

Participez à un concours sur la smartcity

AI Engineer - Projet 2

Ludovic Lafon

December 6, 2023

OpenClassrooms

Exploration des données

Traitement des données manquantes

Détection des outliers

Exploration des données

Nombre de données par colonnes

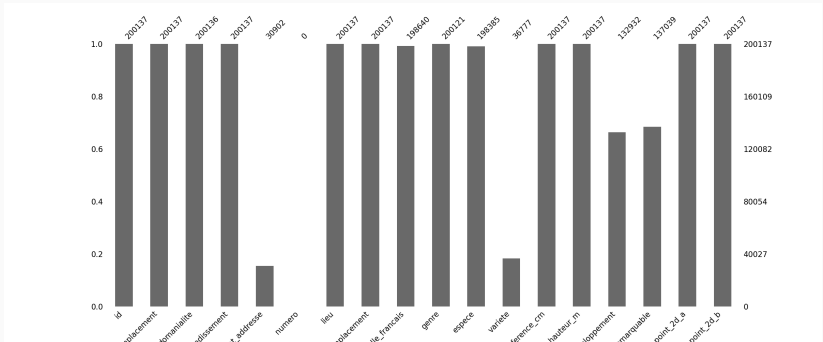


Figure 1: Nombre de données par colonnes

df.shape	
✓ 0.0s	Python
(200137, 18)	

Figure 2: pandas.DataFrame.shape() - 200137 lignes, 18 colonnes

Données Manquantes

Données manquantes :

- domanialite
- complement adresse
- numero
- libelle francais
- genre
- espece
- variete
- stade developpement
- remarquable

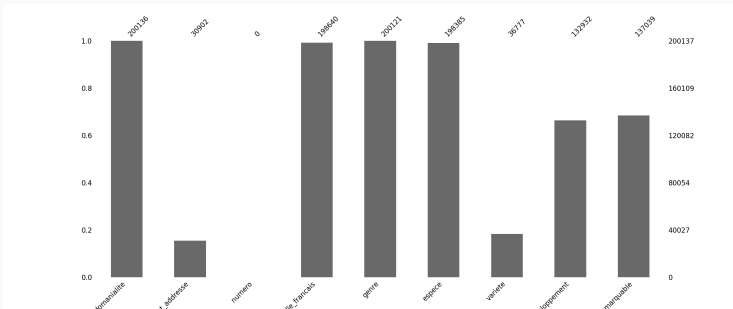


Figure 3: Nombre de données par colonnes pour les colonnes non intégralement remplies

Une colonne est vide !

- numero

Deux colonnes sont presque vide !

- complement adresse
- variete

- domanialite
- arrondissement
- complement adresse
- lieu
- geo point 2d a
- geo point 2d b

Colonnes gardée :

- domanialite
- arrondissement
- geo point 2d a
- geo point 2d b

Colonnes supprimées :

- complement adresse
- lieu

Traitement des données manquantes

- Une seule donnée manquantes

- Une seule donnée manquantes
- Le lieu de cette donnée est connu

- Une seule donnée manquantes
- Le lieu de cette donnée est connu
- L'ensemble des domanialités correspondant à ce lieu sont des Jardin

- Une seule donnée manquantes
- Le lieu de cette donnée est connu
- L'ensemble des domanialités correspondant à ce lieu sont des Jardin

Solution

- Remplacer la donnée manquante par Jardin

Libellé Français

```
sub_df = sub_df[sub_df['libelle_francais'].isna()]
sub_df
```

	id	libelle_francais	genre	espece
528	100589	NaN	Crataegus	japonica
1416	101521	NaN	Cladrastis	lawsoniana
2727	103209	NaN	Pinus	sylvestris
5282	106682	NaN	Sorbus	aria
5287	106687	NaN	Staphylea	colchica
...
200120	2024729	NaN	Magnolia	obovata
200121	2024730	NaN	Magnolia	delavayi subsp. potaninii
200128	2024737	NaN	Parrotia	n. sp.
200129	2024738	NaN	Non spécifié	sinensis
200132	2024741	NaN	Castanea	fargesii

1497 rows x 4 columns

Figure 4: DataFrame des données manquantes pour la colonne libelle francais

Libellé Français

- 1497 données manquantes
- Un libellé français correspond à un couple unique (genre, espece)
- Pour chaque donnée manquante, nous allons essayer de trouver un libellé français correspondant à un couple (genre, espece) parmi les donnée déjà présentes

```
sub_df = sub_df[sub_df['libelle_francais'].isna()]
sub_df
```

	id	libelle_francais	genre	espece
528	100069	Na/N	Coleoptera	geranea
5416	101521	Na/N	Cheliceris	leontiana
2727	103009	Na/N	Pisces	syndensis
5182	105662	Na/N	Scitidae	alpi
5267	105667	Na/N	Staphylinas	calchica
...
200120	2024729	Na/N	Hymenoptera	obsoleta
200121	2024730	Na/N	Hymenoptera	delongi subsp. putanensis
200128	2024737	Na/N	Parasitica	n. sp.
200129	2024738	Na/N	2021 Lygididae	sternus
200132	2024741	Na/N	Cucurbita	fargesi

1497 rows x 5 columns

Figure 5: DataFrame des données manquantes pour la colonne libelle francais

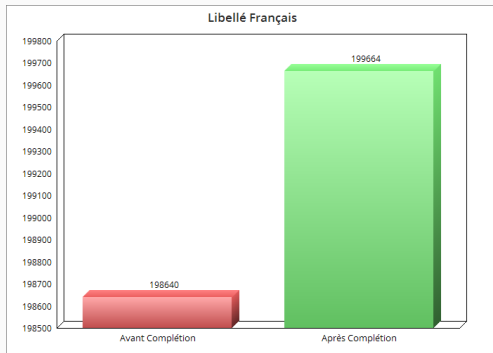


Figure 6: Résultat de la recherche de libellé français correspondant à un couple (genre, espece)

Résultat

- 1024 données manquantes ont pu être retrouvées

Détection des outliers

Nous avons les coordonnées GPS de chaque arbre :

- geo point 2d a
- geo point 2d b

Nous allons utiliser l'algorithme KNN pour détecter les outliers.

Données Géographiques - Résultats KNN



Figure 7: Résultat de la recherche des outliers avec KNN

Pour chaque arbre, nous avons les données suivantes :

- hauteur
- circonference

Nous allons utiliser la méthode des interquartiles pour détecter d'éventuels outliers.

Données numériques de l'arbre - IQR

Situation d'origine :

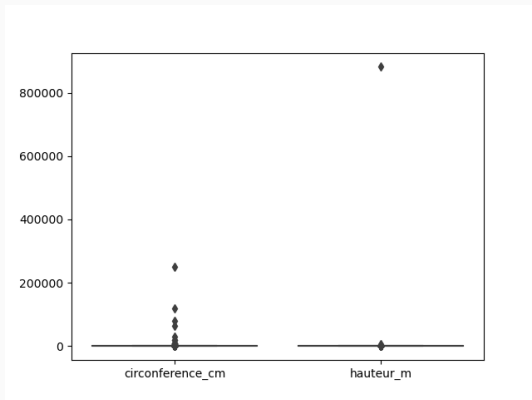


Figure 8: Boxplot des variables **circonference** et **hauteur**

Données numériques de l'arbre - IQR

- $IQR = Q3 - Q1$
- $Limite\ basse = Q1 - 1.5 \times IQR$
- $Limite\ haute = Q3 + 1.5 \times IQR$

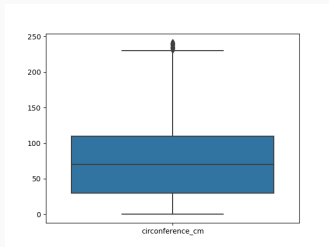


Figure 9: Boxplot de la variable circonference après application de la méthode des IQR

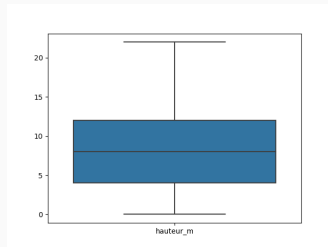


Figure 10: Boxplot de la variable hauteur après application de la méthode des IQR