

Projet 6 : Classifiez des biens de consommation

Lancelot LECLERCQ

7 mars 2022

Sommaire

1. Introduction
2. Classification des descriptions textuelles
3. Classification des images
4. Classification de l'ensemble des données
5. Conclusion

Introduction

Problématique

- L'entreprise Place de marché est un marketplace e-commerce
 - Vendeurs proposent des articles à des acheteurs en postant une photo et une description
 - Attribution de la catégorie d'un article effectuée manuellement par les vendeurs \Rightarrow peu fiable
- Objectif
 - Améliorer l'expérience utilisateur des vendeurs et des acheteurs
 - Automatisation de l'attribution d'une catégorie
 - Étude de la faisabilité d'un moteur de classification

Jeu de données

- Jeu de données textuelles
 - Nom, prix, description, note, pour chaque objet

Données

uniq_id
crawl_timestamp
product_url
product_name
product_category_tree
pid
retail_price
discounted_price
image
is_FK_Advantage_product
description
product_rating
overall_rating
brand
product_specifications

- Jeu d'images
 - Nous avons une image par objet



FIG. : Exemple d'image associée à un objet (ici des rideaux)

Classification des descriptions textuelles

Méthode : Nettoyage et création de "bag of words"

Sélection des stopwords personnalisés

- Nettoyage :

- Retrait des chiffres et caractères spéciaux,
- Retrait de la ponctuation,
- Uniformisation de la casse

- Tokenisation

- conservation des mots pertinents à partir de listes de "stopwords", des mots très récurrents à supprimer

- Lemmatisation

- Similaire à la tokenisation avec la suppression des terminaisons des mots
- Permet d'uniformiser les variations singulier/pluriel, masculin/féminin

- Racinisation (Stemmatisation)

- Similaire à la lemmatisation avec conservation de la racine des mots
- Permet d'uniformiser les variations de vocabulaires en regroupant les mots ayant les mêmes racines



Méthode : Nettoyage et création de "bag of words"

Comparaison des différents traitements du texte

TAB. : Tableau comparatif des différents procédés de nettoyage du texte

Modification	Contenu
Texte brut	Polyester Multicolor Abstract Eyelet Door Curtain (213 cm in Height, Pack of 2) Price : Rs. 899 This curtain enhances the look of the interiors.This curtain is made from 100%
Tokenisation	polyester multicolor abstract eyelet door curtain height pack price curtain enhances look curtain made
Lemmatisation	polyester multicolor abstract eyelet door curtain height pack curtain enhances look curtain made
Racinisation	polyest multicolor abstract eyelet door curtain height pack curtain enhanc look curtain made

Comparaison : validation croisée

- Essais pour les 3 types de nettoyages

- tokenisation
- lemmatisation
- racinisation

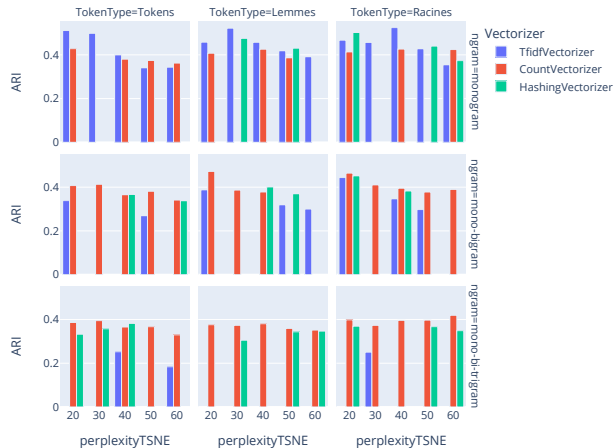
- Essais de trois vectoriseurs :

- TfidfVectorizer
- CountVectorizer
- HashingVectorizer

- Essais de différents découpages :

- mono-grams
- mono et bigrams
- mono, bi et trigrams

Comparaison des scores en fonction du vectoriseur et de la perplexité



Meilleure classification

- Meilleur score en utilisant :

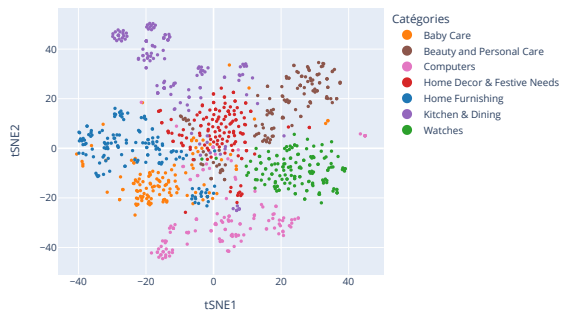
- les racines des mots
- des monograms
- le vectoriseur tf-idf

Matrice de confusion des labels prédits (x) et réels (y)
t-SNE40 TfidfVectorizer(1, 1) Racines

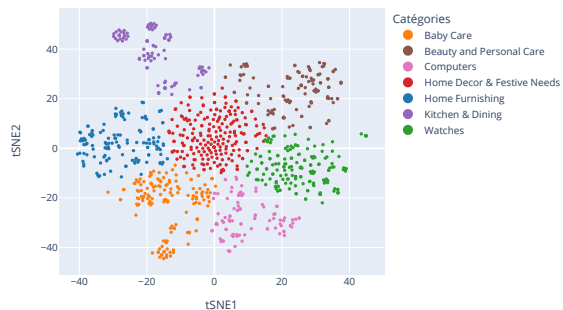
Catégorie réelle	Baby Care	86	3	44	1	19	0	0
	Beauty and Personal Care	5	96	0	7	0	15	0
	Computers	2	2	69	9	1	7	1
	Home Decor & Festive Needs	18	30	16	119	10	17	3
	Home Furnishing	14	0	1	5	110	0	0
	Kitchen & Dining	0	0	0	0	0	102	0
	Watches	0	8	7	1	0	1	146
		Baby Care Beauty and Personal Care Computers Home Decor & Festive Needs Home Furnishing Kitchen & Dining Watches						
		Catégorie prédite						

Visualisation de la meilleure classification

t-SNE40 TfidfVectorizer(1, 1) Racines



KMeans t-SNE40 TfidfVectorizer(1, 1) Racines

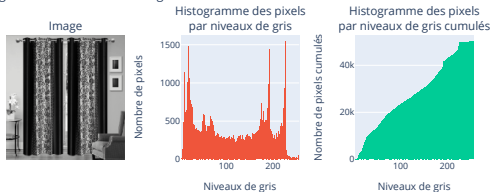


Classification des images

Méthode : Nettoyage

Traitement des images pour les algorithmes SIFT et ORB

Image convertie en niveaux de gris



- ← Utilisation de niveaux de gris
- ✓ Égalisation de l'histogramme
- ↓ Filtration de l'image

Image égalisée par CLAHE

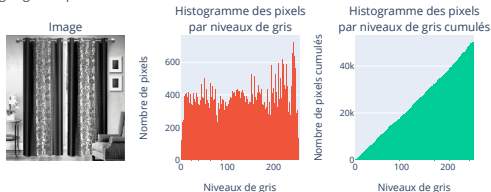
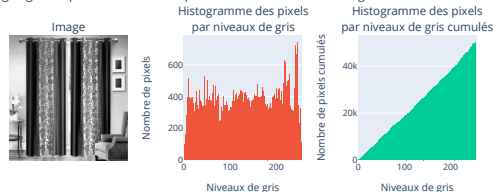


Image égalisée par CLAHE et filtrée par Non-local Means Denoising

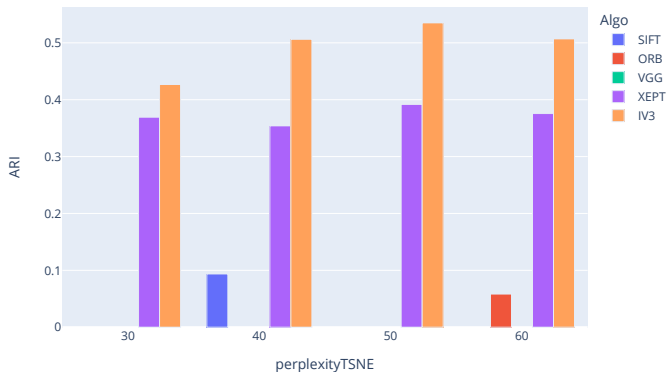


Comparaison : validation croisée

- Essais avec 5 algorithmes :

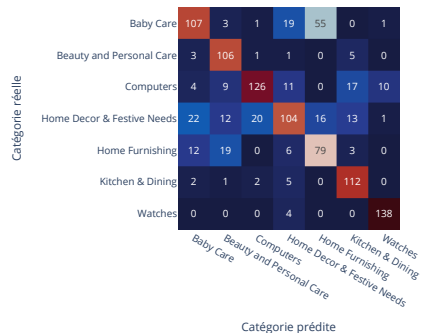
- SIFT
- ORB
- VGG16
- Xception
- InceptionV3

Comparaison des scores en fonction de l'algorithme et de la perplexité



Meilleure classification

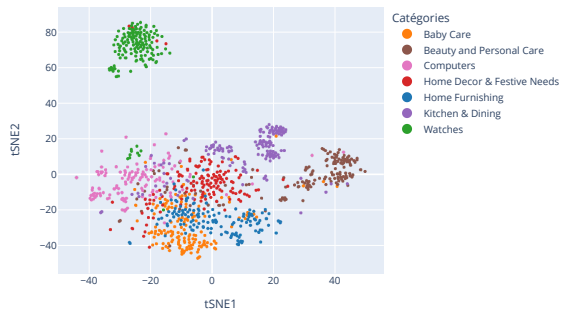
Matrice de confusion des labels prédits (x) et réels (y)
t-SNE40 IV3



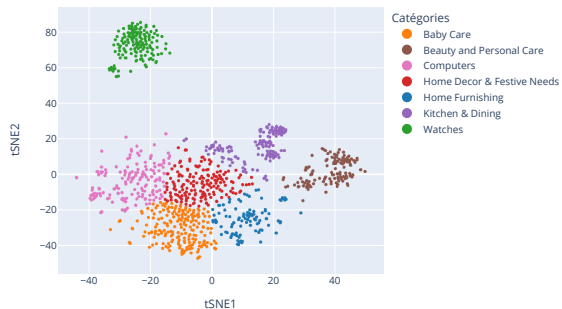
Meilleure classification avec : InceptionV3

Visualisation de la meilleure classification

t-SNE50 IV3



KMeans t-SNE50 IV3



Classification de l'ensemble des données

Méthode

- Combinaison de la matrice "bag of words" pour la meilleure classification du texte avec la matrice "bag of visual words" pour la meilleure classification des images
- Essais avec plusieurs valeurs de composantes utilisées pour la PCA afin d'affiner la classification

Conclusion