

Airbnb Pricing Prediction

Amine OMRI

Résumé

Ce projet vise à traiter des données collectées sur le site Aribnb pour la date du 14 décembre 2020 qui comporte toutes les données relatives aux listings a cette date.

Ce projet comporte trois parties principales :

Anayse

01

le premier est une analyse de données qui vise à effectuer une analyse exploratoire des données et à réaliser des visualisations, ce qui inclut le traitement des données nécessaires.

Pré-traitement

02

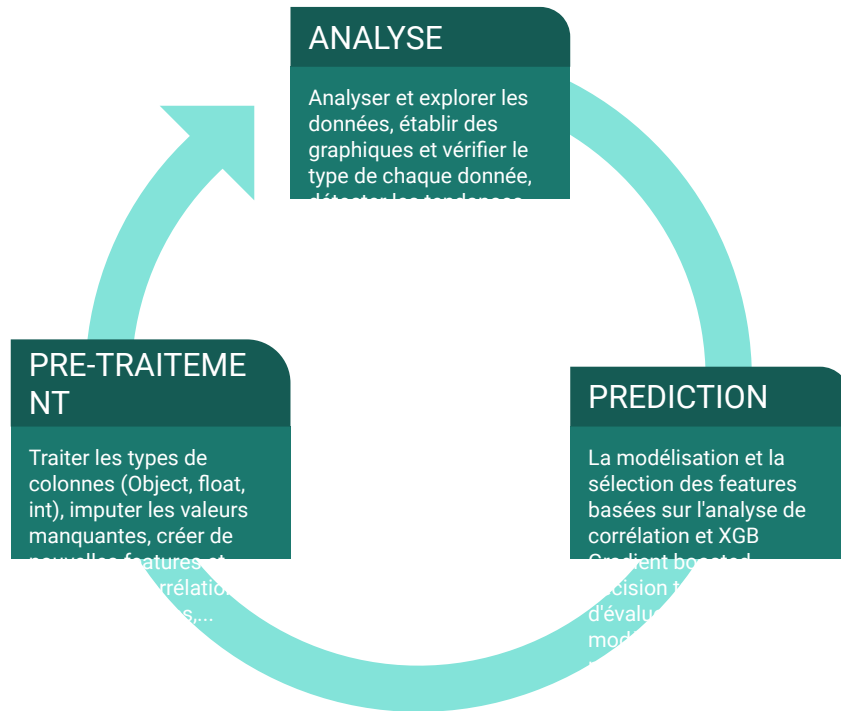
La seconde consiste à effectuer un prétraitement des données (nettoyage, encodage des données, Features engineering).

Prédiction

03

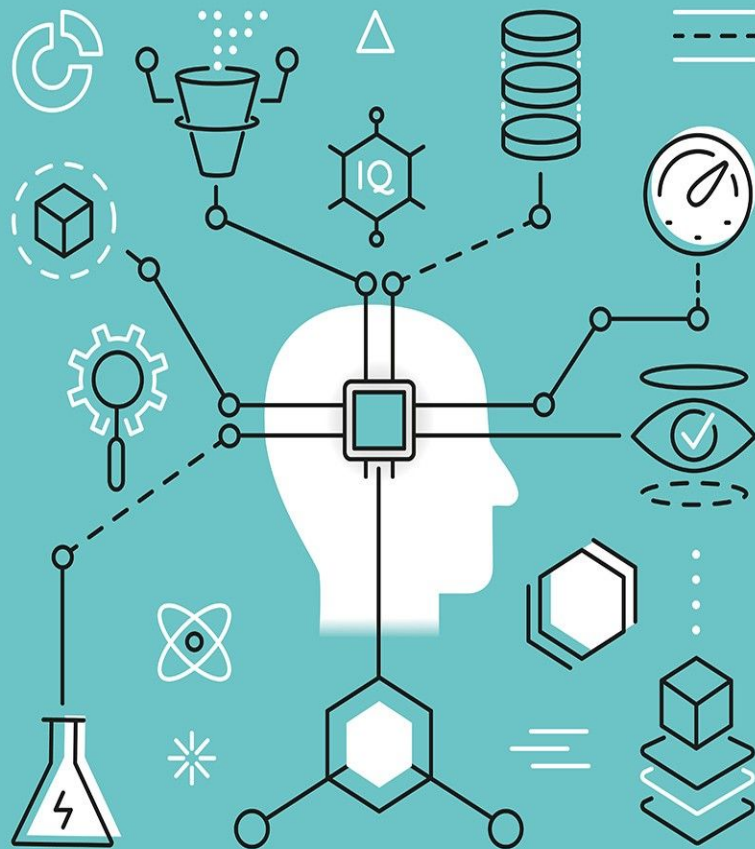
Le dernier est la modélisation et la prédiction qui comprend également certaines étapes de traitement telles que le traitement des valeurs aberrantes et la sélection des features, ainsi que des visualisations.

Cycle du projet



La partie Pré-traitement de données

Nettoyage et encodage et feature
engineering



Pré-traitement

1

Nettoyage du dataset listings

2

Traiter les valeurs manquantes

3

One-hot encode les variables catégorielles

4

Feature Engineering:

- **last_review**: cette colonne servira à filtrer les listes qui ne sont plus active
- **host_location**: nous pouvons l'utiliser pour déterminer si l'hôte est local ou non
- **host_since**: peut être utilisé pour calculer l'expérience des hôtes en fonction du nombres d'années depuis leur première inscription
- **amenities**: créer des features à partir de cette colonne en Prétraiter la colonne amenities pour extraire tous les valeur amenities possibles et affecter un identifiant à chacun d'entre eux, après avoir défini une fonction d'encodage qui mettra 1 dans l'index correspondant dans une matrice. Et enfin, on va avoir la matrice document - terme en appliquant cette fonction d'encodage à tous les documents du corpus.

Pré-traitement

5 Del meme facon (bag-of-words binaire) créer des features à partir de **host_verifications**

6 Créer des Feature à partir de colonnes de texte **description**

```
def nlp_pipeline(book_texts):  
    clean_books = []  
    for book in book_texts:  
        book = remove_hypens(book)  
        book_i = tokenize_text(book)  
        book_i = remove_characters_after_tokenization(book_i)  
        book_i = convert_to_lowercase(book_i)  
        book_i = remove_stopwords(book_i, custom_stopwords)  
        book_i = get_lemma(book_i)  
        book_i = remove_short_tokens(book_i)  
        book_i = keep_only_words_in_wordnet(book_i)  
        book_i = apply_lemmatize(book_i)  
        clean_books.append(book_i)  
    return clean_books
```

Apré avoir passer les documents par le pipline NLP il faut vectoriser le corpus maintenant que nous avons le corpus nettoyé, nous pouvons utiliser **TfidfVectorizer** pour convertir le texte en format vectoriel.

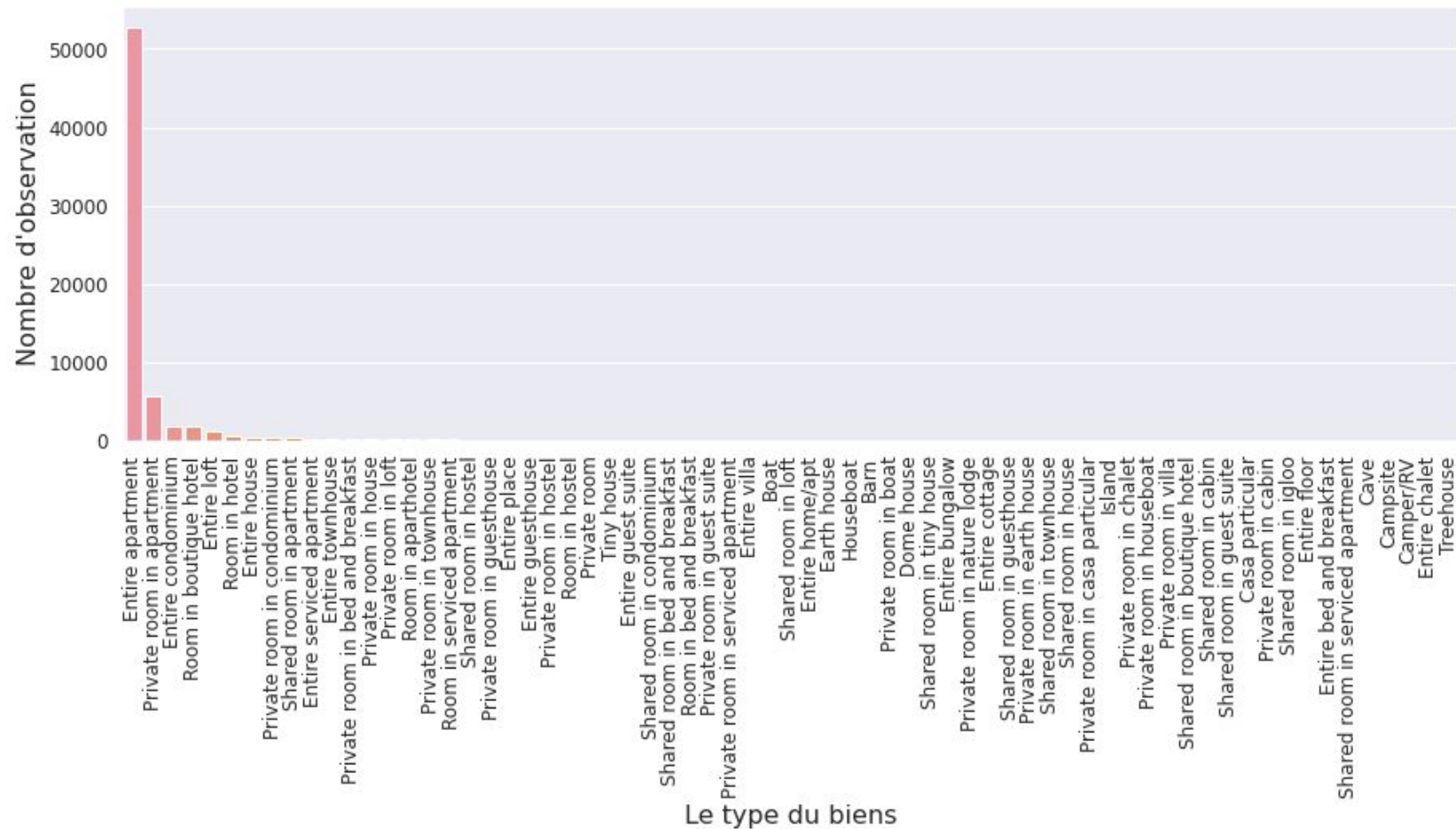
	abbess	able	absence	absolute	absolutely
0	0.0	0.0	0.0	0.0	0.0
1	0.0	0.0	0.0	0.0	0.0
2	0.0	0.0	0.0	0.0	0.0

La partie Analyse de données

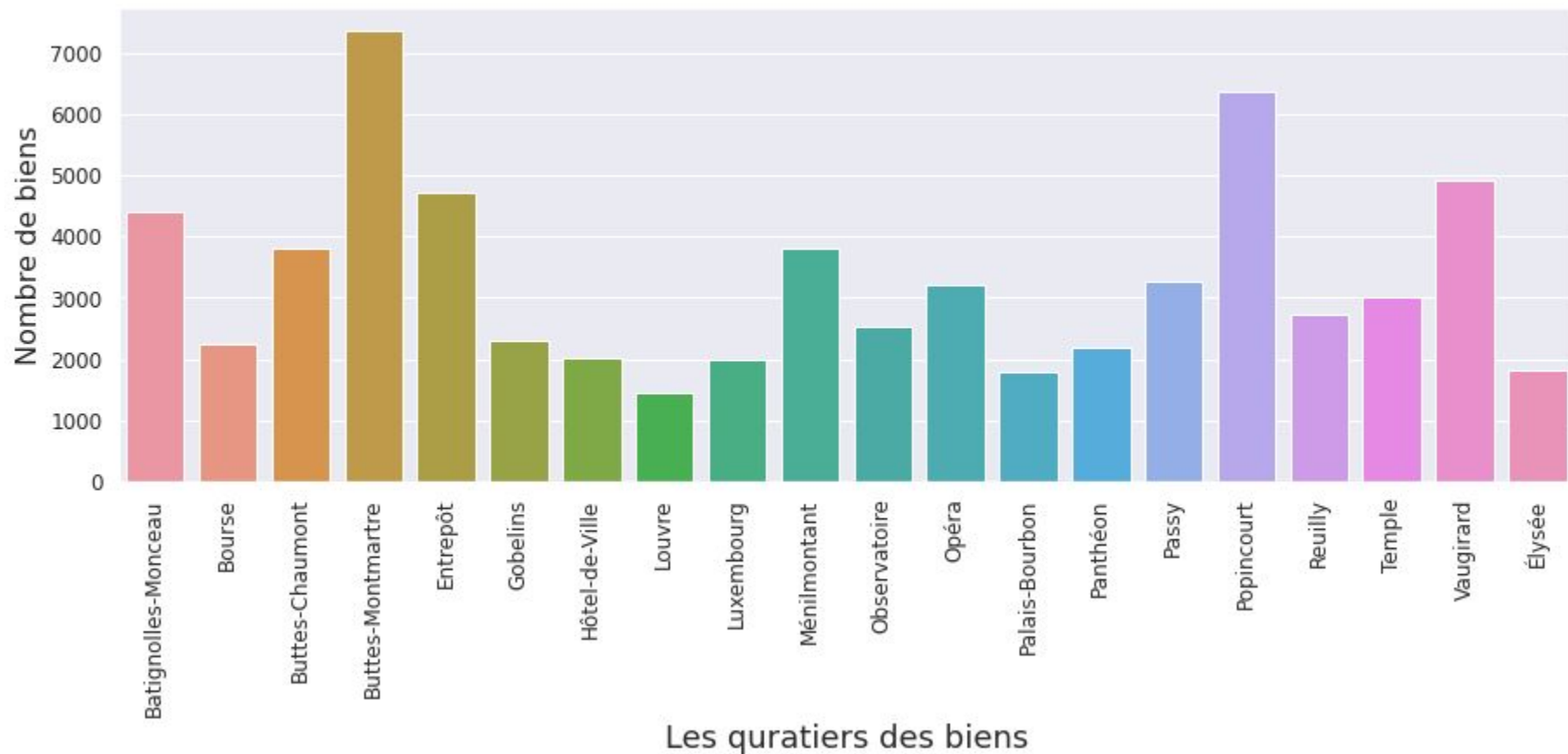
Analyser les données afin de répondre aux questions et d'obtenir des informations sur le prix des listings.



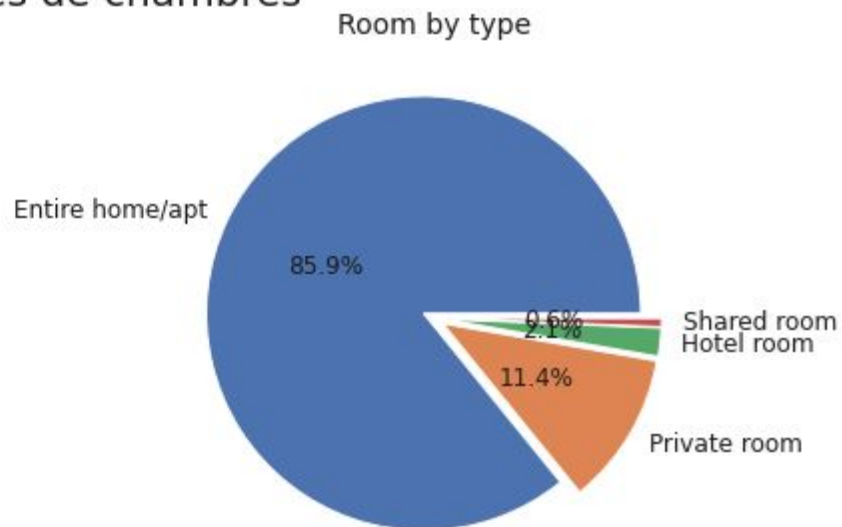
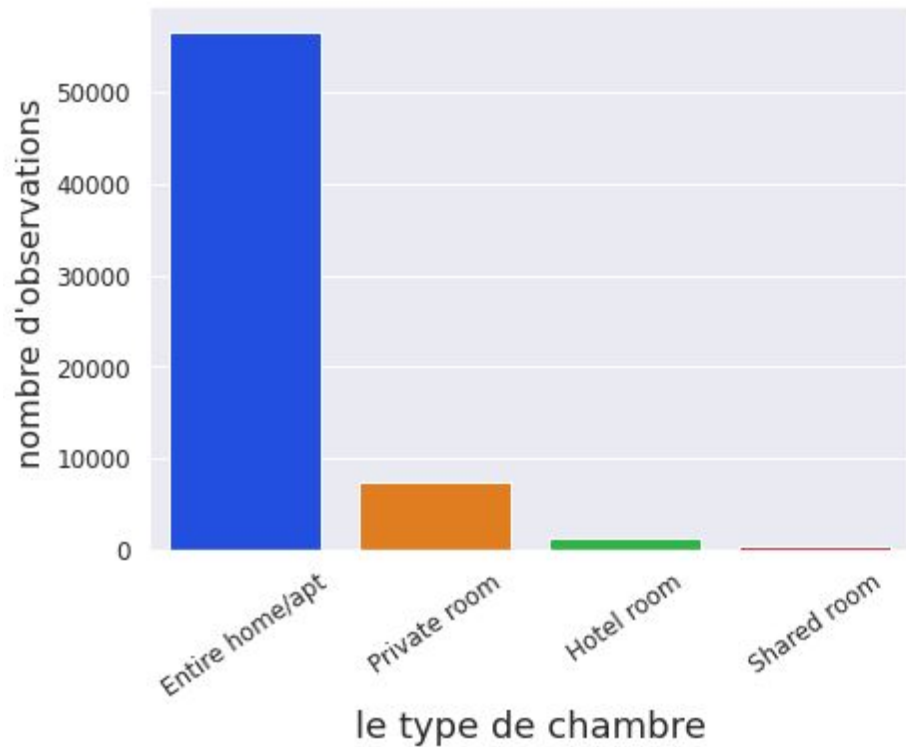
Répartition des listings par property_type



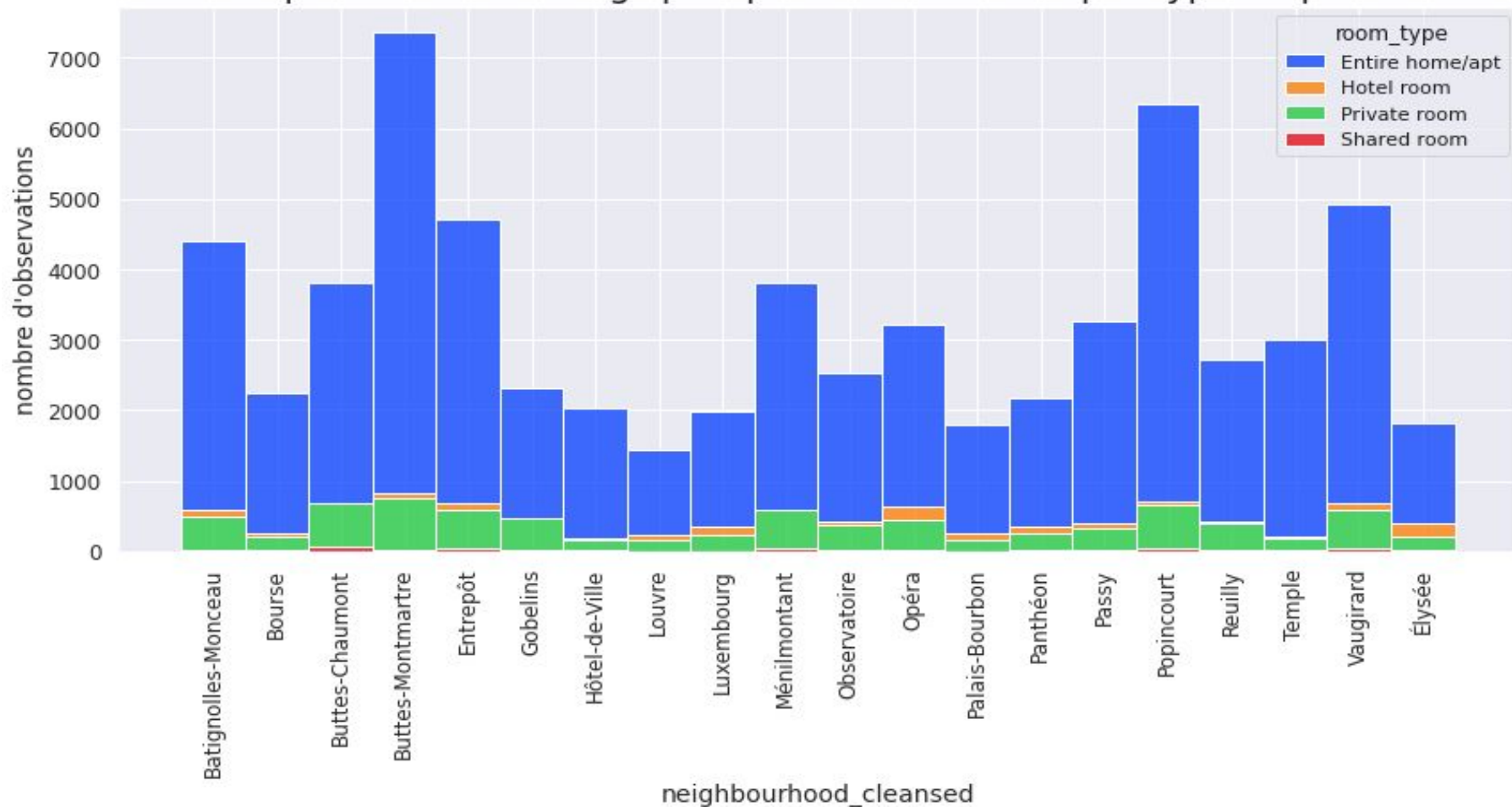
Répartition des biens par quartier



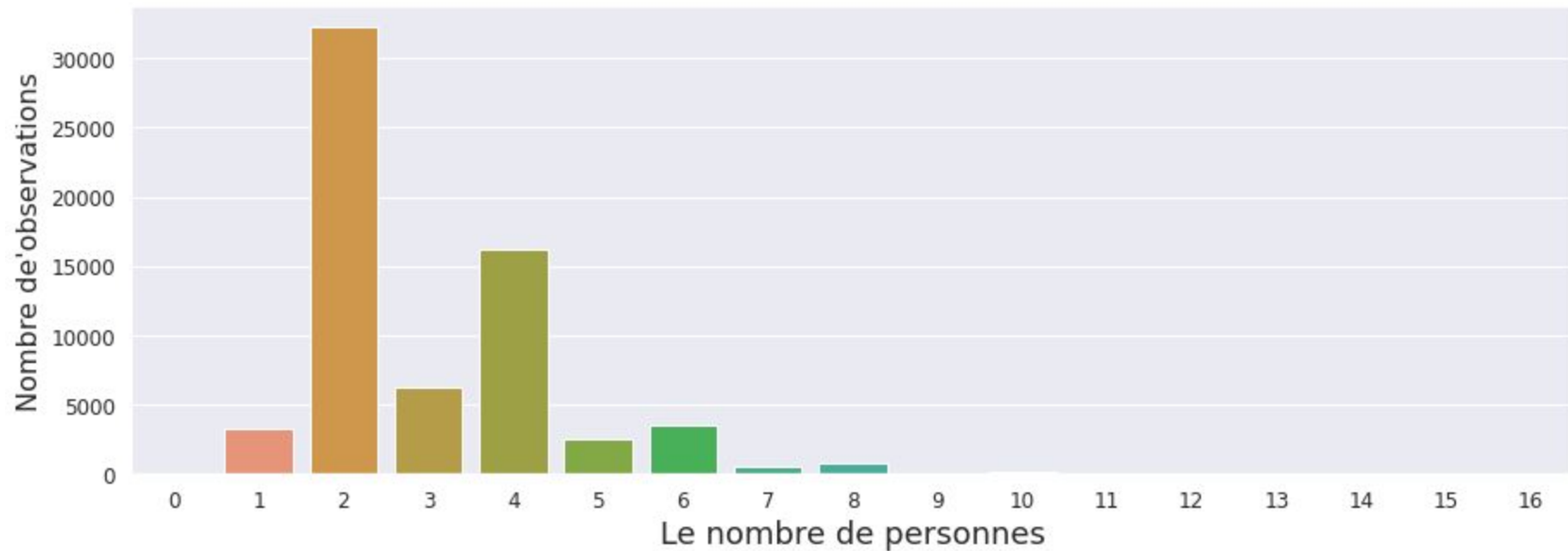
Répartition des types de chambres



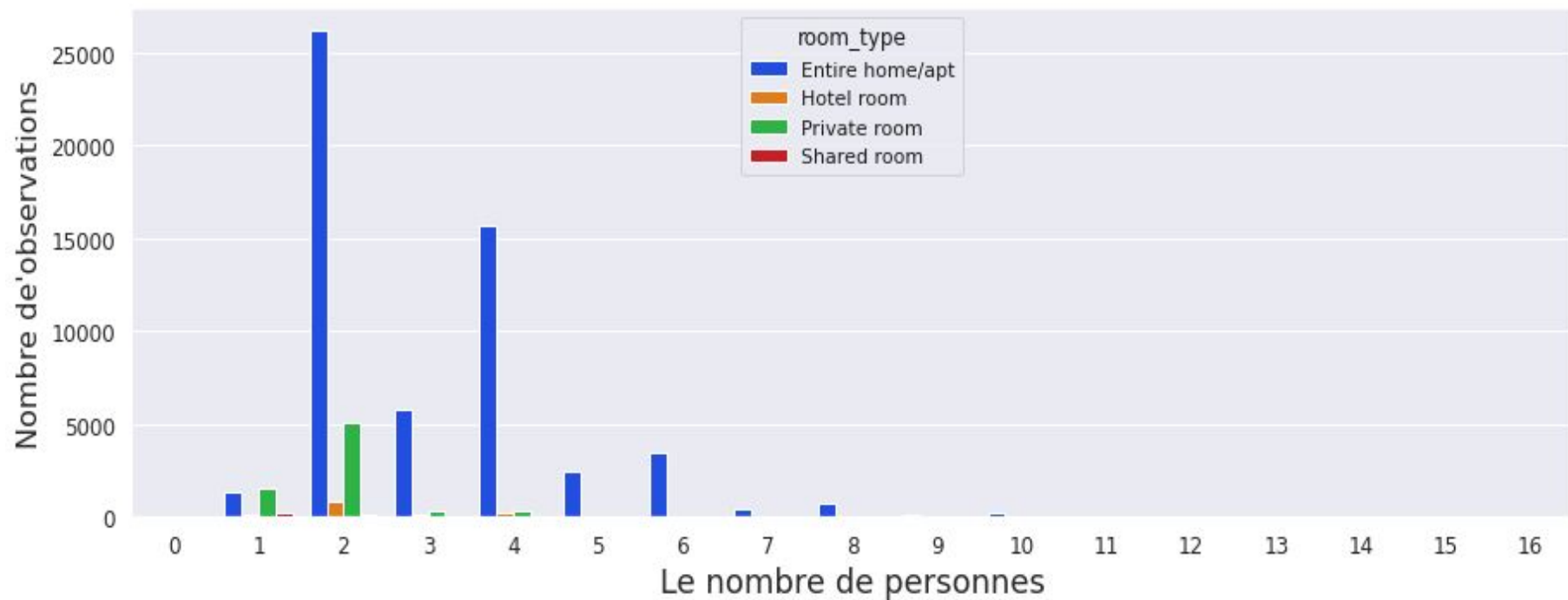
Répartition des listings par quartier et classée par type de pièce



Répartition des listings par nombre de pesonnes (accommodates)

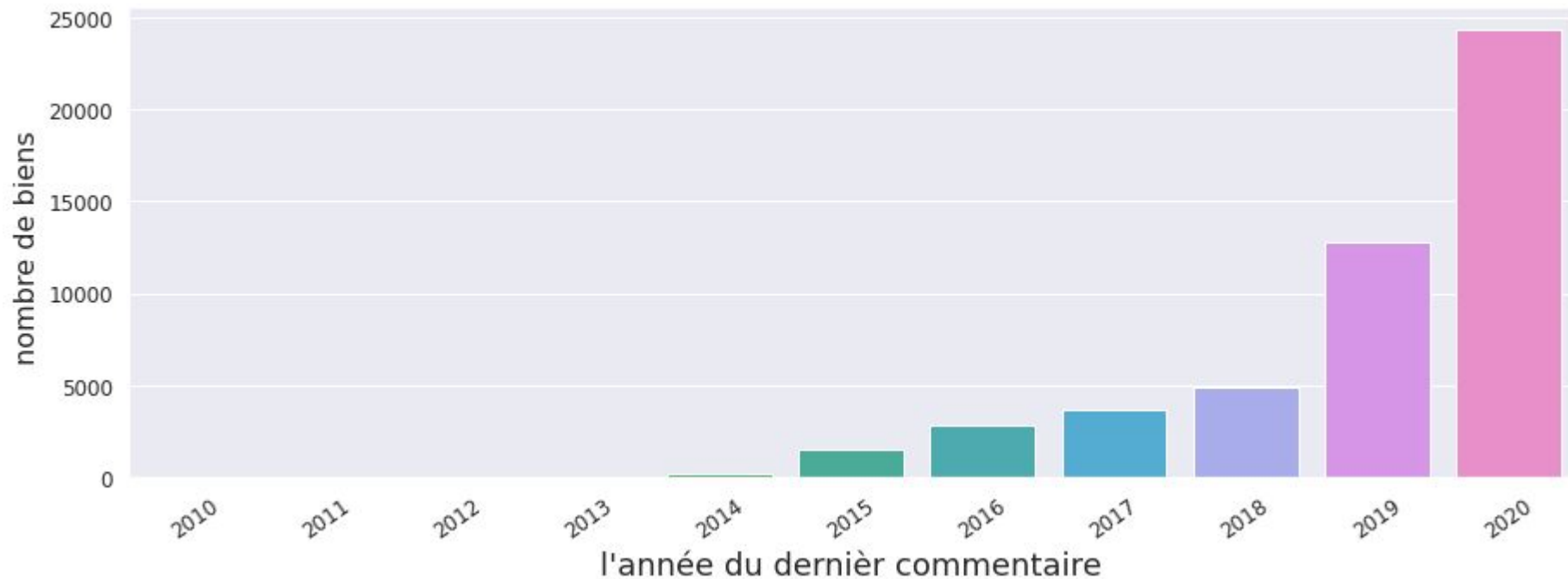


Répartition des listings par nombre de pesonnes (accommodates) classés par room_type



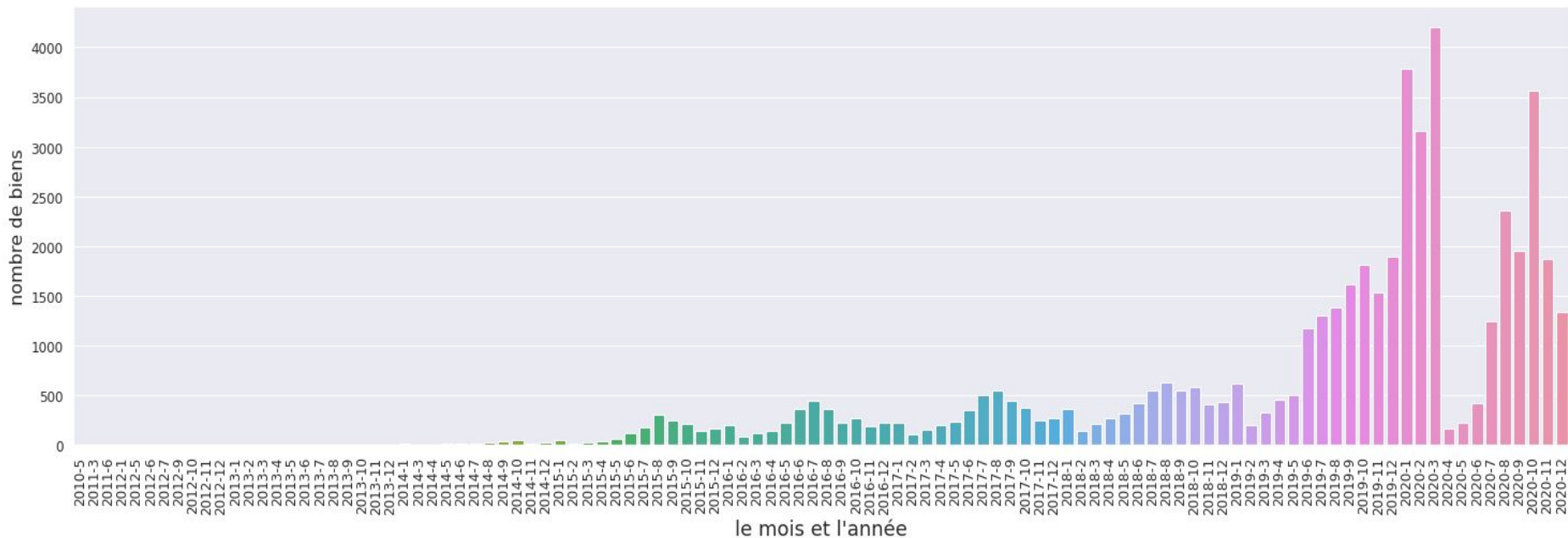
last_review: ce champ servira à filtrer les listes qui ne sont plus actives

Nombre de biens par dernière année où il a été commenté



last_review: ce champ servira à filtrer les listes qui ne sont plus actives

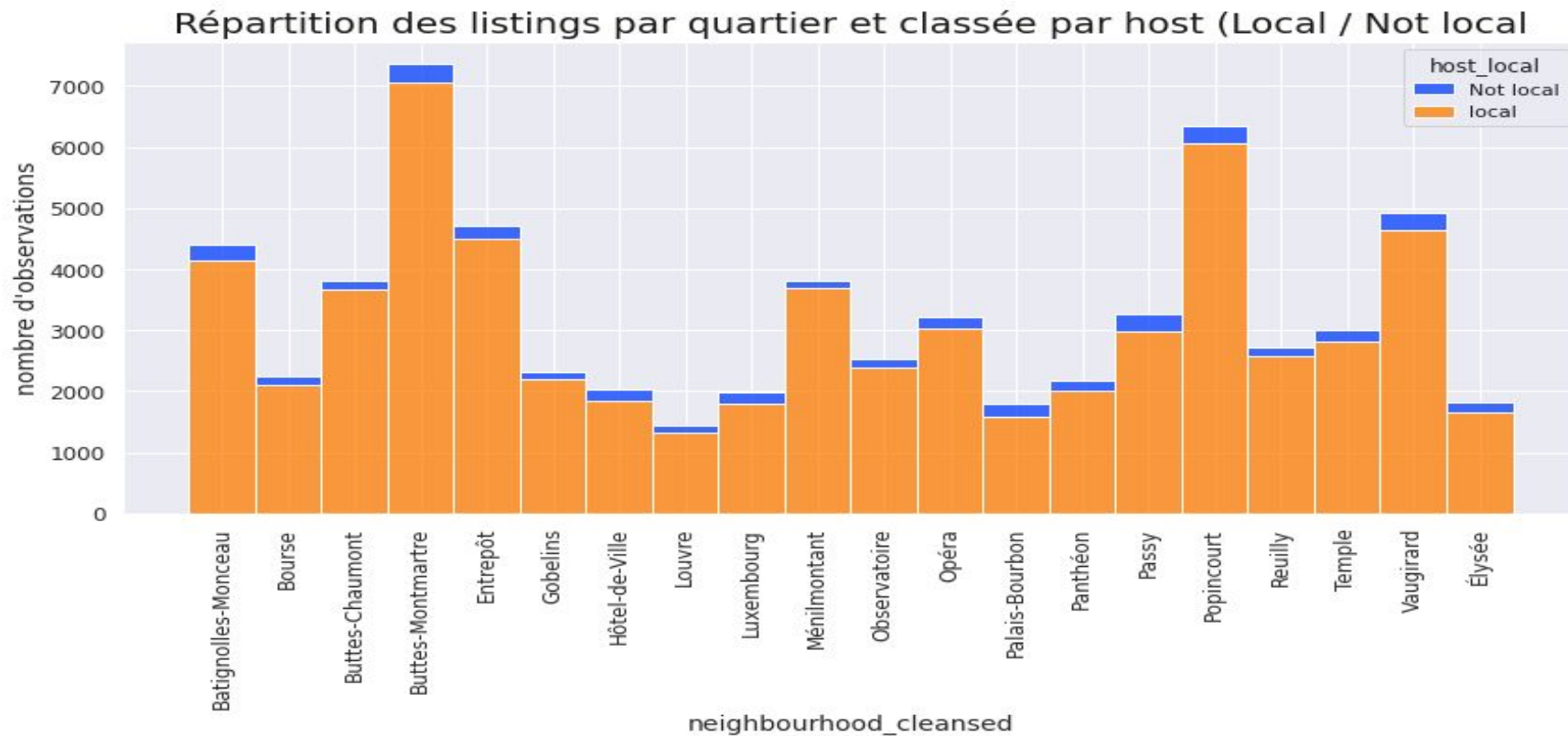
Nombre de biens par dernière année et mois où il a été commenté



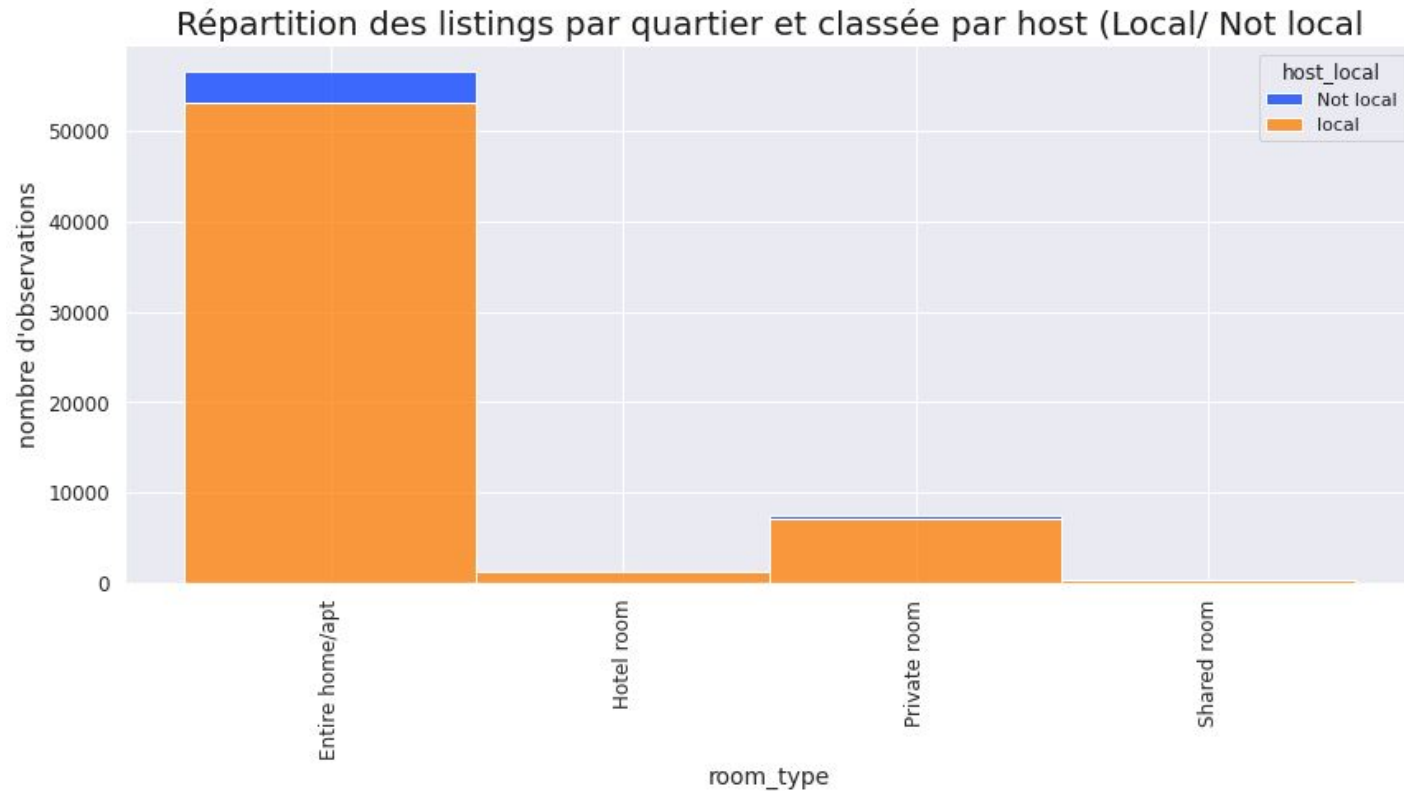
host_location: nous pouvons l'utiliser pour déterminer si l'hôte est local ou non



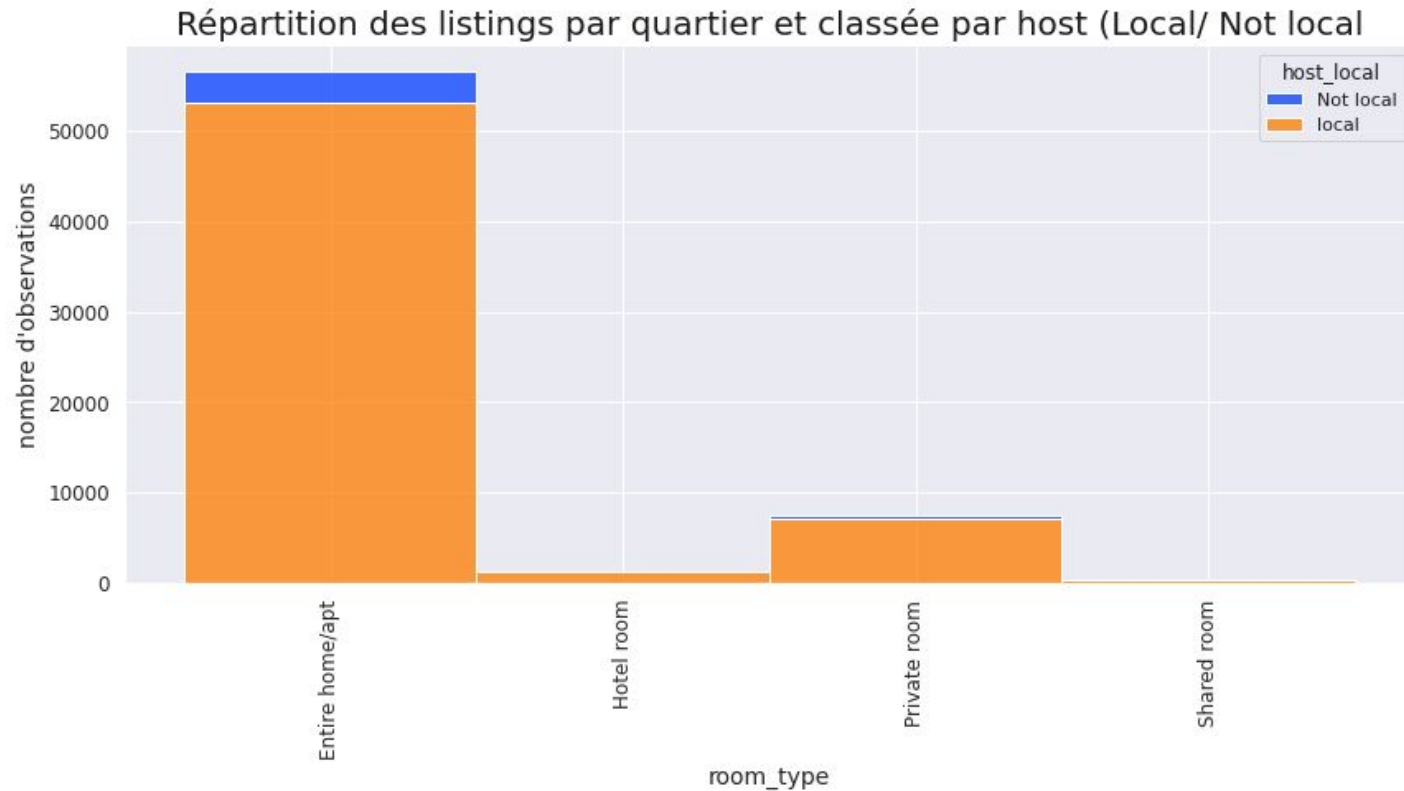
host_location: nous pouvons l'utiliser pour déterminer si l'hôte est local ou non



host_location: nous pouvons l'utiliser pour déterminer si l'hôte est local ou non

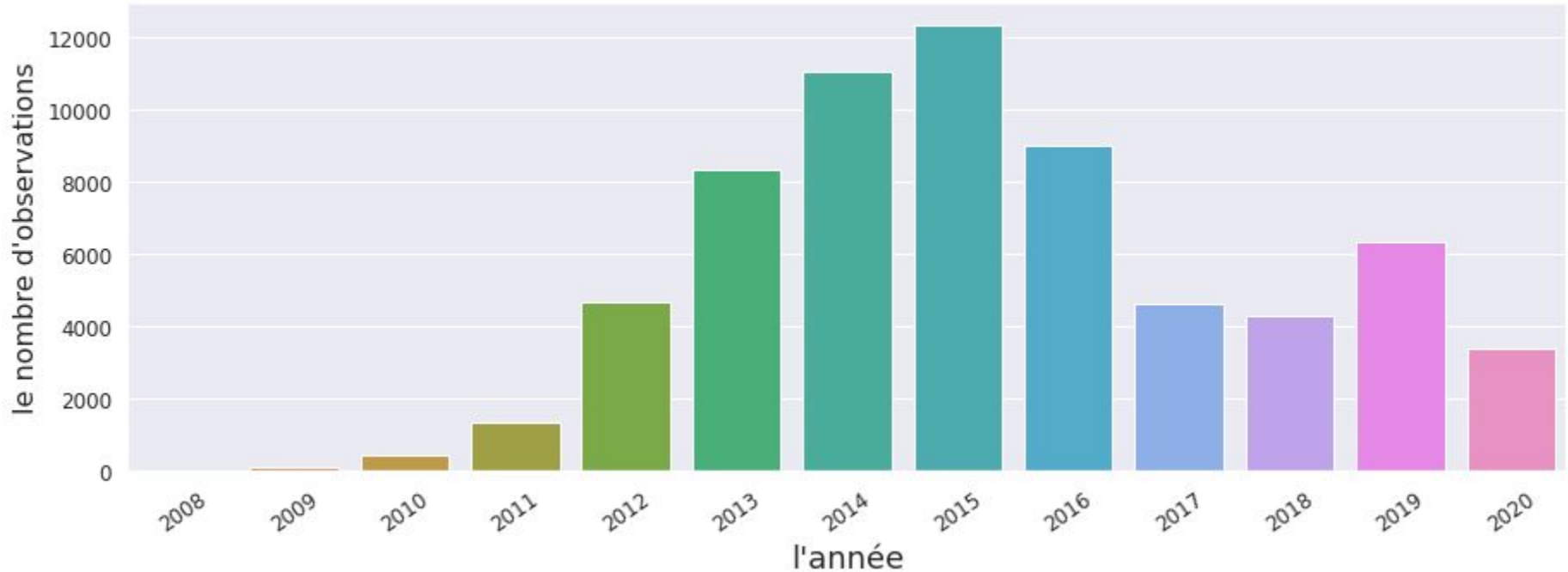


host_location: nous pouvons l'utiliser pour déterminer si l'hôte est local ou non



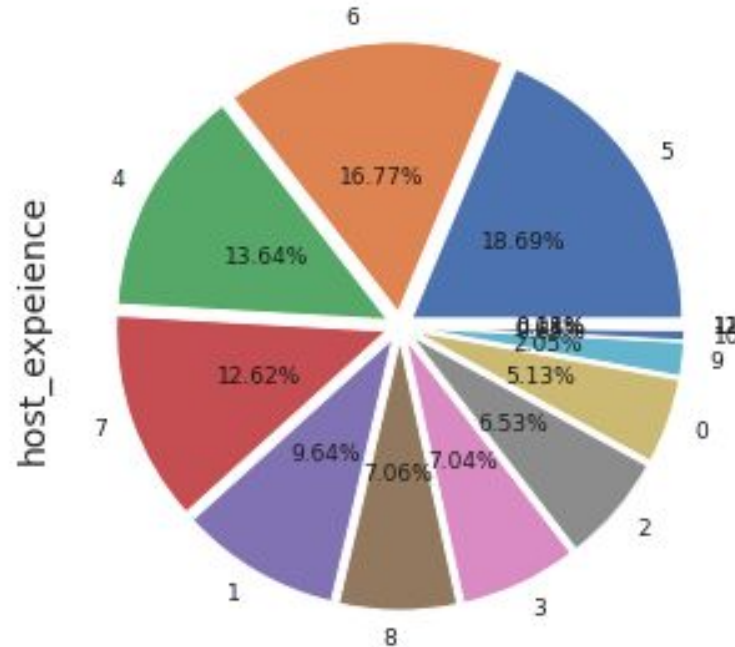
host_since: peut être utilisé pour calculer l'expérience des hôtes en fonction de la durée depuis leur première inscription

Répartition des hots par l'année depuis laquelle ils gèrent le listings

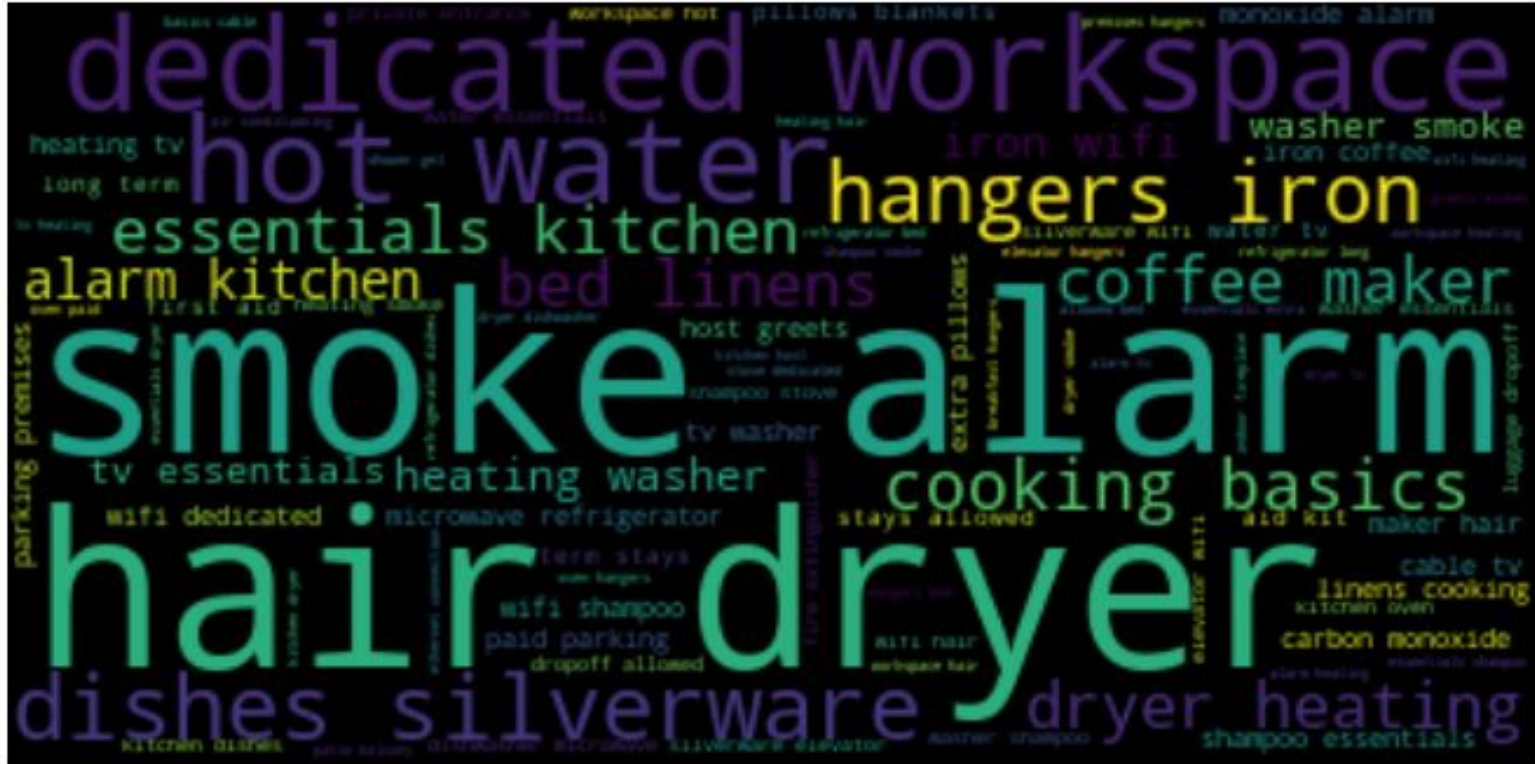


host_since: peut être utilisé pour calculer l'expérience des hôtes en fonction de la durée depuis leur première inscription

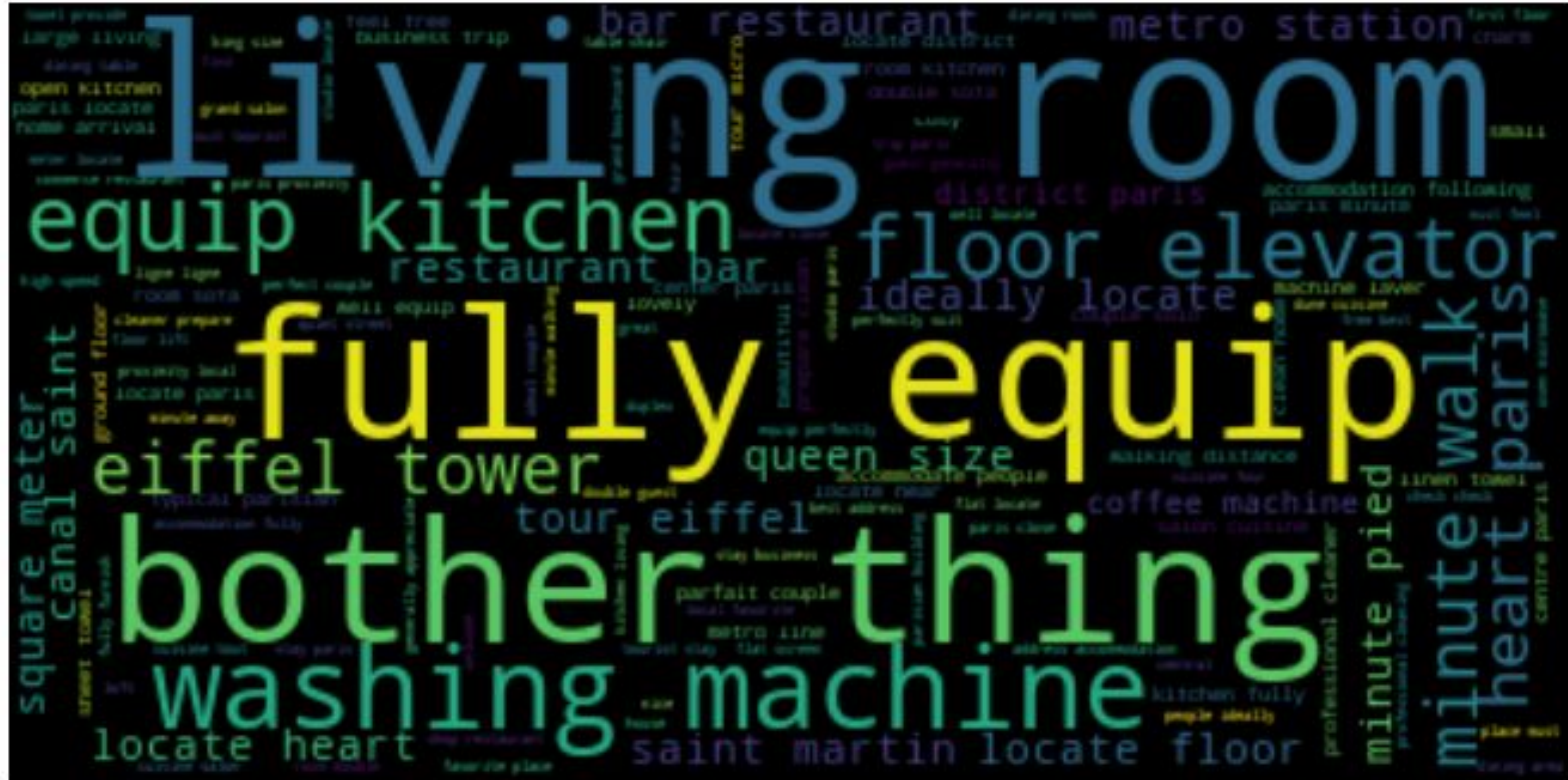
Répartition des hots par l'année depuis laquelle ils gèrent le listings



amenities: Créer des features à partir de amenities (équipements)



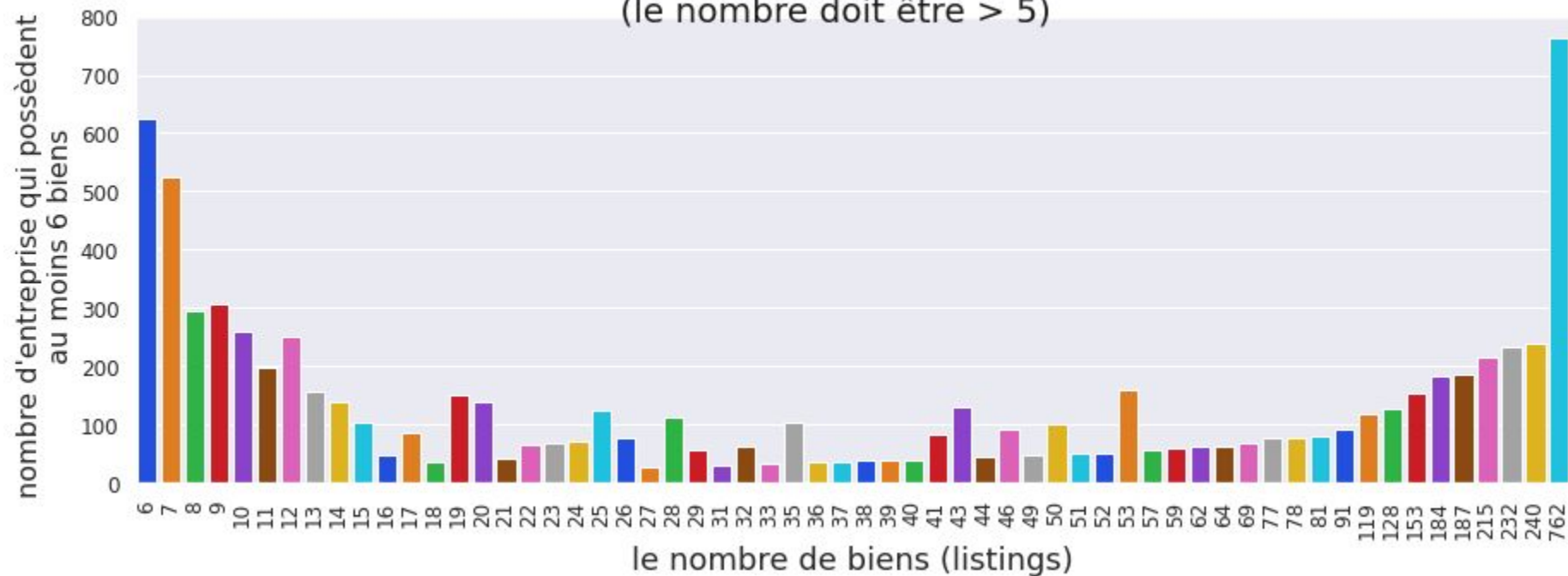
description : extraction de features à partir de la description (colonne textuelle) à l'aide d'un pipeline NLP



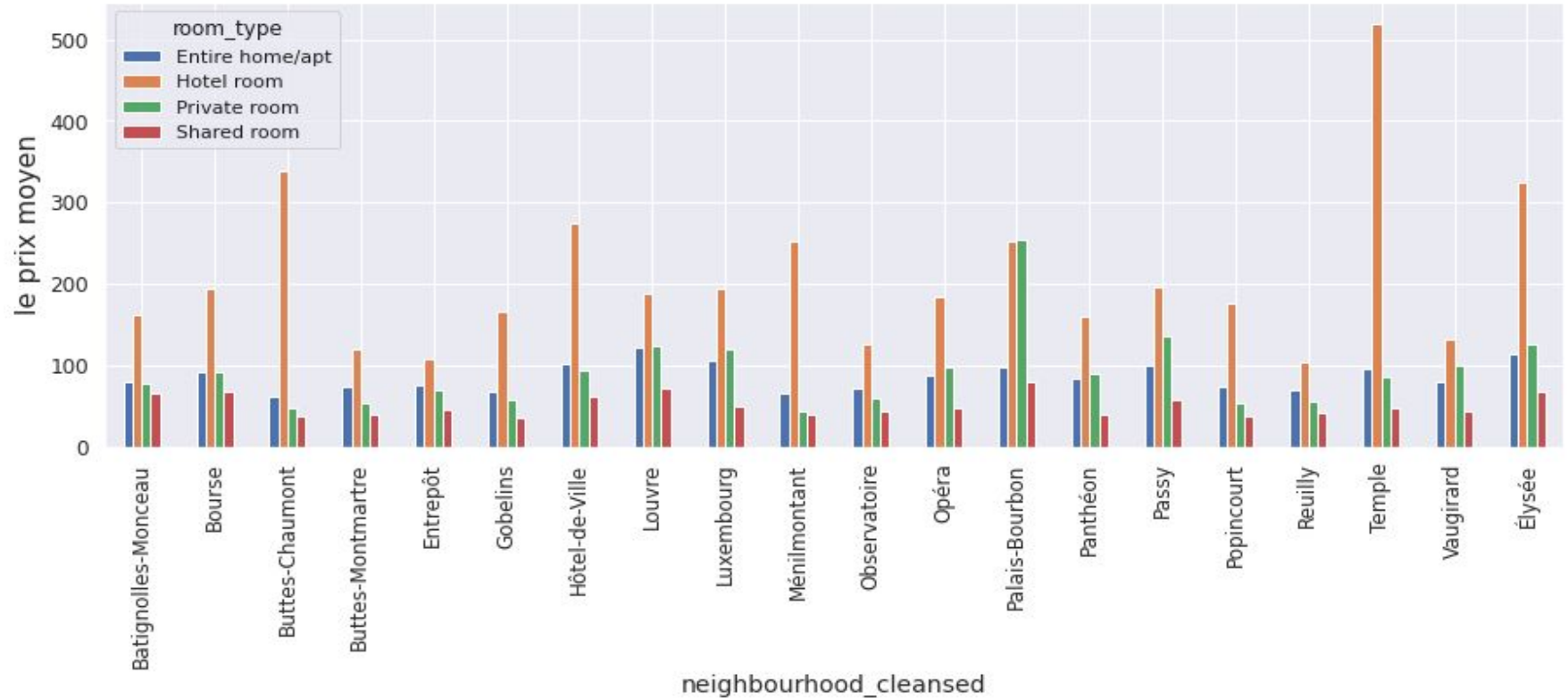
calculated_host_listings_count: valeur continue qui correspond aux nombres effectifs de listings pour les hôtes - mesure permettant de déterminer l'expérience des hôtes ou de distinguer les entreprises et les



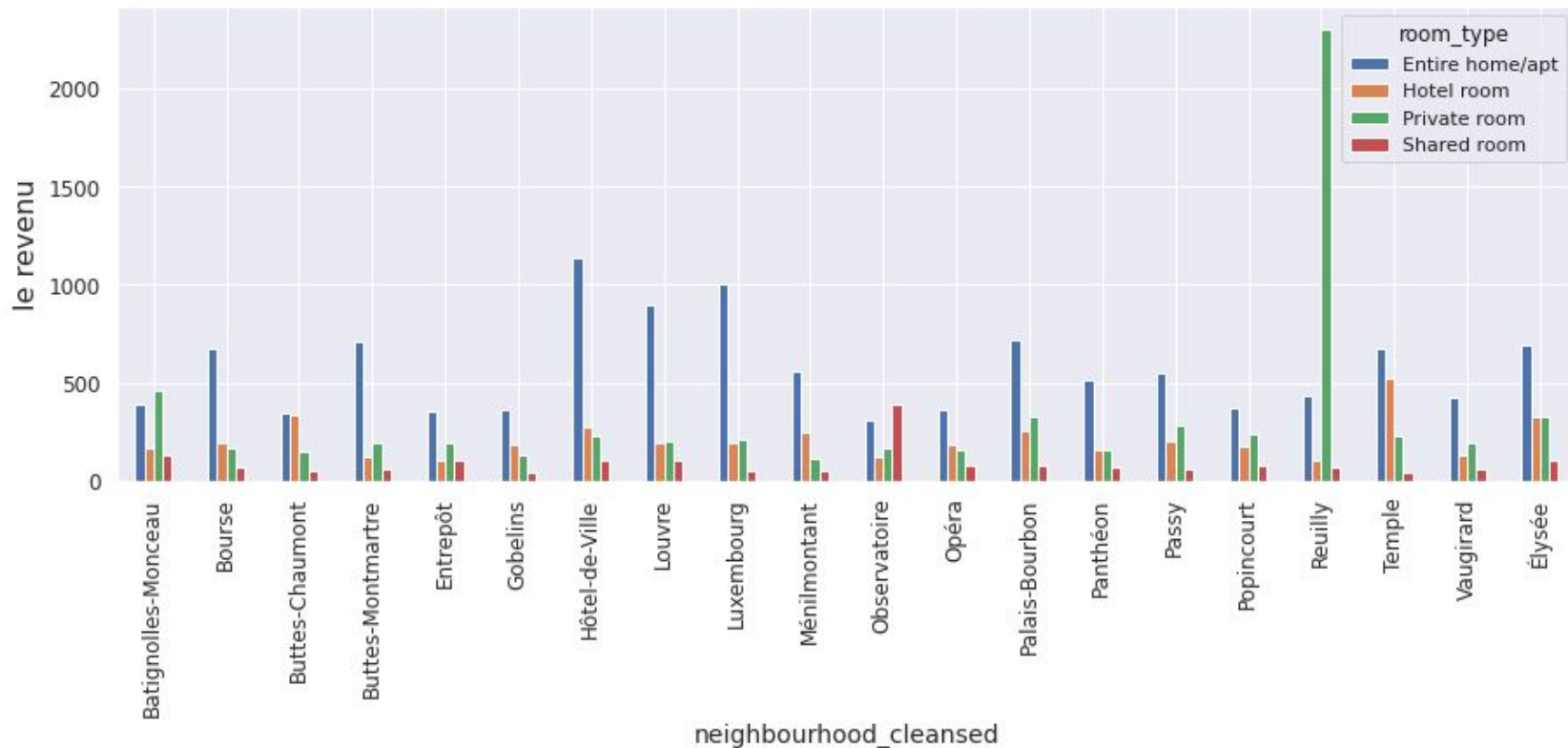
Nombre de listings qui peuvent être détenus par une entreprise
(le nombre doit être > 5)



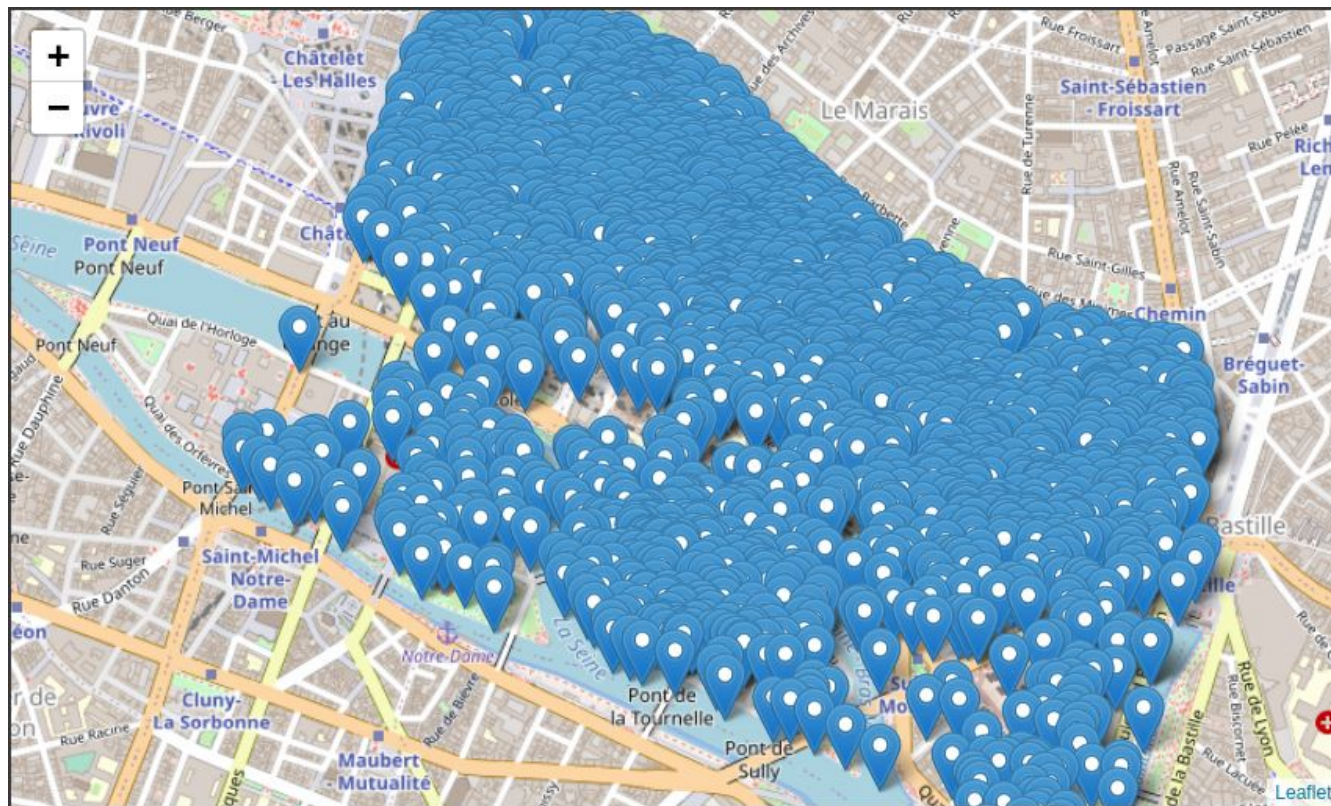
Visualisation du prix moyen du jour pour deux personnes



Calculer les revenus estimés pour chaque listing, les revenus estimés pour chaque listing seront calculés sur la base du prix d'une nuit et du nombre minimum de nuits à partir de la base de données



Les listings dans le quartier de l'Hôtel-de-Ville en utilisant folium

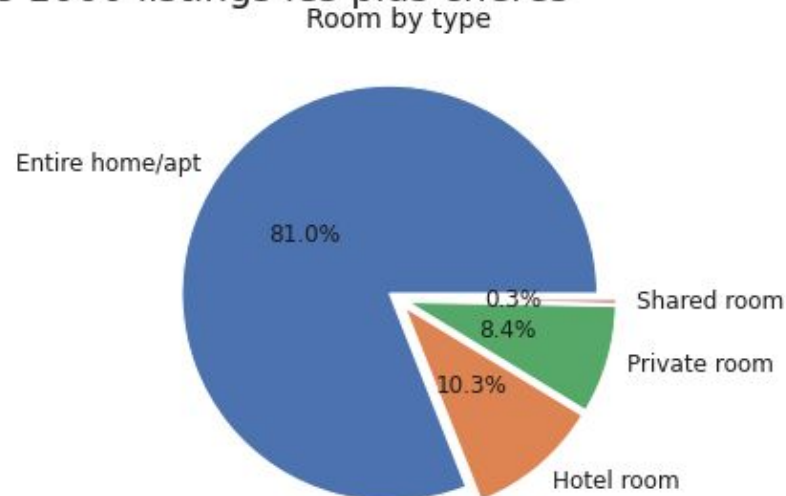
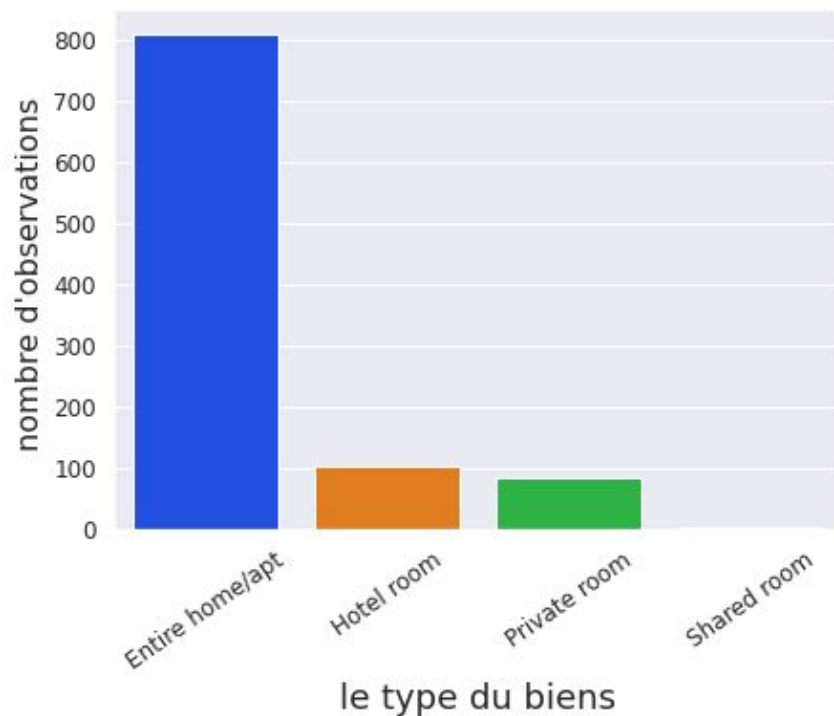


**Cette fois-ci en se basant sur la colonne prix au lieu de revenu
je récupère les 100 listings les plus chers**



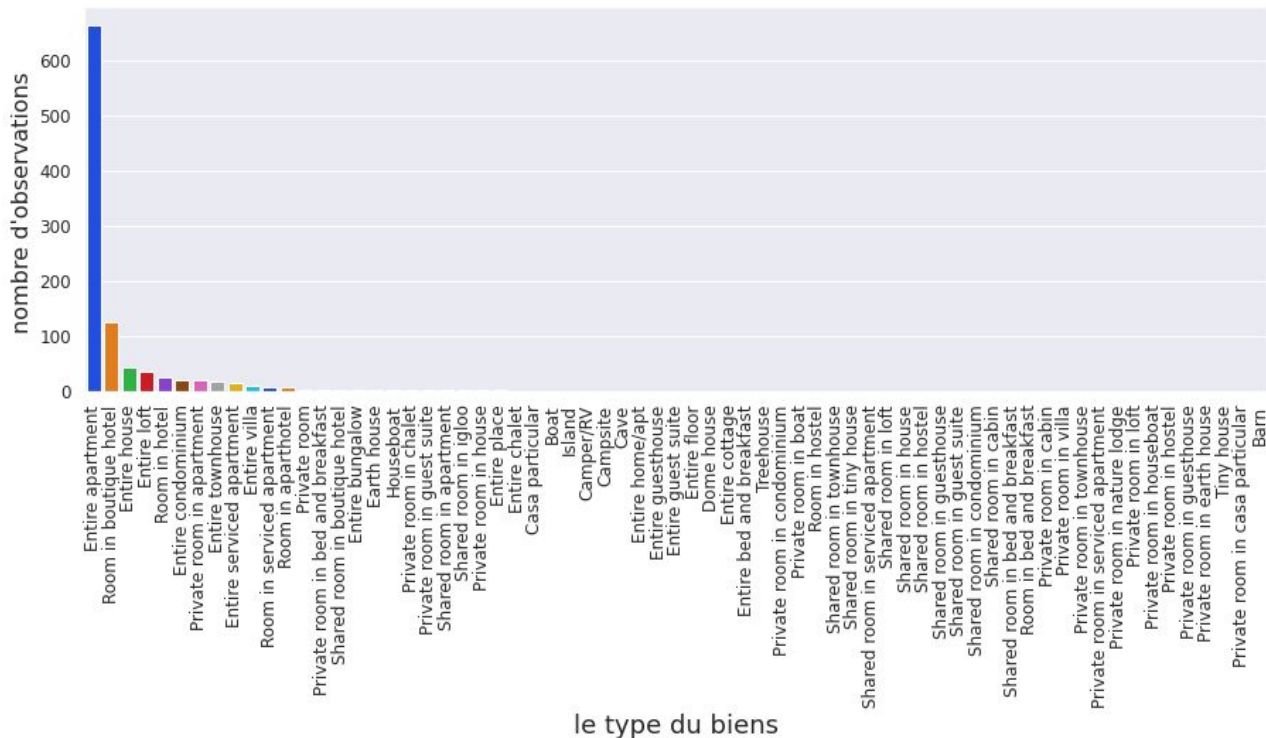
Type de bien a acheter

Répartition des property_type sur les 1000 listings les plus chères



Cette fois-ci en se basant sur la colonne prix au lieu de revenu
je récupère les 100 listings les plus chers

Répartition des property_type sur les 1000 listings les plus chers



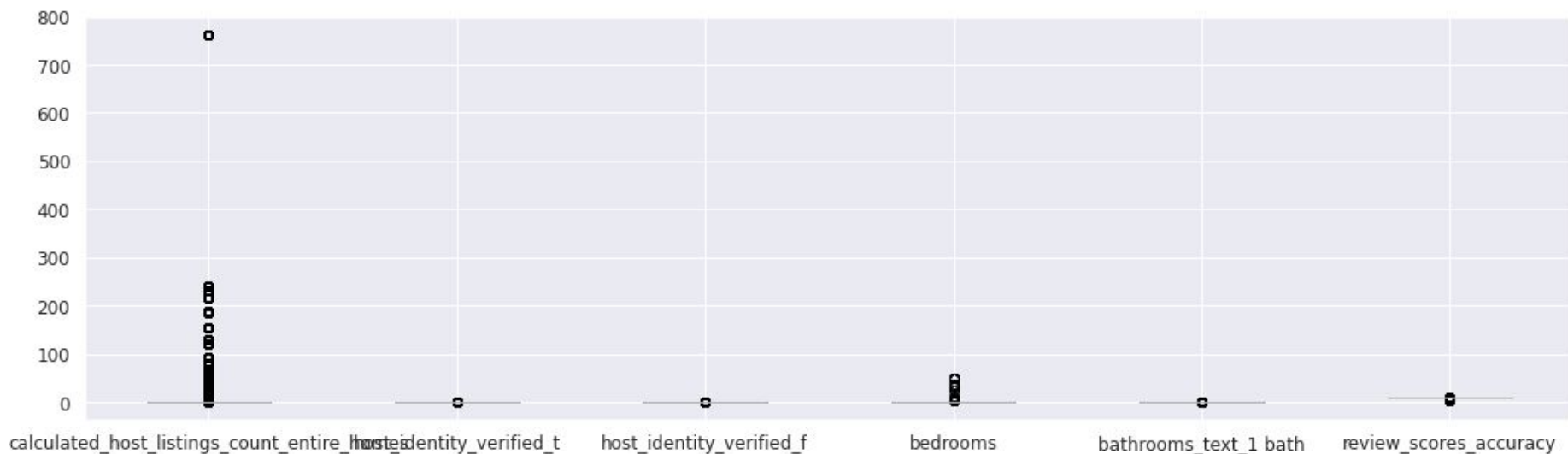
La partie prédiction du prix des listings

Ici, j'ai fait une sélection de features et une analyse de corrélation afin de trouver le meilleur modèle qui pourrait s'adapter à nos données



Le premier essai : j'ai utilisé un ensemble de données de toutes les colonnes prétraitées sauf celles créées dans la partie ingénierie de données (menities, host_verifications, description)

Suppression des colonnes comportant beaucoup de valeurs aberrantes



établir une graphique
collinearity heatmap
pour détecter les
features corrélées qui
pourraient empêcher
notre modèle de
converger



Normalisation et standardisation:

A l'exception de **accommodates** et de **host_experience**, les autres features numériques sont toutes asymétriques et pourraient bénéficier d'une transformation logarithmique

J'ai utilisé différentes méthodes pour standardiser mes données :

- StandardScaler
- MinMaxScaler
- RobustScaler

Et j'ai gardé MinMaxScaler puisqu'il conserve la structure globale des données



R SQUARED et RMSE

	RL	RIDGE	LASSO	XGB trees
TRAIN	Training RMSE: 0.1914 Training r2: 0.5128	Training RMSE: 0.1915 Training r2: 0.5125	Training RMSE: 0.1914 Training r2: 0.5128	Training MSE: 0.1827 Training r2: 0.5351
TEST	Validation RMSE: 0.191 Validation r2: 0.5193	Validation RMSE: 0.191 Validation r2: 0.5193	Validation RMSE: 0.191 Validation r2: 0.5193	Validation MSE: 0.1843 Validation r2: 0.5361

[illegible]

Le deuxième essai : J'ai utilisé un ensemble de données de toutes les colonnes prétraitées et aussi celles créées dans la partie ingénierie des données (menities, host_verifications, description)

Mon jeu de données était de taille (65917, 2673)

R SQUARED et RMSE

	RL	RIDGE	LASSO	XGB trees
TRAIN	Training RMSE: 0.1504 Training r2: 0.6172	Training RMSE: 0.1424 Training r2: 0.6376	Training RMSE: 0.1421 Training r2: 0.6384	Training MSE: 0.1714 Training r2: 0.5637
TEST	Validation RMSE: 2.8569801186007695 e+18 Validation r2: -7.190735871508386 e+18	Validation RMSE: 0.1567 Validation r2: 0.6056	Validation RMSE: 0.1572 Validation r2: 0.6043	Validation MSE: 0.1783 Validation r2: 0.5512

RANDOM FOREST

AVEC n_estimators = 20	AVEC n_estimators = 100
Training RMSE: 0.0271 Validation RMSE: 0.1774 Training r2: 0.931 Validation r2: 0.5535	Training RMSE: 0.0234 Validation RMSE: 0.1688 Training r2: 0.9403 Validation r2: 0.5753

MERCI
POUR VOTRE
ATTENTION

Questions ?