

# De nouvelles fonctionnalités pour la start-up Avis Restau

Avis Restau



# Contexte



**Avis Restau** propose une solution de mise en relation des clients et des restaurants.

Cet outil permet à nos **utilisateurs** de trouver des établissements adaptés à leurs besoins et aux **restaurants** d'attirer de nouveaux clients.

Mais pour les aider et continuer sa croissance notre start-up aimeraient **ajouter de nouvelles fonctionnalités de collaboration**.

# Objectif du projet



Suite à l'ajout de nouvelle fonctionnalités permettant aux utilisateurs de l'application de :

- poster des avis sous forme de commentaires
- poster des photos prises dans le restaurant

Il nous faut déployer deux solutions d'IA:

- détecter les sujets d'insatisfaction présents dans les commentaires postés
- labelliser automatiquement les photos postées sur la plateforme.

# Modélisation thématique (*Topic Modelling*)



COLLECTE



Tokens = [available, capability, cea, challenge, chemistry, chimie, coating, commercial, ...]

EXTRACTION DE FEATURES

Dictionary = {'available': 0, 'capability':  
1, 'cea': 2, 'challenge': 3, 'chemistry':  
4, 'chimie': 5, 'coating': 6,  
'commercial': 7, ...}



Corpus (BoW)

INFERENCE / MODÉLISATION



# Modélisation thématique (*Topic Modelling*)



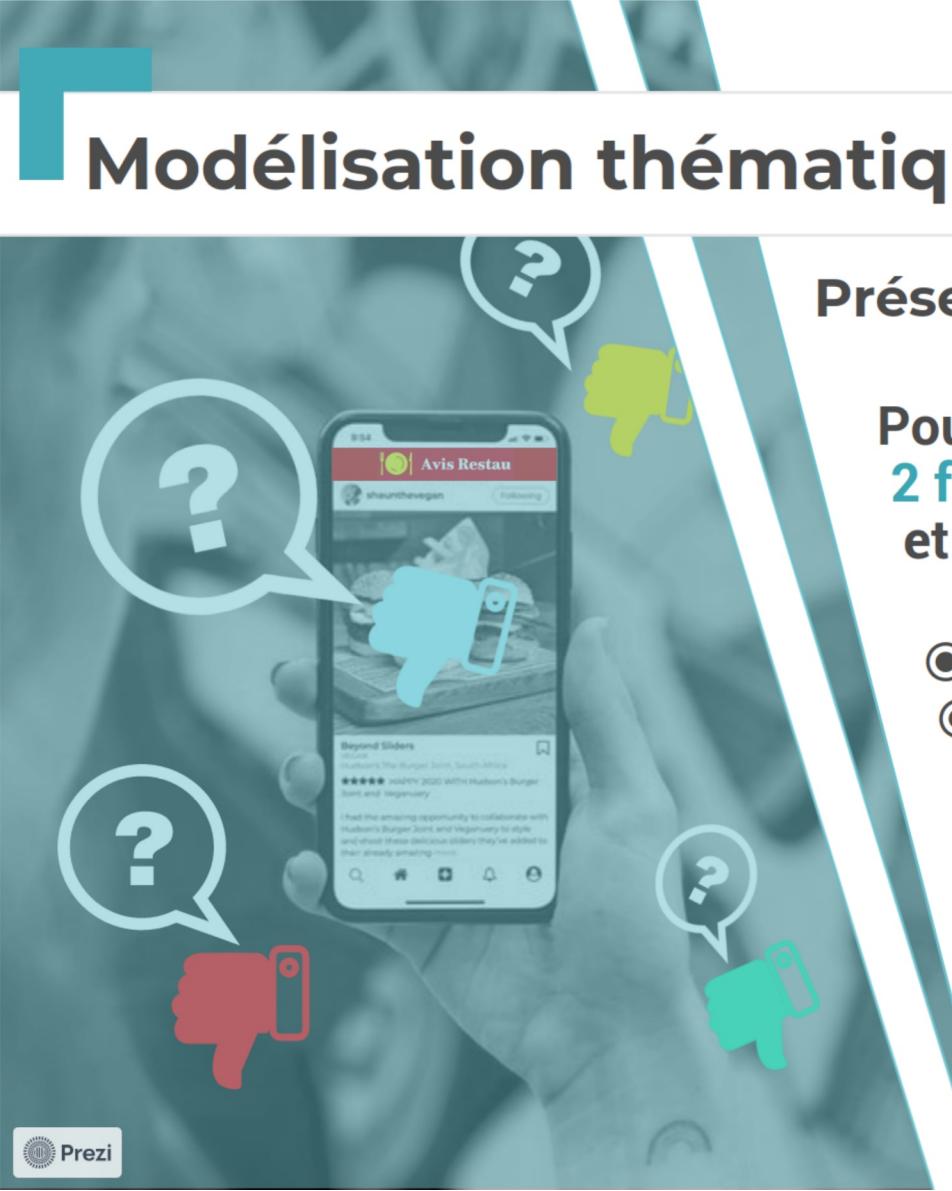
## Présentation du jeu de données

Pour ce problème, nous avons travaillé avec **2 fichiers source** décrivant les **reviews** laissées et les **commerces** concernés.

- **150.346 commerces** décrits par 14 variables
- **6.990.280 reviews** décrites par 9 variables

Mais nous avons constatés qu'il y avait de nombreux types de commerces...

- **51.864 commerces** restants
- **677.372 reviews** sur 1 millions chargées
- **10.000 reviews** choisies au hasard



# Modélisation thématique (*Topic Modelling*)



## EDA - vérifications de base

Les **vérifications** des erreurs les plus fréquentes (valeurs manquantes, doublons, outliers, erreurs de format, erreurs lexicales, contenus multiples) puis les **analyses univariées et multivariées** ont permis de déceler quelques problèmes qui ont été corrigés avec les actions suivantes:

- Suppression des colonnes inutiles au projet
- Traitement des outliers
- Traitement du déséquilibre des notes

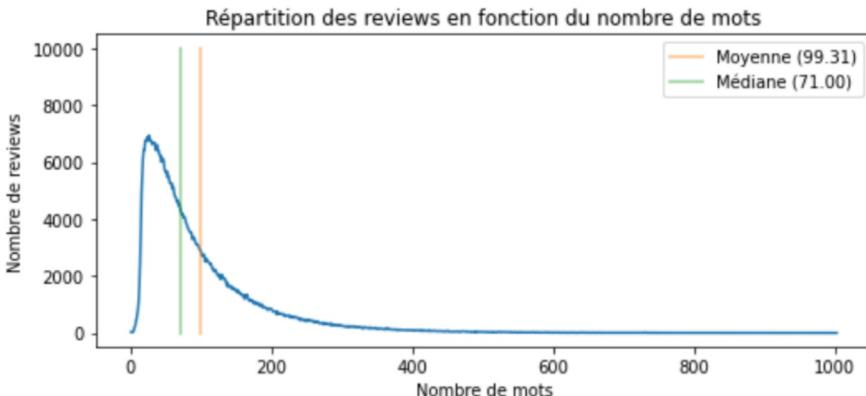


# Modélisation thématique (*Topic Modelling*)

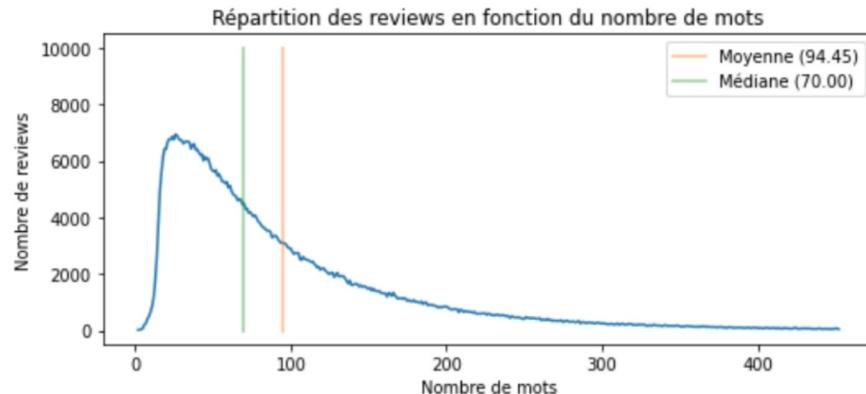


## EDA - Nombre de mots / reviews > outliers

AVANT



APRÈS



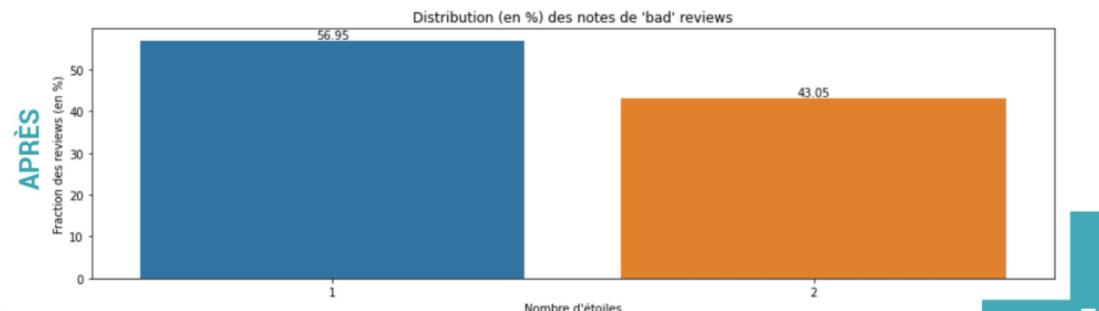
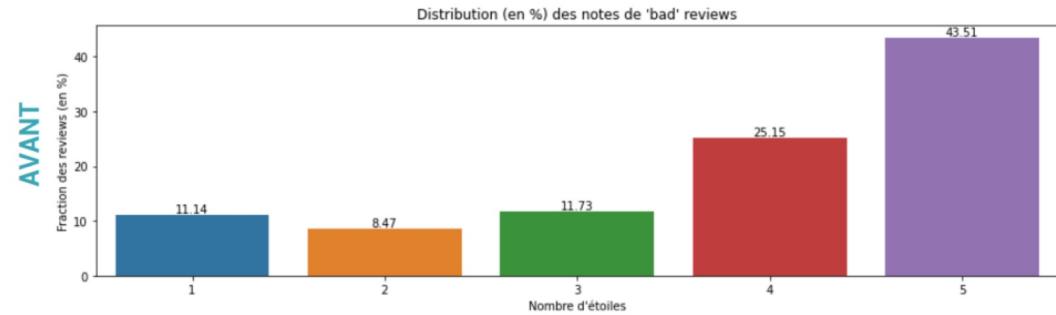
Boo  
酒很好，这里有关意大利和加州  
Awesome  
少有的鲜美，而环境又非常不  
C'mon!!!!  
Lil  
Wow  
!期四特供的炒饭就是剩饭炒饭  
.トラン。ディナーを食べに行  
Booo  
.com/2015/03/chees...  
!饭都比较一般。达不到惊喜  
好食，點啖牛肉烏東同埋芝士  
比都毫不逊色，我最爱他家的  
Nice  
!馆对华人来讲很方便，旁边就  
Great  
てとても美味しかったです。

Nombre de reviews

# Modélisation thématique (*Topic Modelling*)



## EDA - Répartition des notes > binarisation



# Modélisation thématique (*Topic Modelling*)



Pre-processing & mise en application > démo



# Classification d'images



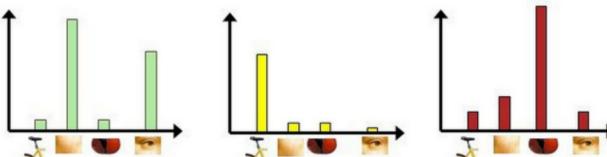
Avis Restau



COLLECTE



EXTRACTION DE FEATURES



SIFT  
ORB

python  
numpy  
pandas

MODELISATION

**Image Classifier**  
k-means / Convolutional Neural Network / ...

k-means

CNN

INFÉRENCE



# Classification d'images



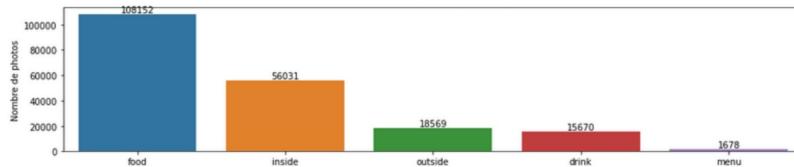
Avis Restau



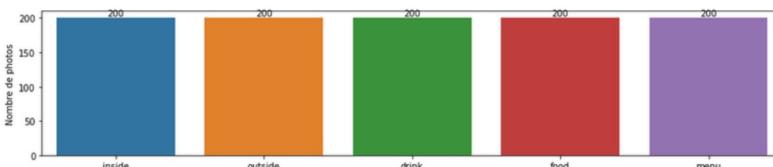
## Présentation du jeu de données

Pour ce problème, nous avons travaillé avec un **fichier source** donnant des informations sur les **photos collectées** par les utilisateurs de Yelp.

- **200.100 photos associées à 4 variables**
- **5 catégories**



- **200 \* 5 photos sélectionnées**



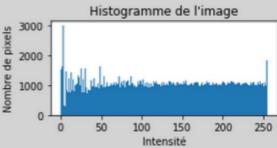
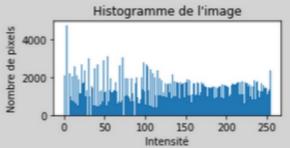
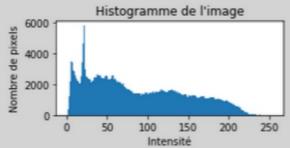
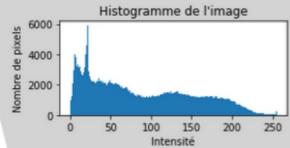
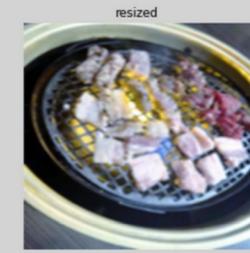
# Classification d'images



Avis Restau

## Pré-traitement des images

data/photos/MLIL0eZ0mN1kIgagfRE-A.jpg FOOD



(1)

(2)

(3)

1. suppression du bruit
2. égalisation de l'histogramme
3. redimensionnement



# Classification d'images



Feature extraction SIFT & CNN > [démonstration](#)



# Axes d'amélioration



Avis Restau



## ● Essayer avec plus de données

(pour voir si l'on peut améliorer le modèle de classification et jusqu'à quel point ou pour essayer de nouvelles distributions de topics)

## ● Essayer d'autres algorithmes / architectures

(EfficientNet, ResNet pour les CNN, Latent Semantic Analysis pour les topics)

## ● Essayer des tests statistiques

(Hopkins statistic, Distance distribution)

## ● Surveiller les données collectées avec les nouvelles fonctionnalités

(pour éviter les problèmes de data-drift ou de concept-drift)

Merci de m'avoir écouté, évalué et conseillé.

