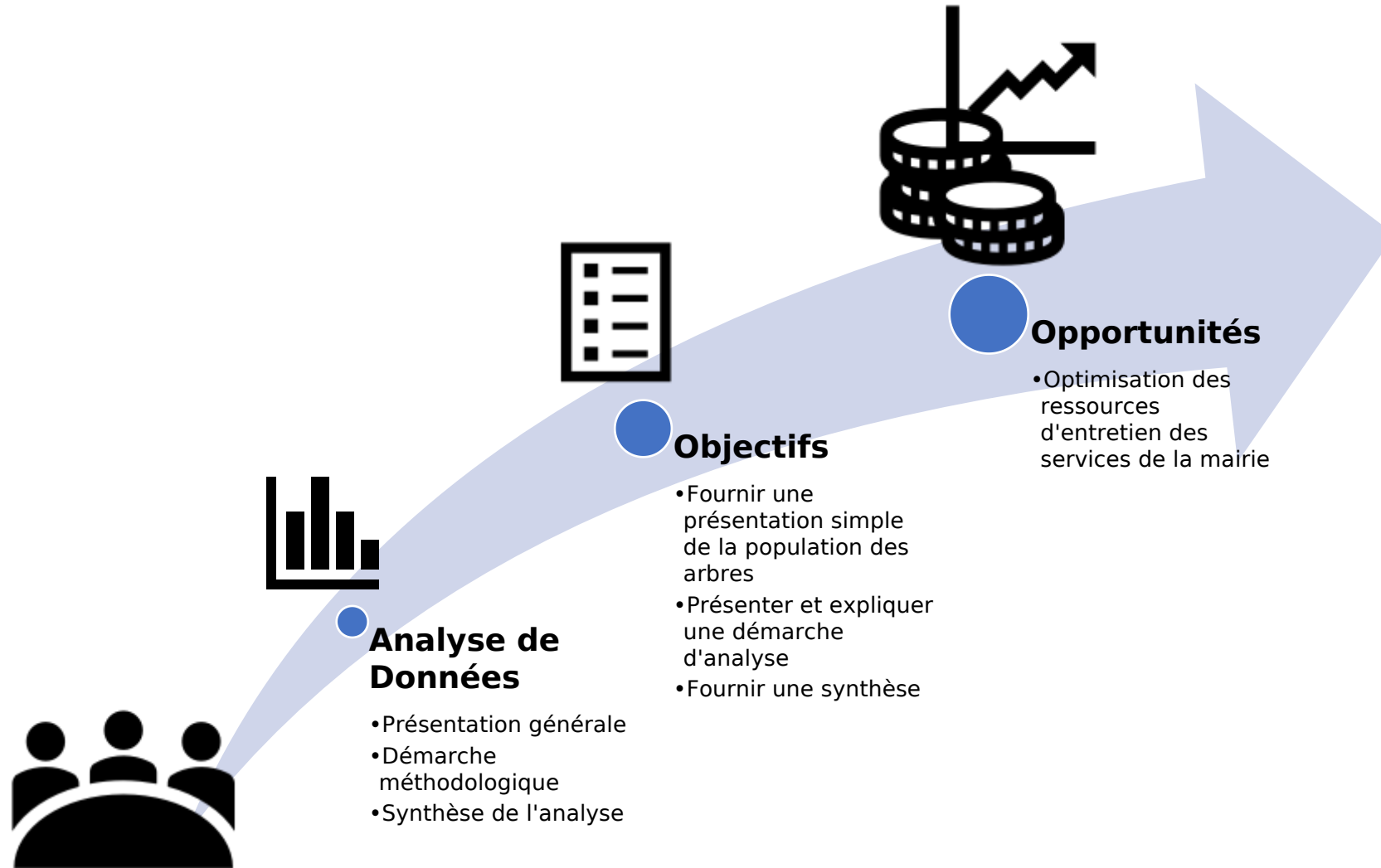


# Projet **Végéталisons la ville**

Analyse des données  
sur les arbres de la ville  
de Paris

Bailly DIOUNOU –  
11/08/2020

# Contexte & périmètre du projet



# Présentation générale du jeu de données

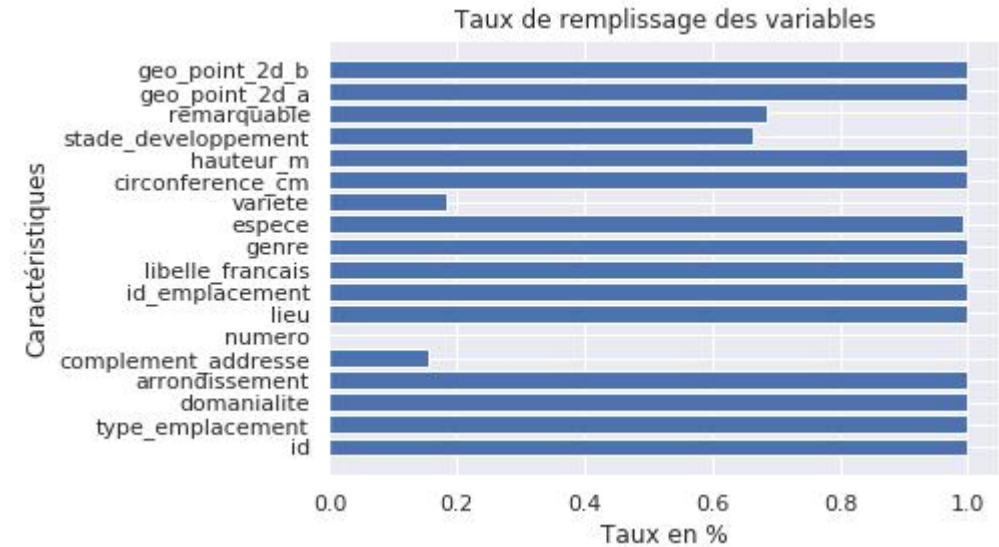
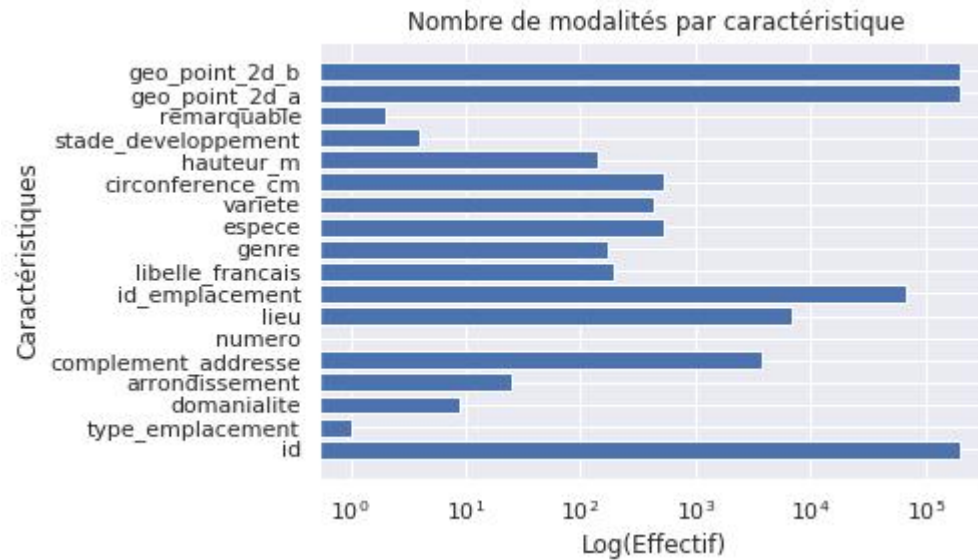
## **Données statistiques de base - Données brutes**

- Taille de la population **N = 200 137** individus
- Nombre de caractéristiques **18 + 1**
- Caractéristiques quantitatives **q = 3** (circ\_cm, hauteur\_m, remarquable toutes *continues*)
- Caractéristiques qualitatives **p = 16** - dont 4 *ordinales*.

# Présentation générale du jeu de données

Indicateurs statistiques de base des variables numériques (quantitatives et qualitatives)							
	id	numero	circonfere nce_cm	hauteur_m	remarquab le	geo_point_ 2d_a	geo_point_ 2d_b
Count	200137	0	200137	200137	137039	200137	200137
Mean	387202,68 23		83,380479 37	13,110509 3	0,0013426 83	48,854490 79	2,3482078 96
Std	545603,24 24		673,19021 3	1971,2173 87	0,0366181 71	0,0302341 45	0,0512196 4
Min	99874		0	0	0	48,742290 3	2,2102412 29
25%	155927		30	5	0	48,835020 74	2,3075303 21
50%	221078		70	8	0	48,854162 09	2,3510951 66
75%	274102		115	12	0	48,876447 14	2,3868380 72
Max	2024745		250255	881818	1	48,911484 76	2,4697594 72

# Présentation générale du jeu de données



# Démarche méthodologique d'analyse de données

**Nettoyage**

**Imputation**

**Analyse Univariée**

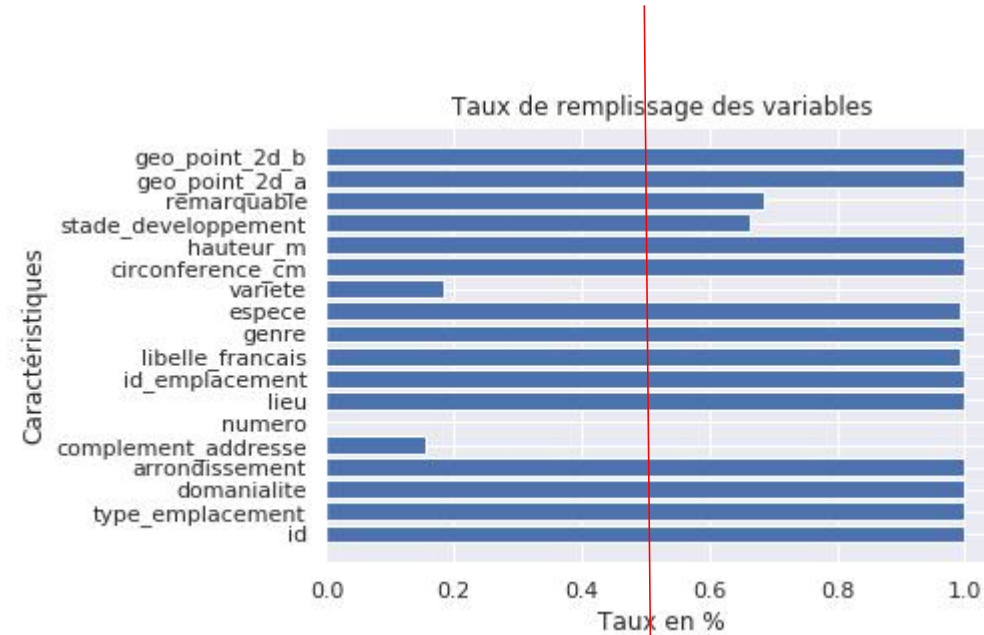
# Démarche méthodologique

## 1. Nettoyage

- **70% de plancher pour le taux de remplissage**

- Pertinence évaluée en fonction de qualité et nombre de modalités

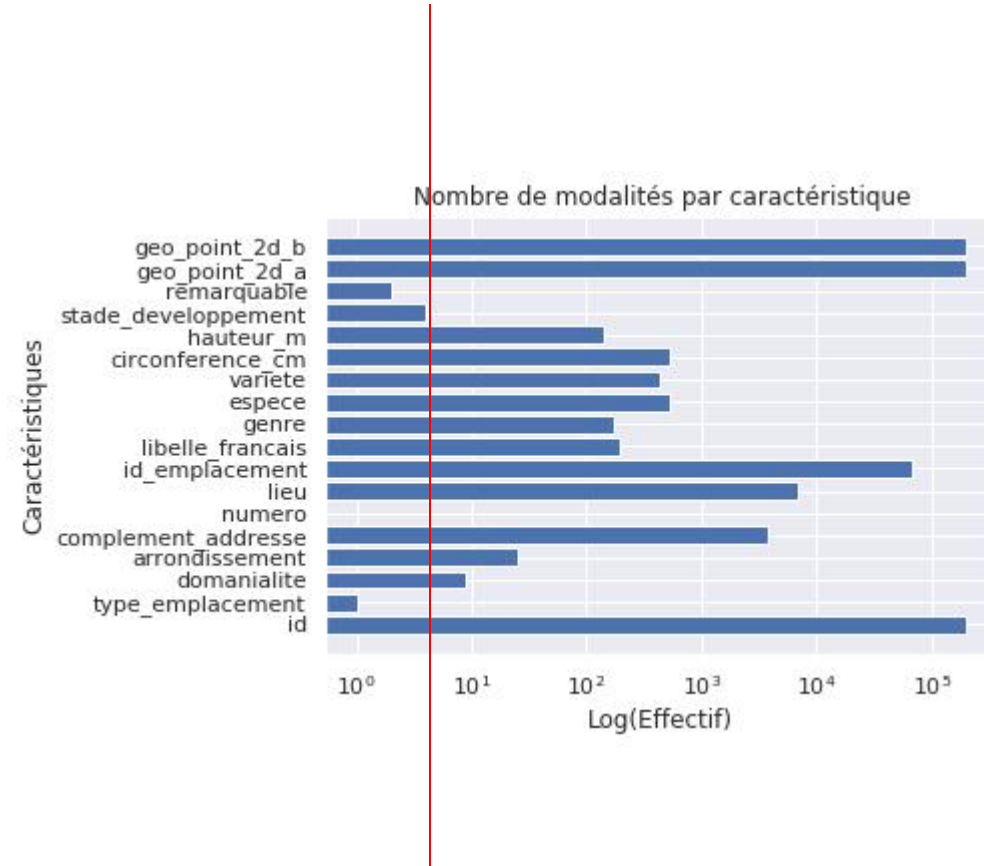
- Filtre par identification des valeurs aberrantes (quantiles)



# Démarche méthodologique

## 1. Nettoyage

- 70% de plancher pour le taux de remplissage
- **Pertinence évaluée en fonction de qualité et nombre de modalités**
- Filtre par identification des valeurs aberrantes (quantiles)

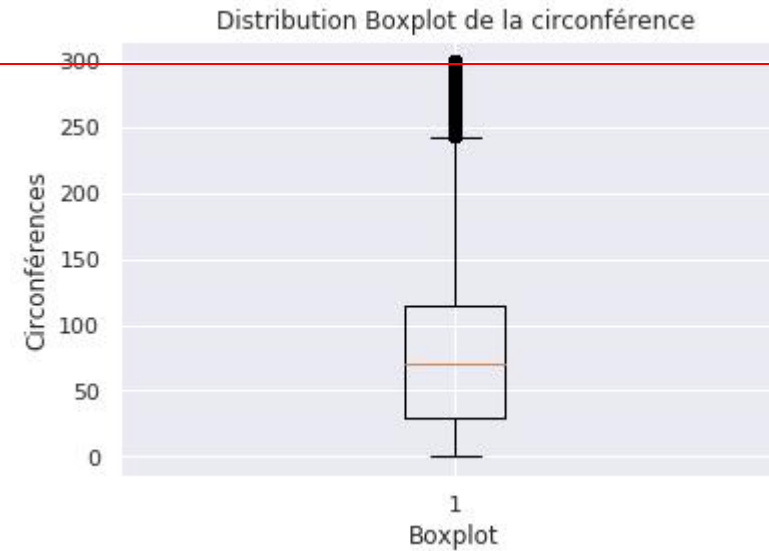




# Démarche méthodologique

## 1. Nettoyage

- 70% de plancher pour le taux de remplissage
- Pertinence évaluée en fonction de qualité et nombre de modalités
- **Filtre par identification des valeurs aberrantes (quantiles)**



# Démarche méthodologique

## 2. *Imputation*

- **Critère dendrologique**

- Imputation basée sur un apprentissage interne vs Imputation simple/itérative
- Gain -> libellé 0.69% / 0.75%, espèce: 0.31 / 0.88

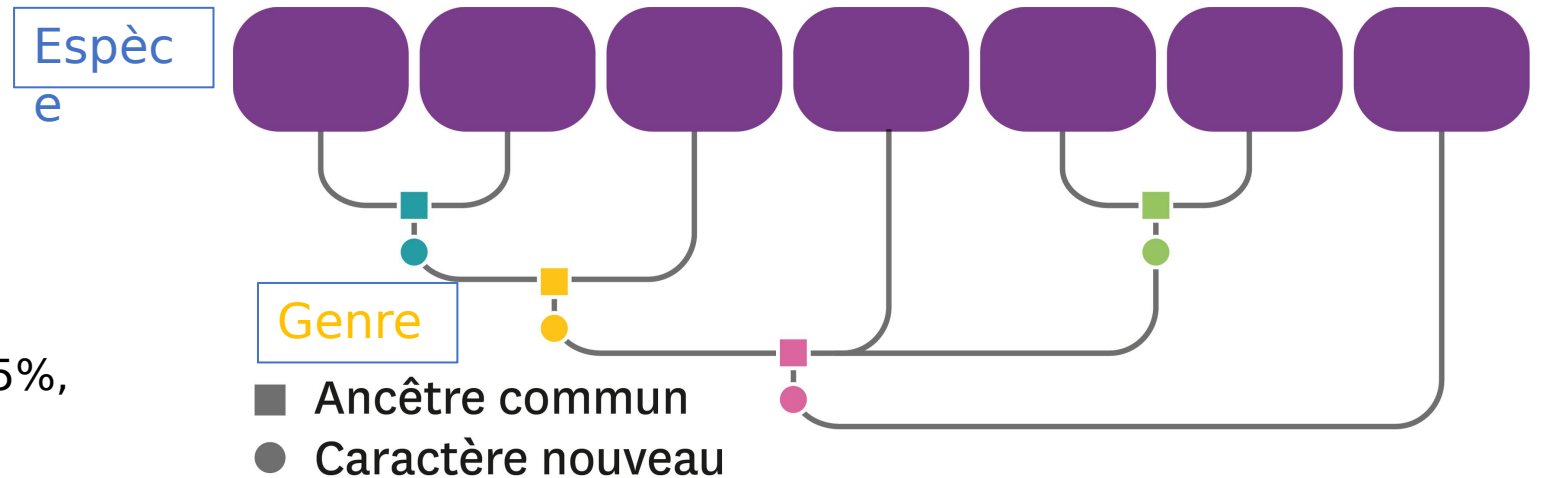
Domaine Ex.: *Bacteria*  
Règne *Procaryotae*  
Phylum  
Classe *Schizomycetes*  
Ordre *Micrococcales*  
Famille *Micrococcaceae*  
Genre *Staphylococcus*  
Espèce *S. aureus*

# Démarche méthodologique

## 2. *Imputation*

- **Critère dendrologique**

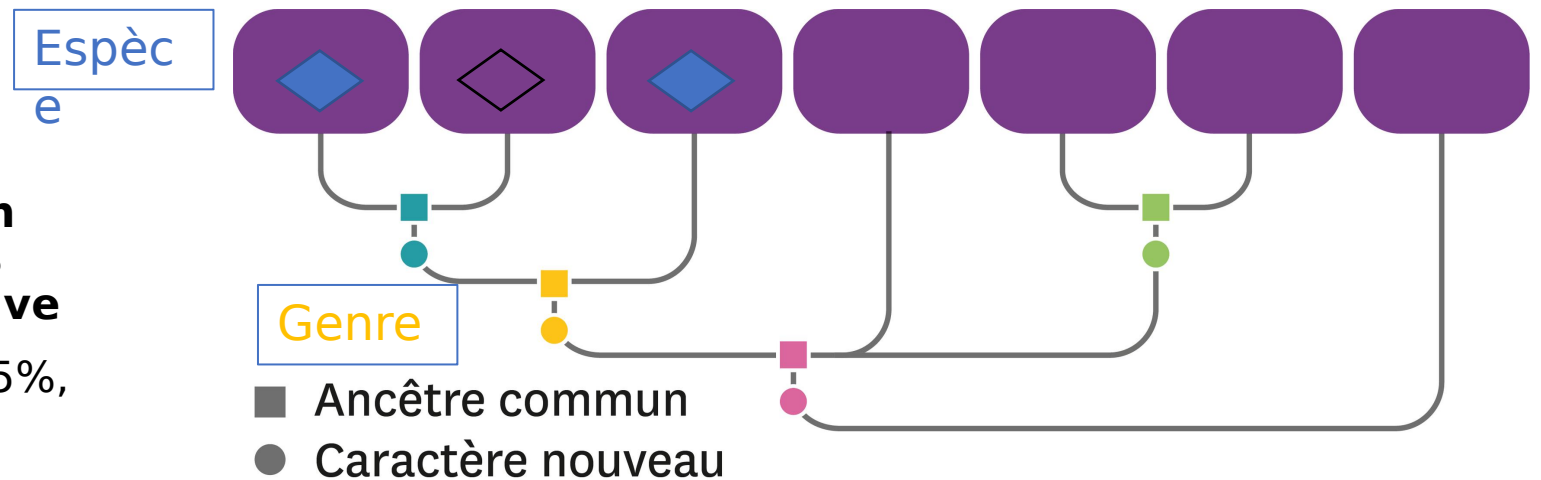
- Imputation basée sur un apprentissage interne vs Imputation simple/itérative
- Gain -> libellé 0.69% / 0.75%, espèce: 0.31 / 0.88



# Démarche méthodologique

## 2. *Imputation*

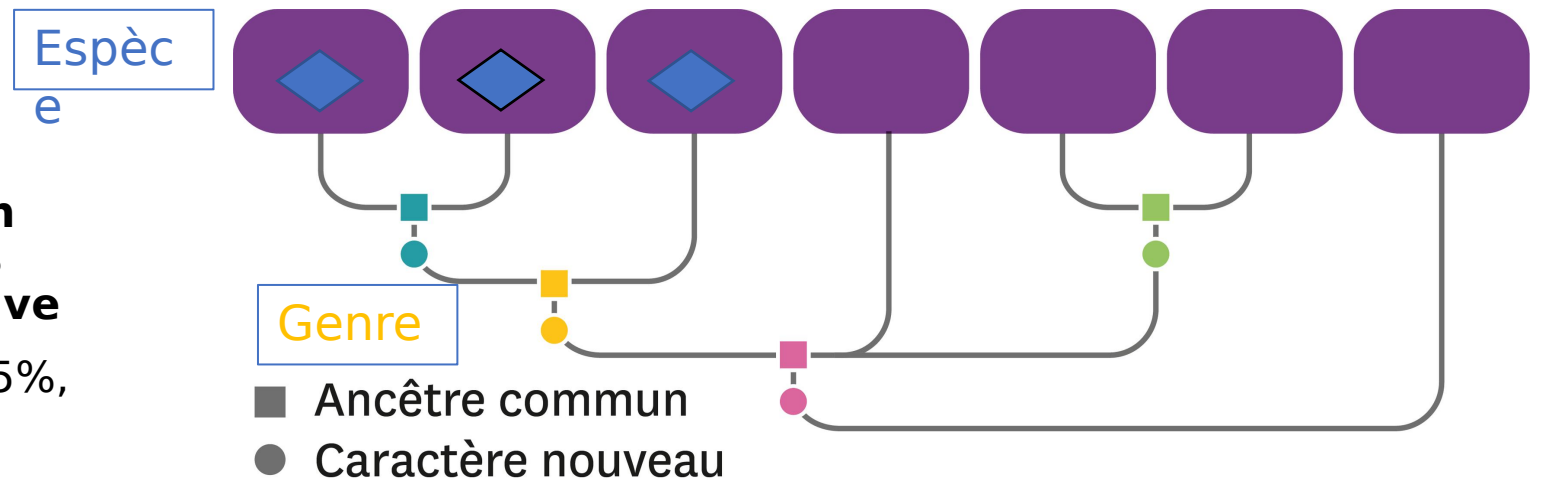
- Critère dendrologique
- **Imputation basée sur un apprentissage interne vs Imputation simple/itérative**
- Gain -> libellé 0.69% / 0.75%,  
espèce: 0.31 / 0.88



# Démarche méthodologique

## 2. *Imputation*

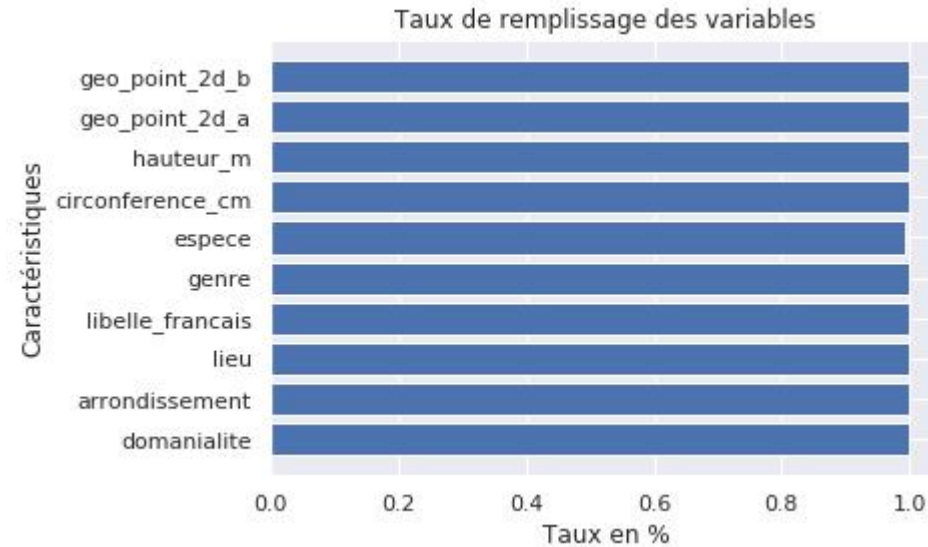
- Critère dendrologique
- **Imputation basée sur un apprentissage interne vs Imputation simple/itérative**
- Gain -> libellé 0.69% / 0.75%, espèce: 0.31 / 0.88



# Démarche méthodologique

## 2. *Imputation*

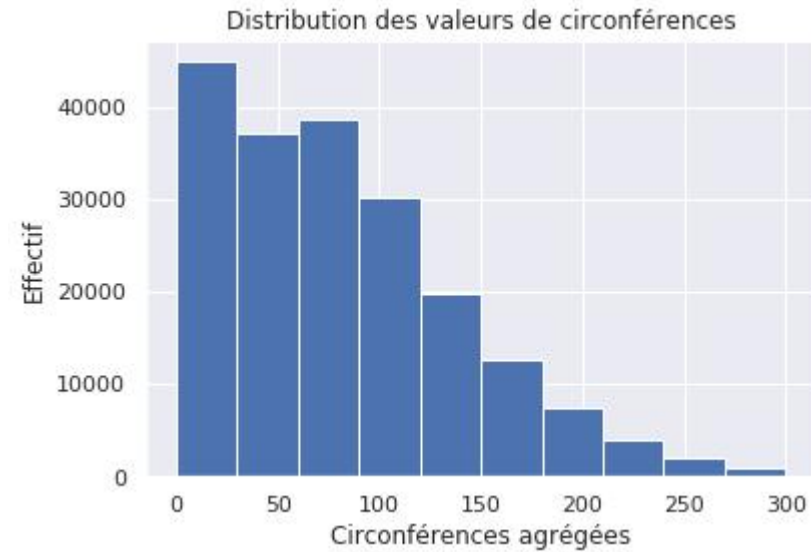
- Critère dendrologique
- Imputation basée sur un apprentissage interne vs Imputation simple/itérative
- **Gain**
  - > **libellé**:  $0.69\% / 0.75\% = 92\%$
  - > **espèce**:  $0.31 / 0.88 = 35\%$



# Démarche méthodologique

## 3. *Analyse Univariée*

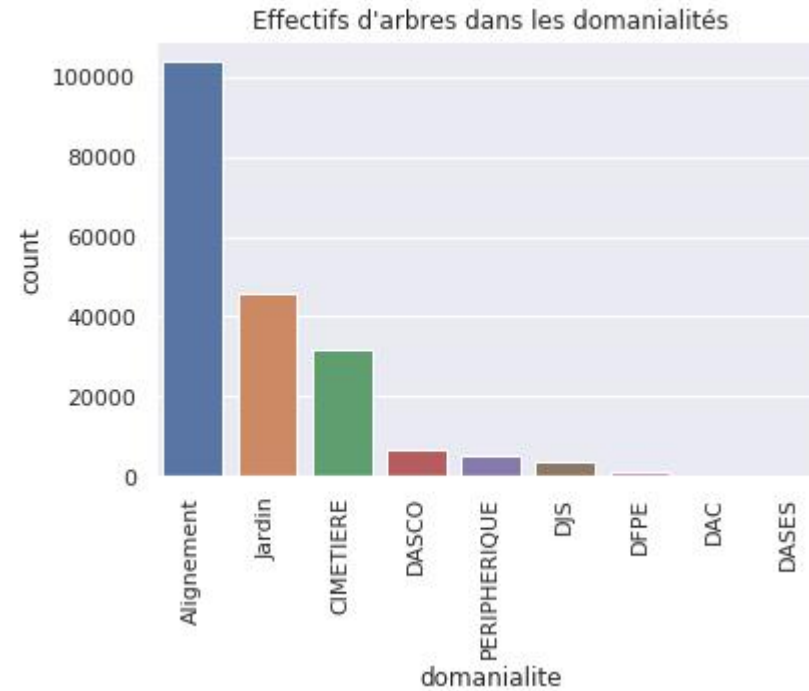
- **Circonférence**
  - Forme quasi gaussienne
- Domanialité
  - Trimodale



# Démarche méthodologique

## 3. *Analyse Univariée*

- Circonférence
  - Forme quasi gaussienne
- **Domanialité**
  - Trimodale



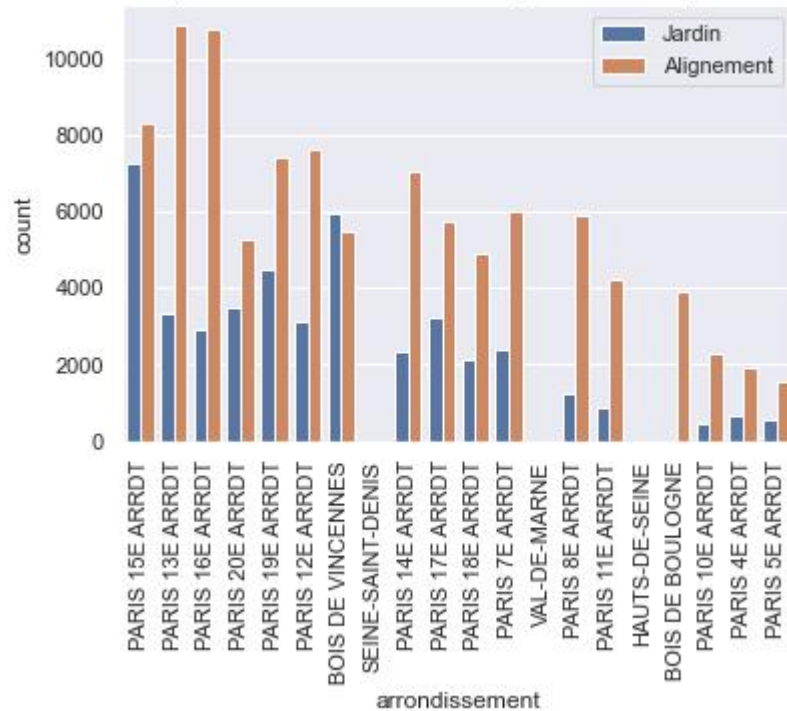


# Démarche méthodologique

## 4. Analyse bivariable & Stratégie

### Domanialité - Arrondissement

Effectifs d'arbres par arrondissement différenciés par domanialité (Jardin vs Alignement)



### Densité d'arbres par arrondissement





# Synthèse de l'analyse de données

## Jeu de données des arbres de la ville de Paris

N = 200 137 individus, 19 caractéristiques (3 quantitatives)

### Démarche méthodologique d'analyse

#### (1) Nettoyage

- Filtre "passe-haut" à taux de remplissage de 70%,
- Filtre de non-pertinence: critère du nombre et de la qualité des modalités
- Filtre sur valeurs aberrantes de dimensions : multiples des quantiles de leurs distributions

#### (2) Imputation

- Critère dendrologique sur relations genre-espèce et genre-libellé français. Valeurs statistique et dendrologique
- Gains sur valeurs manquantes: 92% sur libellé français et 35% sur espèce



**Bilan: Suppression: 13% de lignes, ~50% des colonnes**

#### (3) Analyse univariée et bivariée

- **Analyse univariée** sur la grande majorité des caractéristiques post-nettoyage.
  - Tracé des distributions,
  - Calcul des indicateurs statistique de base (moyenne, médiane, écart-type).
- **Analyse bivariée**
  - Croisement de la domanialité (2 plus fortes modalités: Jardin et Alignement), avec Arrondissement, Libellé français, Espèce et Circonférence.
  - Tracé des distributions différenciées

# Merci pour votre attention

Temps de questions/réponses