

Projet 6: Classification automatique des biens de consommation

Maodo FALL

OpenClassrooms

Soutenance du projet
14 décembre 2024



Sommaire

- 1 Analyse exploratoire des données
 - Problématique et jeu de données
 - Traitement des données
- 2 Etude de la faisabilité d'un moteur de classification d'articles
 - Etude sur les données textuelles
 - Approches classiques
 - Approches Deep Learning
 - Etude sur les données d'image
 - Approche classique
 - Approche Deep Learning
- 3 Classification supervisée d'images
- 4 Collecte de nouveaux produits à base de "Champagne" via une API
- 5 Conclusion

Sommaire

- 1 **Analyse exploratoire des données**
 - Problématique et jeu de données
 - Traitement des données
- 2 Etude de la faisabilité d'un moteur de classification d'articles
 - Etude sur les données textuelles
 - Approches classiques
 - Approches Deep Learning
 - Etude sur les données d'image
 - Approche classique
 - Approche Deep Learning
- 3 Classification supervisée d'images
- 4 Collecte de nouveaux produits à base de "Champagne" via une API
- 5 Conclusion

Place de Marché



Problématique

Problématique :

"Place de Marché" souhaite lancer une marketplace e-commerce. Sur leur site, des vendeurs proposent des articles à des acheteurs en postant une photo et une description. Pour l'instant, la catégorisation d'un article est effectuée manuellement par les vendeurs, et est donc peu fiable. De plus, le volume des articles est pour l'instant très petit. L'entreprise a pour objectif d'automatiser cette tâche d'attribution de la catégorie.

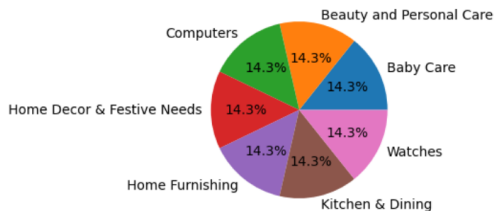
Missions :

- Etudier la faisabilité d'un moteur de classification.
- Classification supervisée d'images.
- Collecte de nouveaux produits via une API.

Jeu de données

- Le jeu de données contient 1050 produits et 15 colonnes.
- On se focalise plutôt sur la description et l'image des produits.
- La colonne des catégories est équilibrée.

Distribution of the categories



Nettoyage des données textes

Les données brutes de texte ont été passées sous plusieurs étapes de traitement pour les nettoyer et les traiter.

- nettoyage de texte.
- Suppression des mots les plus fréquents du vocabulaire.
- Ramener les tokens à leur racine.



Exemple de données d'image

Exemple d'images de la catégorie Watches :

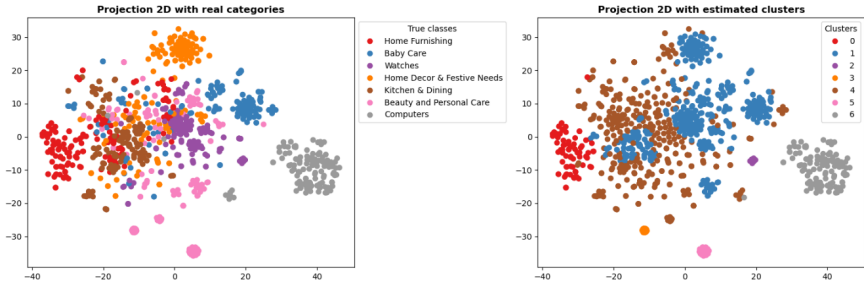


Sommaire

- 1 Analyse exploratoire des données
 - Problématique et jeu de données
 - Traitement des données
- 2 Etude de la faisabilité d'un moteur de classification d'articles
 - Etude sur les données textuelles
 - Approches classiques
 - Approches Deep Learning
 - Etude sur les données d'image
 - Approche classique
 - Approche Deep Learning
- 3 Classification supervisée d'images
- 4 Collecte de nouveaux produits à base de "Champagne" via une API
- 5 Conclusion

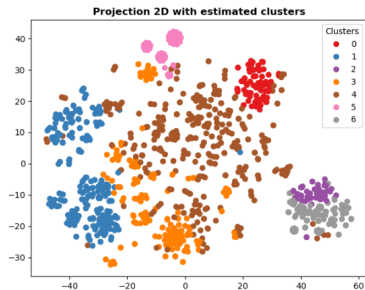
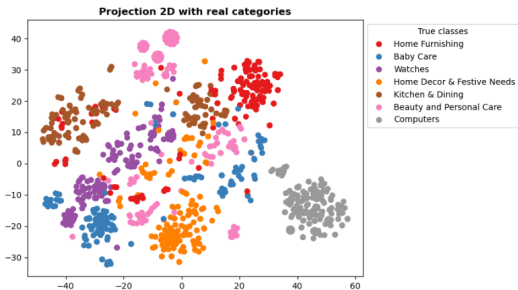
Approche Count-Vectorizer

Le score ARI de mesure de similarité est de : 0.162



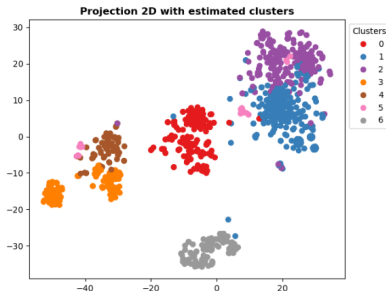
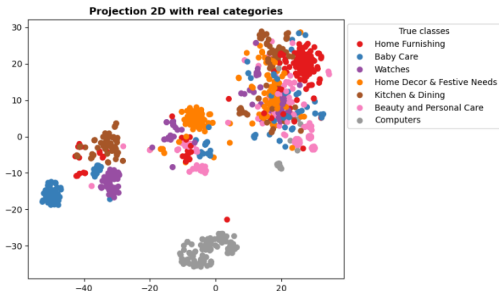
Approche TF-IDF

Le score ARI de mesure de similarité est de : 0.191



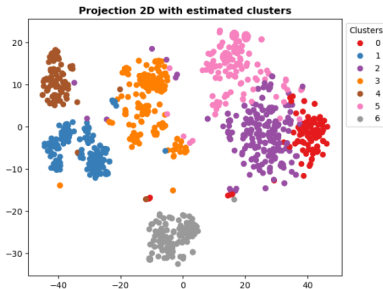
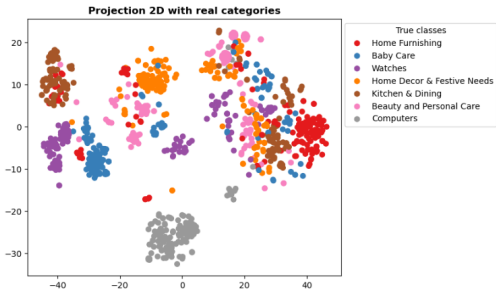
Approche Word2Vec

Le score ARI de mesure de similarité est de : 0.261



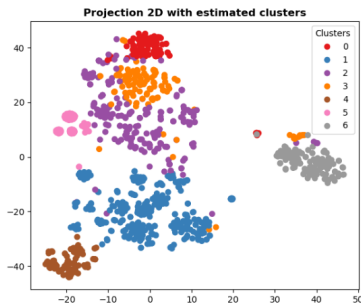
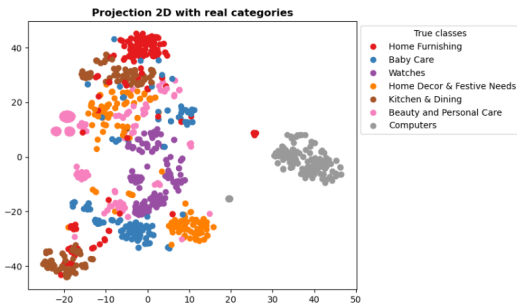
Approche BERT

Le score ARI de mesure de similarité est de : 0.294



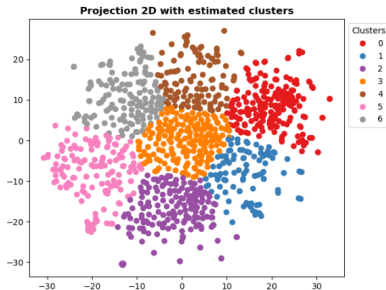
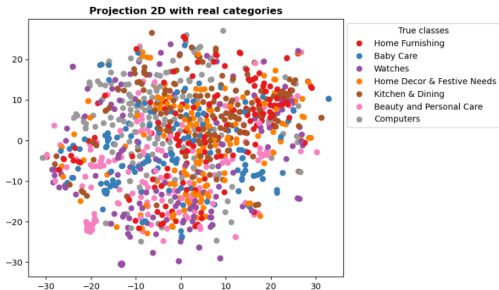
Approche USE

Le score ARI de mesure de similarité est de : 0.263



Approche SIFT

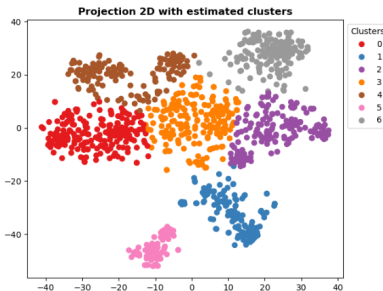
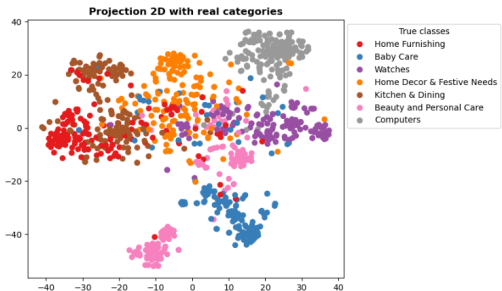
Le score ARI de mesure de similarité est de : 0.061



Approche ConvNet - Transfert

Learning avec VGG16

Le score ARI de mesure de similarité est de : 0.451



Sommaire

- 1 Analyse exploratoire des données
 - Problématique et jeu de données
 - Traitement des données
- 2 Etude de la faisabilité d'un moteur de classification d'articles
 - Etude sur les données textuelles
 - Approches classiques
 - Approches Deep Learning
 - Etude sur les données d'image
 - Approche classique
 - Approche Deep Learning
- 3 Classification supervisée d'images
- 4 Collecte de nouveaux produits à base de "Champagne" via une API
- 5 Conclusion

Approche simple par préparation initiale des images

Classification avec Transfert Learning avec VGG16 :

	precision	recall	f1-score	support
0	0.67	0.73	0.70	30
1	0.88	0.73	0.80	30
2	0.88	0.70	0.78	30
3	0.63	0.80	0.71	30
4	0.79	0.77	0.78	30
5	0.87	0.90	0.89	30
6	0.87	0.87	0.87	30
accuracy			0.79	210
macro avg	0.80	0.79	0.79	210
weighted avg	0.80	0.79	0.79	210

Supervised classification confusion matrix with VGG16 transfer learning

True classes	Home Furnishing	22	1	0	2	4	1	0
	Baby Care	0	22	1	3	2	0	2
	Watches	4	0	21	3	0	1	1
	Home Decor & Festive Needs	1	1	2	24	0	2	0
	Kitchen & Dining	5	0	0	2	23	0	0
	Beauty and Personal Care	1	0	0	1	0	27	1
	Computers	0	1	0	3	0	0	26
		0	1	2	3	4	5	6
		Predicted classes						

Modèle VGG16 avec data augmentation

	precision	recall	f1-score	support
0	0.62	0.70	0.66	30
1	0.84	0.70	0.76	30
2	0.62	0.67	0.65	30
3	0.67	0.67	0.67	30
4	0.62	0.43	0.51	30
5	0.75	0.90	0.82	30
6	0.84	0.90	0.87	30
accuracy			0.71	210
macro avg	0.71	0.71	0.70	210
weighted avg	0.71	0.71	0.70	210

Supervised classification confusion matrix with data augmentation

True classes	Home Furnishing	21	0	2	0	4	3	0
	Baby Care	0	21	3	1	2	0	3
	Watches	0	1	20	3	1	3	2
	Home Decor & Festive Needs	1	0	5	20	1	3	0
	Kitchen & Dining	11	1	1	4	13	0	0
	Beauty and Personal Care	1	0	1	1	0	27	0
	Computers	0	2	0	1	0	0	27
		0	1	2	3	4	5	6
		Predicted classes						

Sommaire

- 1 Analyse exploratoire des données
 - Problématique et jeu de données
 - Traitement des données
- 2 Etude de la faisabilité d'un moteur de classification d'articles
 - Etude sur les données textuelles
 - Approches classiques
 - Approches Deep Learning
 - Etude sur les données d'image
 - Approche classique
 - Approche Deep Learning
- 3 Classification supervisée d'images
- 4 Collecte de nouveaux produits à base de "Champagne" via une API
- 5 Conclusion

Résultats de l'appel de l'API food-database du site edamam

	foodId	label	category	foodContentsLabel	image
0	food_a656mk2a5dmqb2adiamu6beihduu	Champagne	Generic foods	None	https://www.edamam.com/food-img/a71/a718cf3c52...
1	food_b753ithamdb8psbt0v2k9aquo06c	Champagne Vinaigrette, Champagne	Packaged foods	OLIVE OIL; BALSAMIC VINEGAR; CHAMPAGNE VINEGAR...	None
2	food_b3dyababjo54xobm6r8jzbghjqe	Champagne Vinaigrette, Champagne	Packaged foods	INGREDIENTS: WATER; CANOLA OIL; CHAMPAGNE VINE...	https://www.edamam.com/food-img/d88/d88b64d973...
3	food_a9e0ghsamvoc45bwa2ybsa3gken9	Champagne Vinaigrette, Champagne	Packaged foods	CANOLA AND SOYBEAN OIL; WHITE WINE (CONTAINS S...	None
4	food_an4jjueaucpus2a3u1ni8auhe7q9	Champagne Vinaigrette, Champagne	Packaged foods	WATER; CANOLA AND SOYBEAN OIL; WHITE WINE (CON...	None
5	food_bmu5dmkazwuvpaa5prh1daa8js0	Champagne Dressing, Champagne	Packaged foods	SOYBEAN OIL; WHITE WINE (PRESERVED WITH SULFIT...	https://www.edamam.com/food-img/ab2/ab2459fc2a...
6	food_alpl44taoyv11ra0lic1qa8xculi	Champagne Buttercream	Generic meals	sugar; butter; shortening; vanilla; champagne;...	None
7	food_am5egz6aq3fpjlaf8xpdkbc2asis	Champagne Truffles	Generic meals	butter; cocoa; sweetened condensed milk; vanilla...	None
8	food_bcz8rhiajk1fuva0vkfmeakbouc0	Champagne Vinaigrette	Generic meals	champagne vinegar; olive oil; Dijon mustard; s...	None
9	food_a79xmnya6ftogreaeukbroa0thhh0	Champagne Chicken	Generic meals	Flour; Salt; Pepper; Boneless, Skinless Chicke...	None

Sommaire

- 1 Analyse exploratoire des données
 - Problématique et jeu de données
 - Traitement des données
- 2 Etude de la faisabilité d'un moteur de classification d'articles
 - Etude sur les données textuelles
 - Approches classiques
 - Approches Deep Learning
 - Etude sur les données d'image
 - Approche classique
 - Approche Deep Learning
- 3 Classification supervisée d'images
- 4 Collecte de nouveaux produits à base de "Champagne" via une API
- 5 Conclusion

Conclusion

- Le moteur de classification est bien réalisable au vu des résultats obtenus.
- L'approche CNN est l'algorithme qui obtient de meilleurs clusters et le meilleur score ARI, et elle fonctionne sur le principe de Transfert Learning. Le temps d'exécution de cette méthode reste néanmoins assez long.
- La classification supervisée avec data augmentation n'a pas permis d'améliorer les résultats de la classification sans augmentation.
- Les nouvelles données collectées pourront nous servir de données de test de notre modèle. Ce qui nous permettra d'utiliser l'ensemble du jeu de données initial pour entraîner le modèle.

MERCI BEAUCOUP !