

Projet 6 : Classifiez automatiquement des biens de consommation

Formation DATA SCIENTIST - Victor BARBIER

Sommaire

- Contexte et problématique
- Préparation des données
- Étude du texte
- Étude des images
- Faisabilité et conclusion

Contexte et problématique

Contexte et problématique

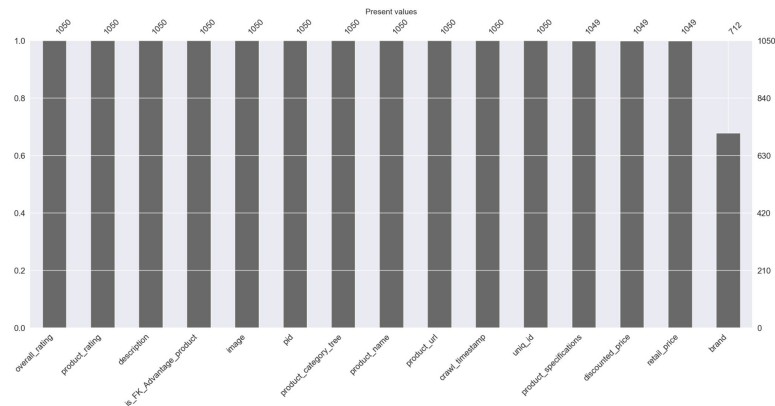
- Entreprise :
 - Place de Marché
 - Plateforme de e-commerce
 - Articles définis par une description et une image
- Mission :
 - Étude de faisabilité
 - Classification automatique
 - Basé sur une description et une image
- Méthode :
 - Réaliser un pré-traitement des descriptions et des images
 - Utiliser différents algorithmes d'analyse d'images et de texte
 - Effectuer un clustering et/ou des prédictions
 - Afficher une visualisation en 2D



Préparation des données

Préparation des données

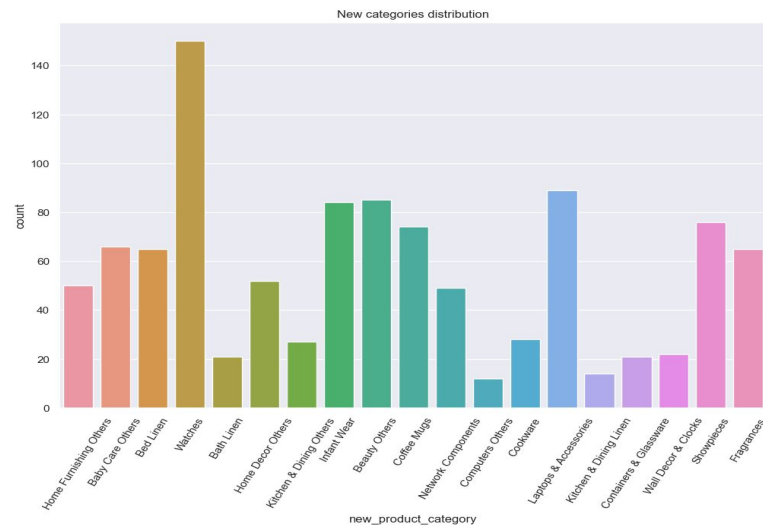
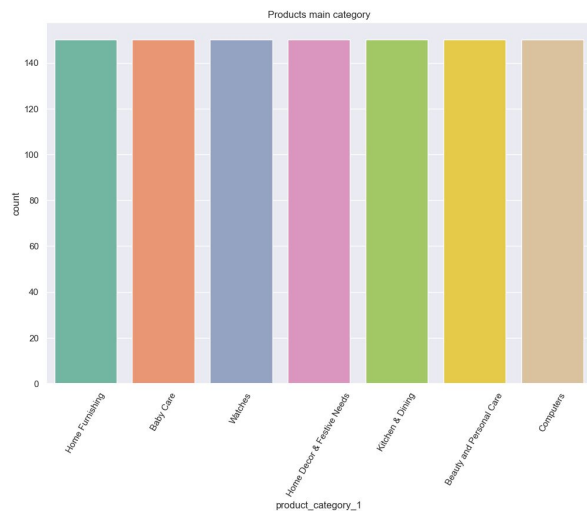
- 1050 articles
- 15 variables



VARIABLE	product_name	product_category_tree	description	image	...
EXEMPLE	Multicolor Door Curtain	Home Furnishing >> Curtains & Accessories >> Curtains	This curtain enhances the look of the interiors. This curtain is made from 100% high quality polyester fabric...	55b85ea15a1536d46b719oad6fff8ce7.jpg	...

Nouvelles catégories

- Division de la variable **product_category_tree**
- Création de nouvelles catégories pour affiner classification



Étude du texte

Algorithmes utilisés

Comptage simple	TF-IDF	fastText	BERT	USE
Bag of words	Bag of words	Word embedding	Sentence embedding	Sentence embedding
<p>Vecteurs de fréquence d'apparition des différents mots utilisés dans la description</p> <p>Unigrams</p> <p>Bigrams</p>	<p>Vecteurs de fréquence d'apparition des mots pondérés par leur importance dans le corpus</p> <p>poids=fréquence du terme×indicateur similarité</p>	<p>Librairie open-source</p> <p>Représente les mots en n-grams</p> <p>Matrice de vecteurs de mots</p> <p>Classification par logistic regression</p>	<p>Pré-entraîné</p> <p>Réseaux de transformers</p> <p>12 encodeurs</p> <p>512 tokens maximum</p> <p>Représentation des mots avec contexte</p>	<p>Pré-entraîné</p> <p>Réseaux de transformers</p> <p>Pas de longueur de texte limite</p> <p>512 dimensions en sortie</p>

Pré-traitement du texte

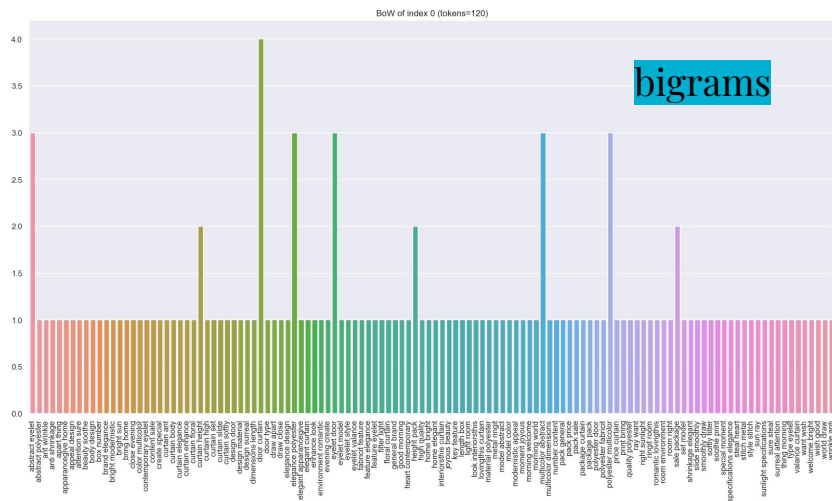
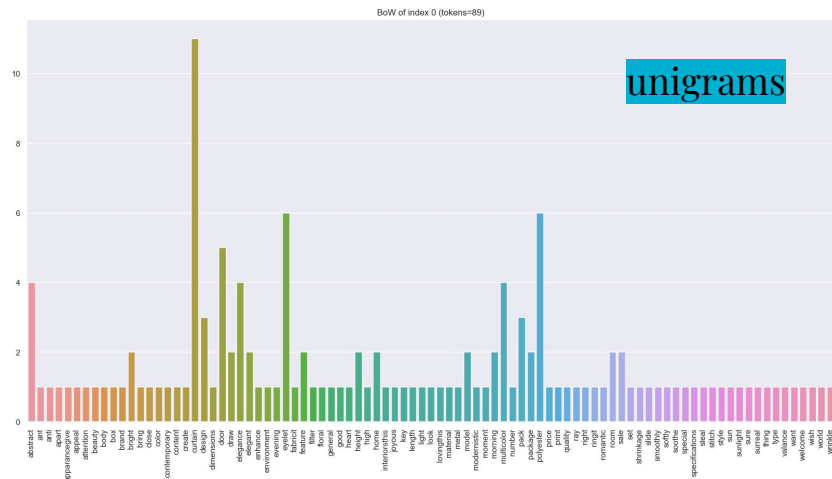
- Librairie **spaCy**
- Tokenization
- Suppression de la ponctuation et des nombres
- Suppression de certains mots :
 - Stopwords
 - Contiennent ponctuation / nombres
 - Noms de moins de deux lettres
- Lemmatization

Specifications of Eternity Handcrafted unique Mosaic Glass Table Lamp (35 cm, Blue) In The Box Sales Package 1 Table Lamp Number of Contents in Sales Package Pack of 1 General Type Buffet, Desk Color Blue Lamp Body Material Glass Lamp Base Material Metal Bulb Type 1x15 LED, CFSL ,Bulb Light Color Blue Assembly Required No Model Name Handcrafted unique Mosaic Glass Model Number Eternity001007 Power Features Power Requirement 110-240V, 50/60Hz Power Source Ac Additional Features Handcrafted lamps Dimensions Width 15 cm Height 35 cm Weight 1000 g

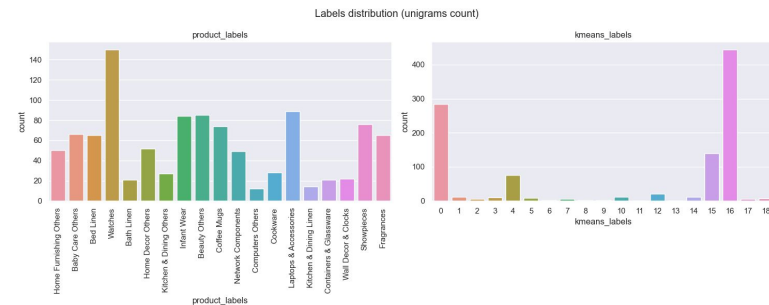
```
['specification', 'eternity', 'handcrafted',  
'unique', 'mosaic', 'glass', 'table', 'lamp',  
'blue', 'box', 'sales', 'package', 'table', 'lamp',  
'number', 'contents', 'sales', 'package', 'pack',  
'general', 'type', 'buffet', 'desk', 'color',  
'blue', 'lamp', 'body', 'material', 'glass',  
'lamp', 'base', 'material', 'metal', 'bulb',  
'type', 'led', 'cfs1', 'bulb', 'light', 'color',  
'blue', 'assembly', 'required', 'model',  
'handcraft', 'unique', 'mosaic', 'glass', 'model',  
'number', 'power', 'features', 'power',  
'requirement', 'power', 'source', 'additional',  
'features', 'handcrafted', 'lamp', 'dimensions',  
'width', 'height', 'weight']
```

Comptage simple

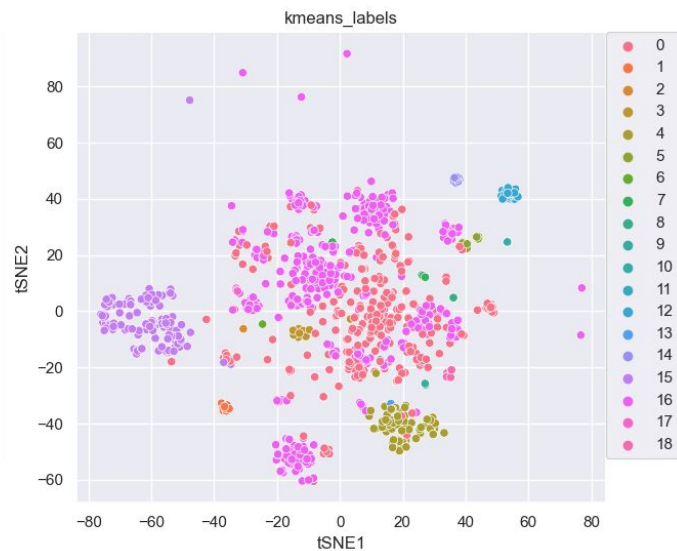
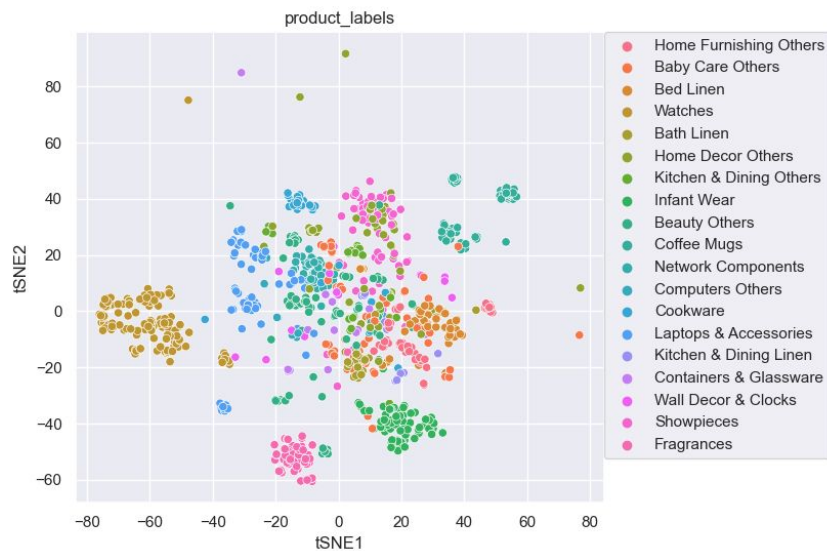
- Prétraitement du texte (**spaCy**)
- Features extraction :
 - **CountVectorizer** (*sklearn*)
 - Unigrams et bigrams
- Clustering avec **KMeans**
- Visualisation avec **t-SNE**



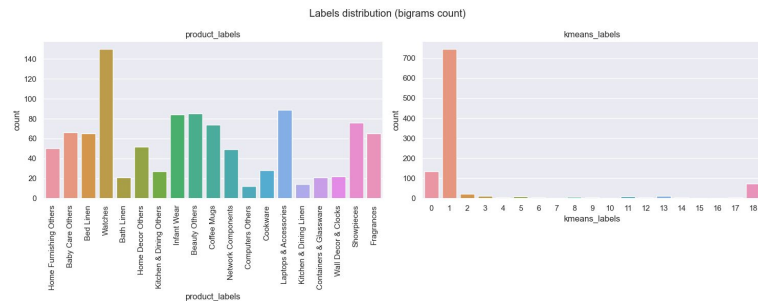
Comptage simple : unigrams



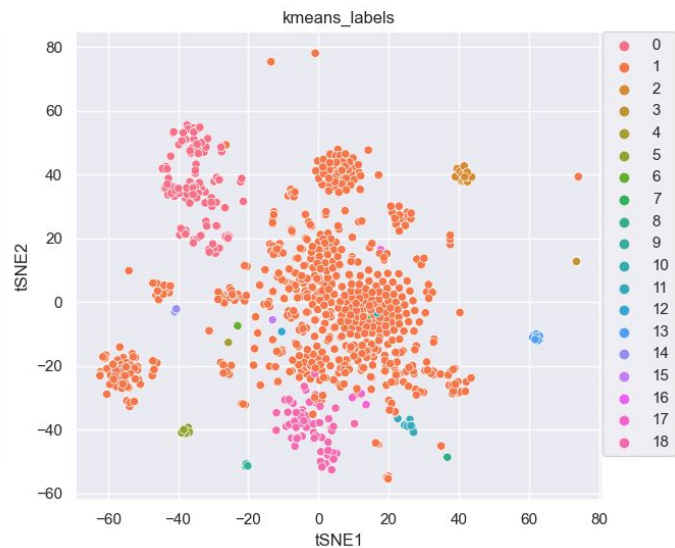
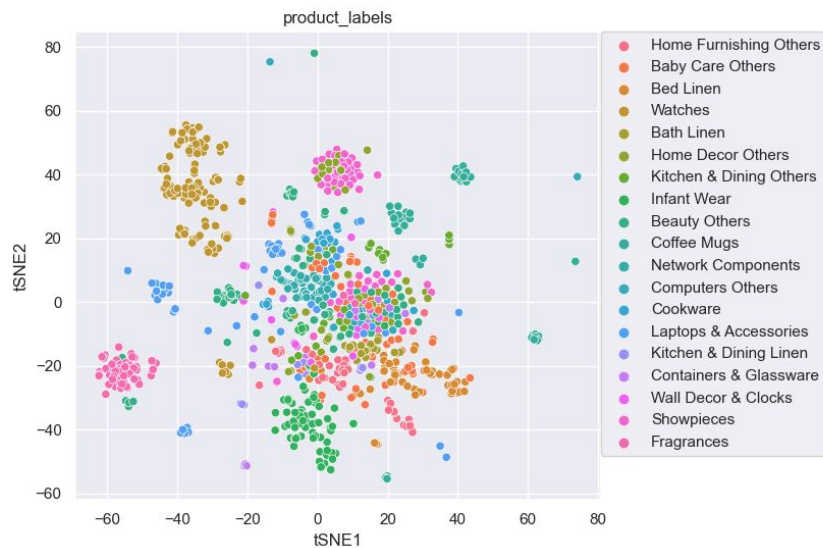
t-SNE visualizations (unigrams count)



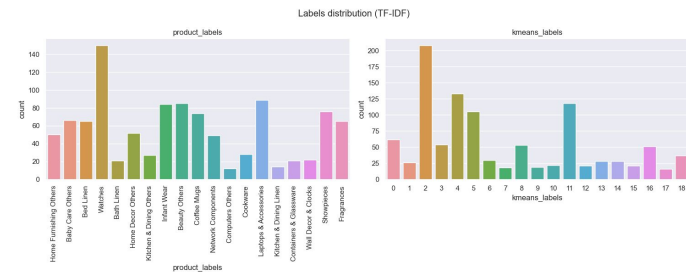
Comptage simple : bigrams



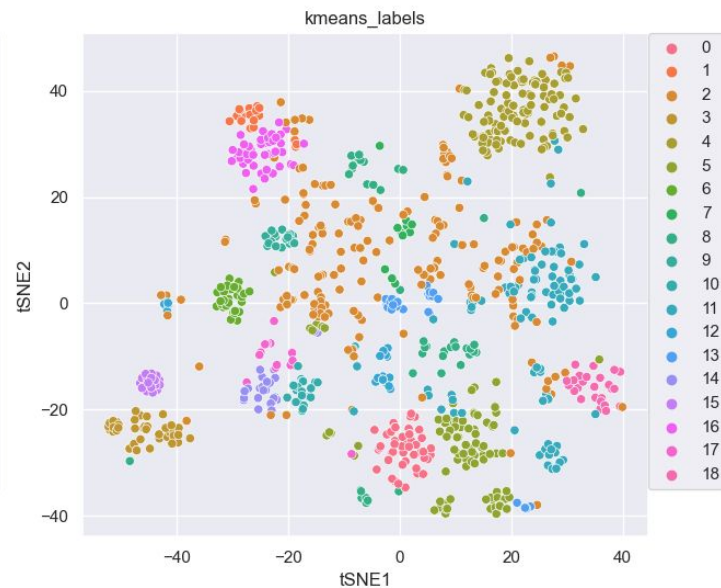
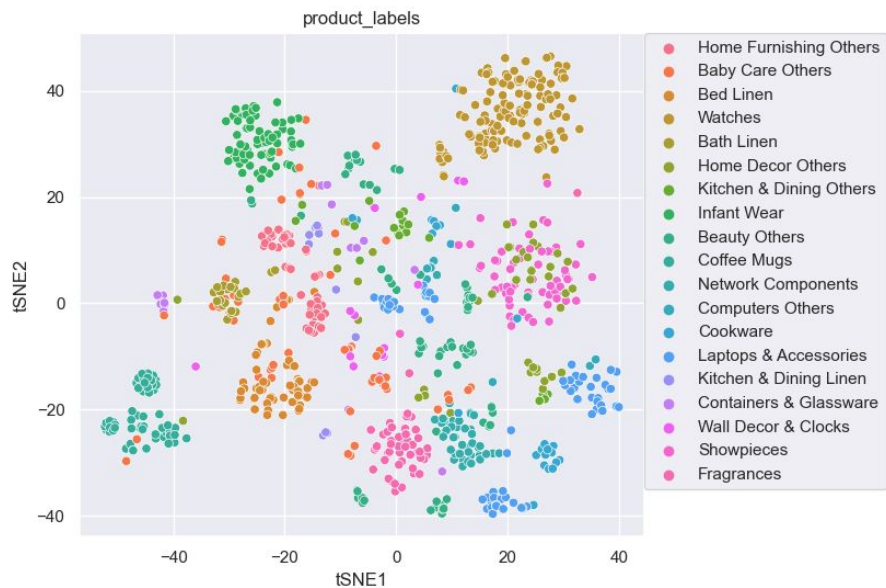
t-SNE visualizations (bigrams count)



TF-IDF



t-SNE visualizations (TF-IDF)



fastText

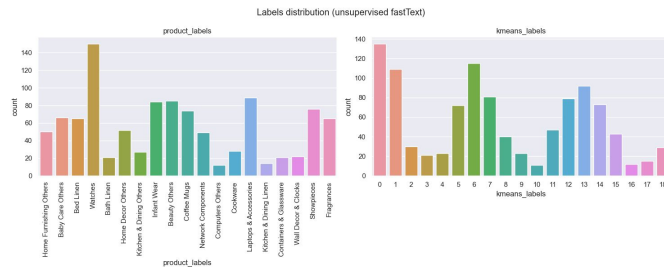
Non supervisé :

- Prétraitement du texte (**spaCy**)
- Features extraction :
 - Algorithme **fastText**
 - Words embedding
- Clustering avec **KMeans**
- Visualisation avec **t-SNE**

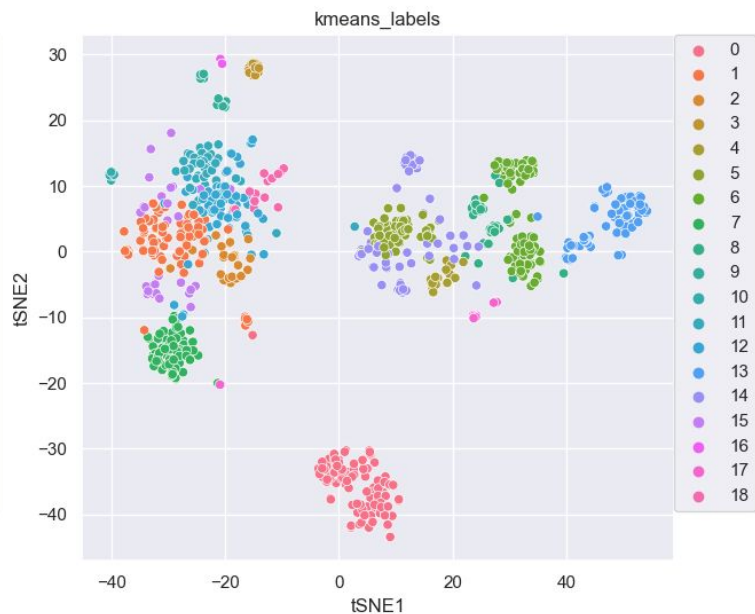
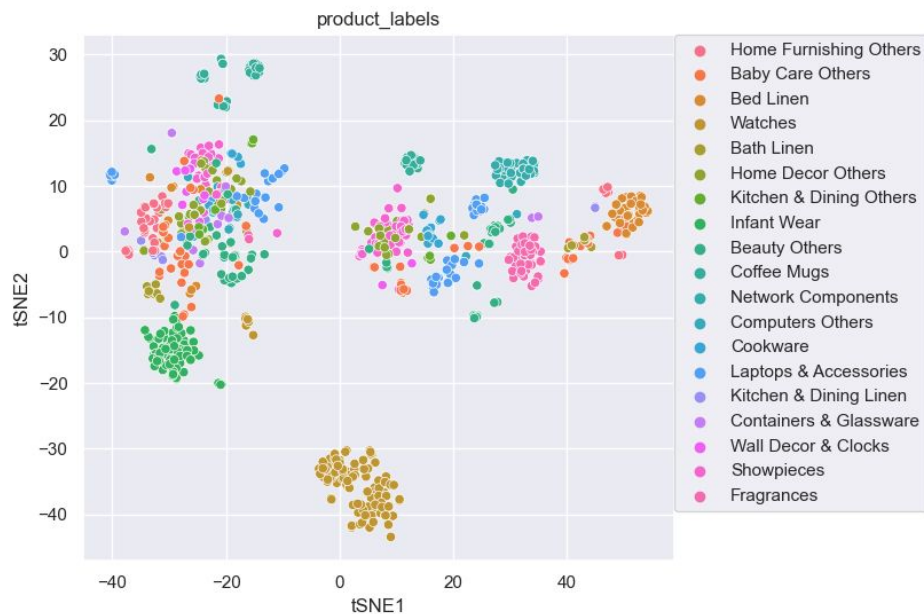
Supervisé :

- Prétraitement du texte (**spaCy**)
- Apprentissage **fastText**
sur jeu d'entraînement
(70% - 735 *produits*)
- Prédications sur jeu de test
(30% - 315 *produits*)
 - Features extraction
 - Prédiction des catégories
- Visualisation avec **t-SNE**

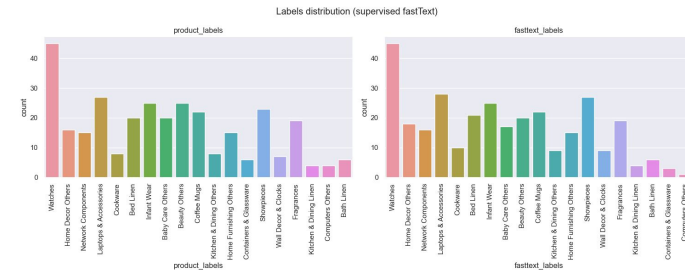
fastText : non supervisé



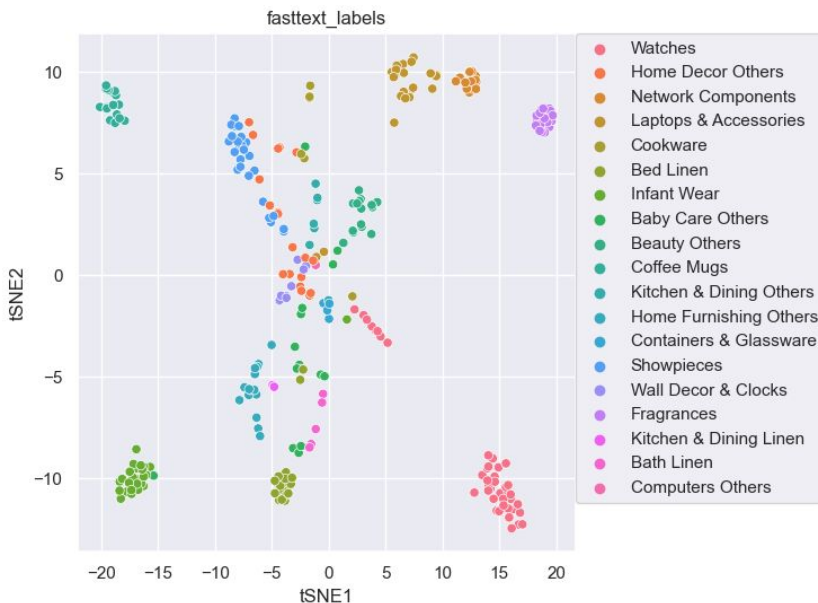
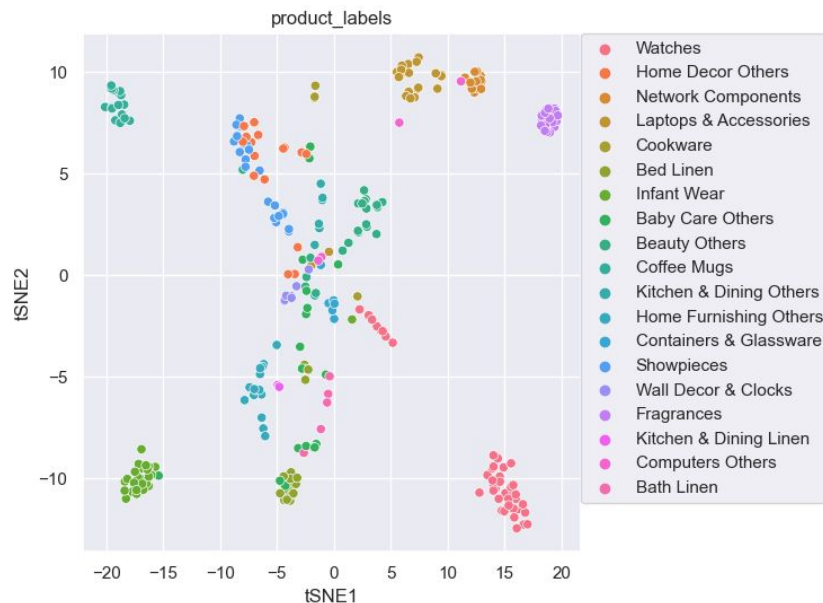
t-SNE visualizations (unsupervised fastText)



fastText : supervisé



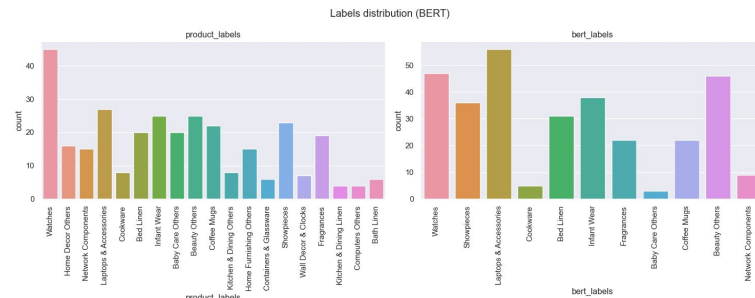
t-SNE visualizations (supervised fastText)



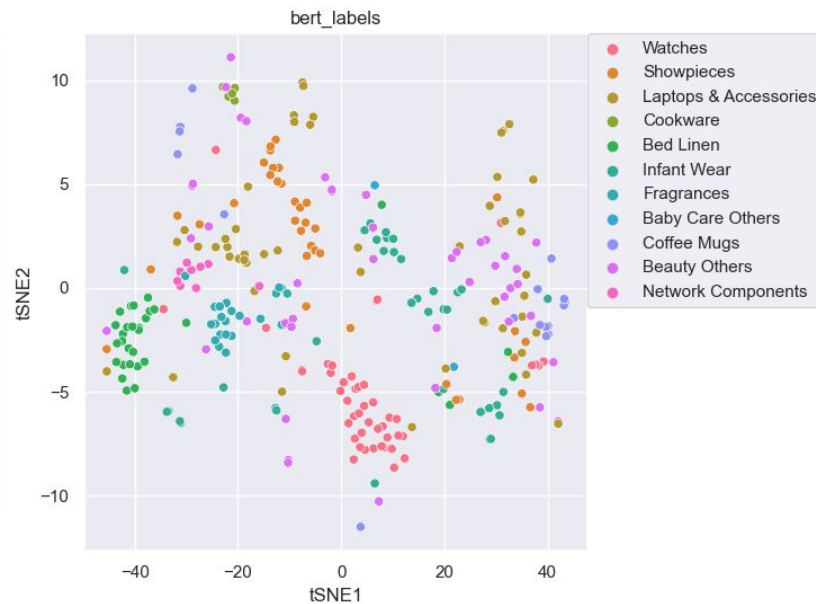
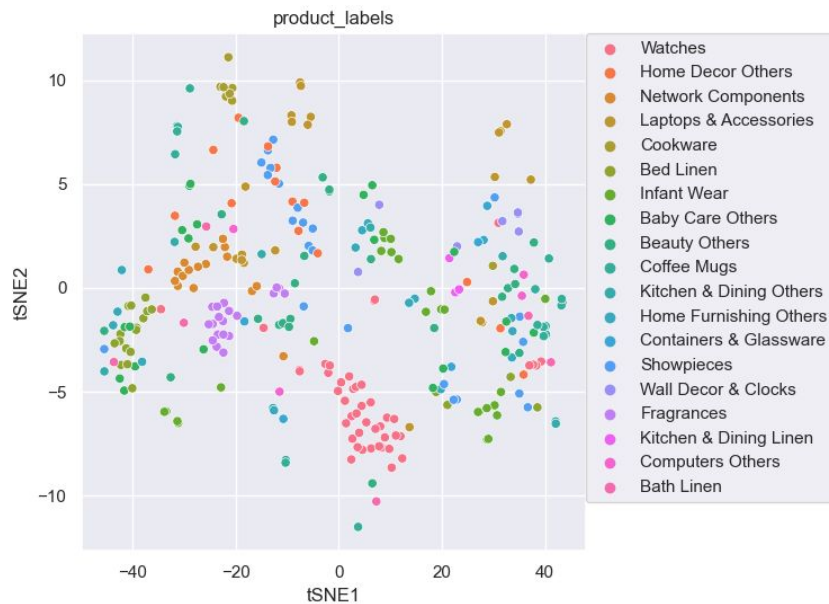
BERT

- Algorithme **BERT** :
 - Pre-trained model
- Fine tuning sur jeu d'entraînement
- Prédiction sur jeu de test
 - Features extraction
 - Prédiction des catégories
- Visualisation avec **t-SNE**

BERT



t-SNE visualizations (BERT)

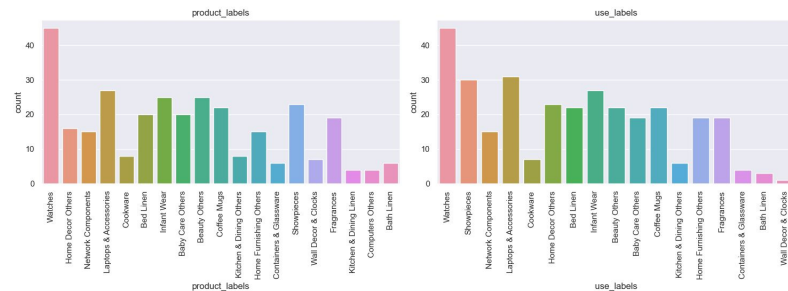


Universal Sentence Encoder

- Algorithme **USE** :
 - Pre-trained model
- Fine-tuning sur jeu d'entraînement
- Prédiction sur jeu de test
 - Features extraction
 - Prédiction des labels
- Visualisation avec **t-SNE**

Universal Sentence Encoder

Labels distribution (USE)



t-SNE visualizations (USE)



Étude des images

Algorithmes utilisés

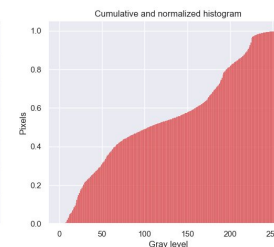
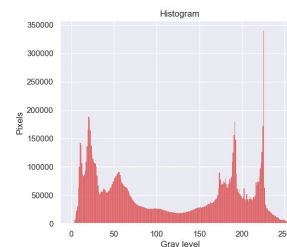
ORB	ResNet50
Algorithme OpenCV	Réseau neuronal convolutif (50 couches)
Détection des points d'intérêt	Pré-entraîné sur des millions d'images
Assigne des informations numériques (descriptors) aux points d'intérêt	Taille d'image : 224 par 224 pixels

ORB

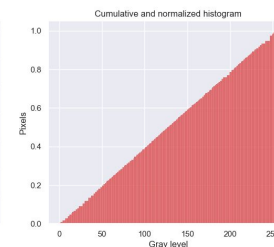
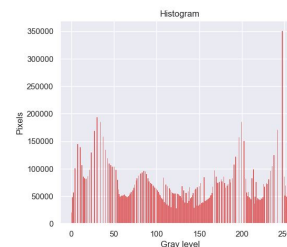
- Prétraitement des images
 - Histograms
- Algorithme **ORB** :
 - Détection des *descriptors*
- Clustering des *descriptors* pour définir des *visual words*
- Création de bag of visual words
- Clustering avec **KMeans** (*sklearn*)
- Visualisation avec **t-SNE**



55b85ea15a1536d46b7190ad6ff8ce7.jpg (3600, 3600)

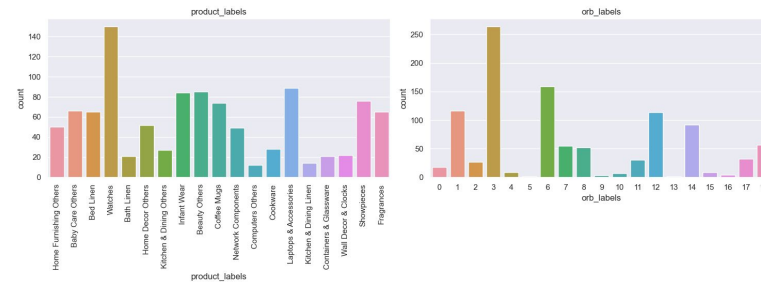


55b85ea15a1536d46b7190ad6ff8ce7.jpg (3600, 3600)

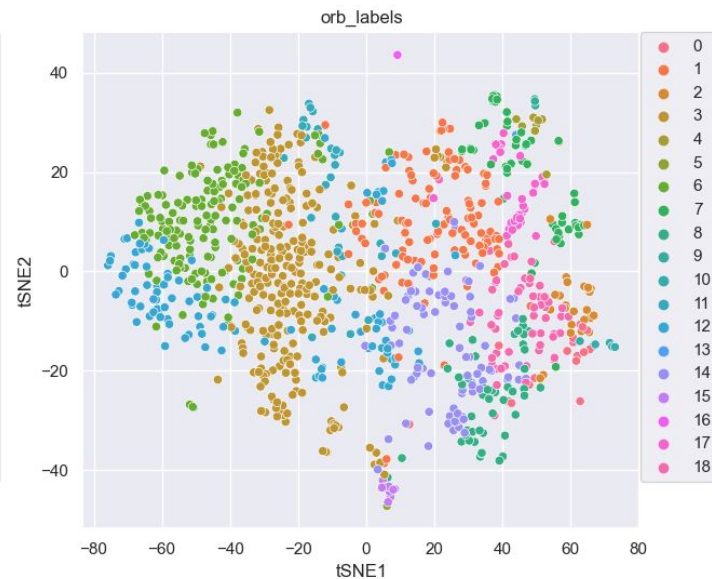
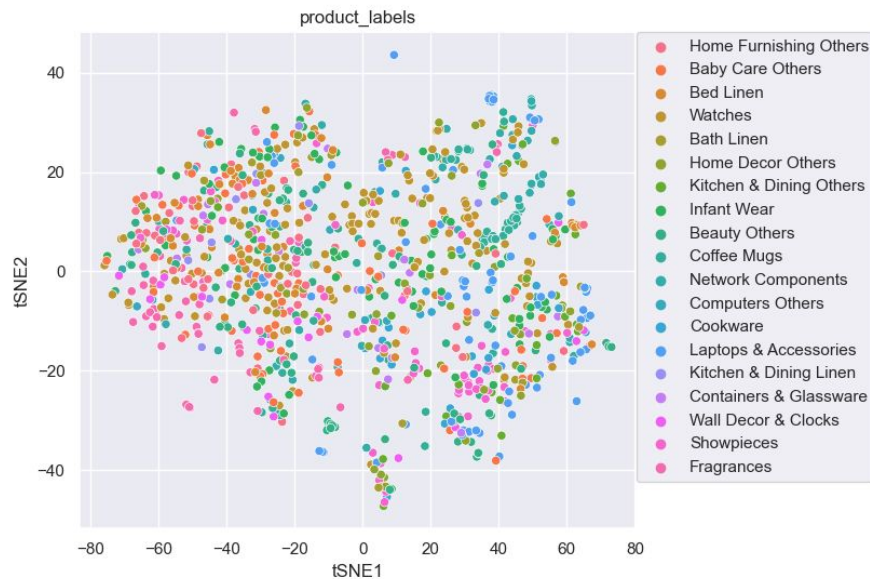


ORB

Labels distribution (ORB)



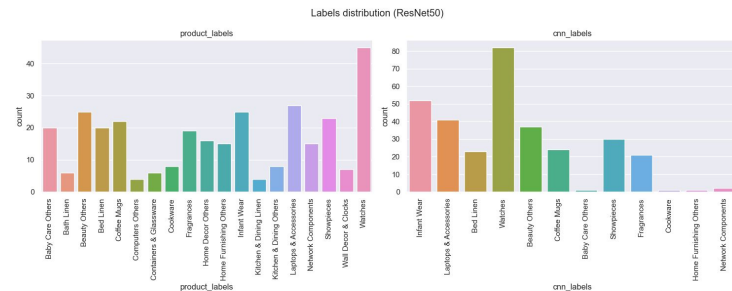
t-SNE visualizations (ORB)



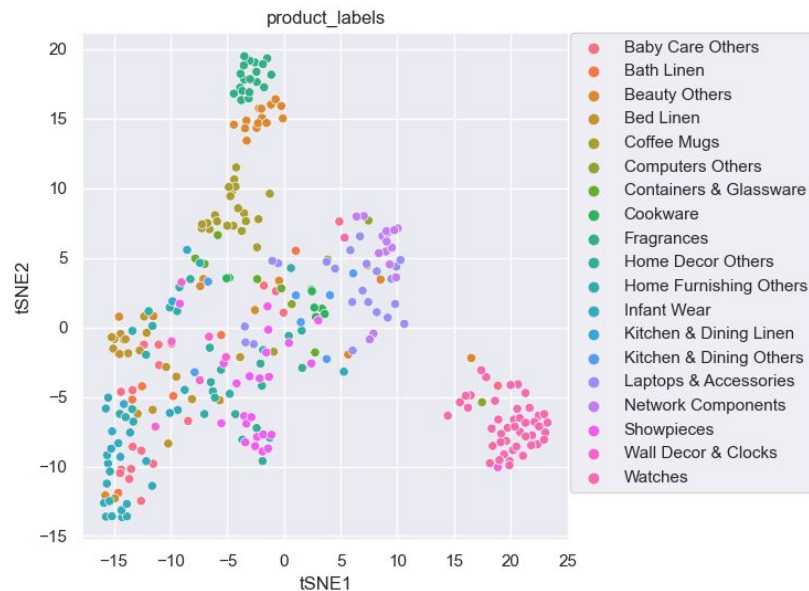
CNN transfer learning

- Algorithme **ResNet50** :
 - Pre-trained model
- Fine-tuning sur jeu d'entraînement
- Prédiction sur jeu de test
 - Features extraction
 - Prédiction des labels
- Visualisation avec **t-SNE**

CNN transfer learning



t-SNE visualizations (ResNet50)



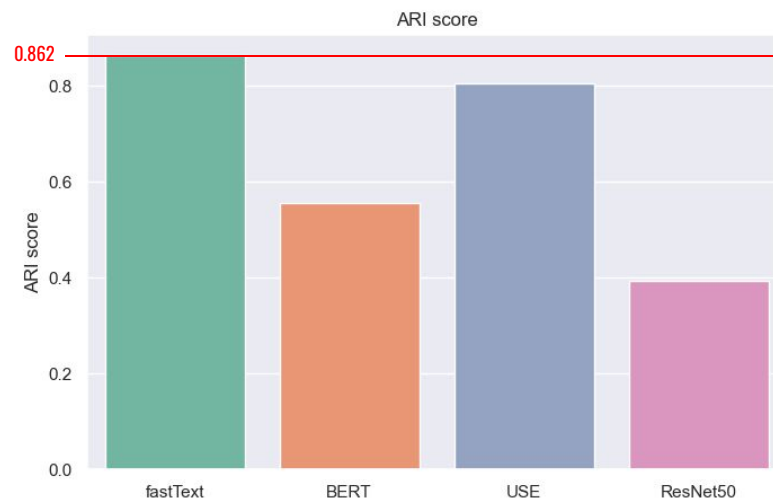
Faisabilité et conclusion

Comparaison des algorithmes : ARI score

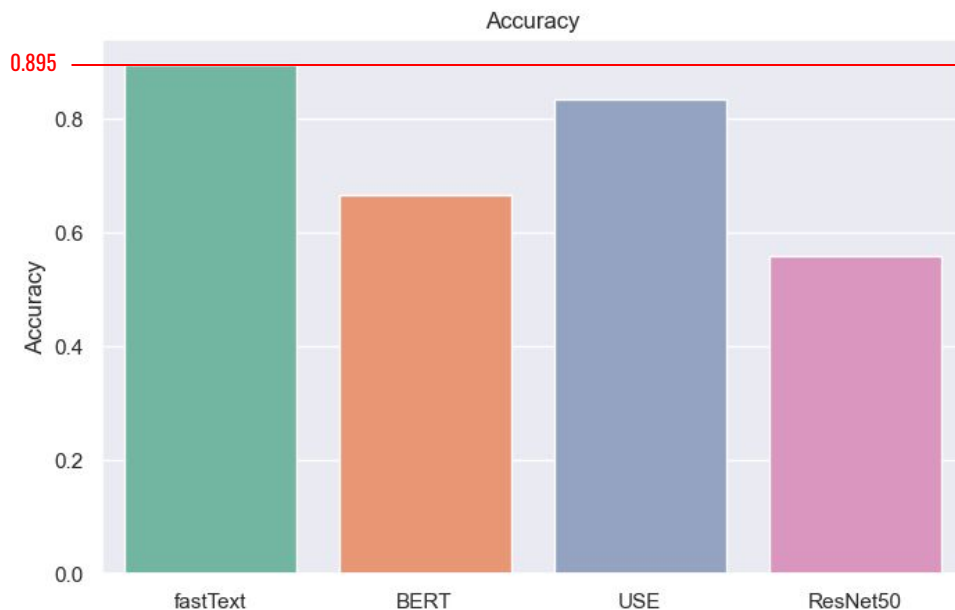
Non supervisé



Supervisé

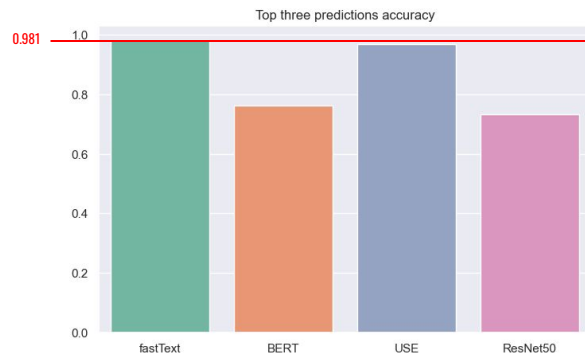


Comparaison des algorithmes : précision



Prédictions : probabilités d'appartenance à une catégorie

Précision en regardant les trois plus grandes probabilités :



Conclusion

- Utilisation des descriptions et images
- Réduction de dimensions pertinente
- Test de différentes approches
- ARI scores et précisions concluantes :
 - Un moteur de classification est faisable !

Améliorations :

- Manque de données d'entraînement
- Affiner la catégorisation (biaisée)
- Prédiction textes + images

Merci !
Avez-vous des
questions ?