

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"

- 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

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

Classification des images

Classification de l'ensemble des données

Conclusion