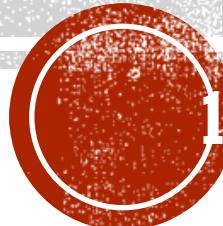


CLASSIFIEZ AUTOMATIQUEMENT DES BIENS DE CONSOMMATION

Etude de la faisabilité d'un moteur de classification



Formation Data Scientist Projet 6
Openclassrooms / Res Entreprise
Tulomba Divine
Octobre 2023

SOMMAIRE

- 1. Introduction, contexte & présentation de la problématique**
- 2. Analyse & préparation du jeu de données**
- 3. Faisabilité classification Texte (Bags of Words et Text Embedding)**
- 4. Faisabilité classification Image (Sift et Transfert Learning)**
- 5. Classification supervisée avec ou sans Data Augmentation**
- 6. Présentation test API**
- 7. Conclusion**



INTRODUCTION

- Objectif:

- une étude de faisabilité d'un moteur de classification
- classification supervisée et non supervisée à partir des images
- tester la collecte de produits à base de "champagne" via API

- Données:

- informations sur les produits disponibles sur le site e-commerce
- API

- Approche:

- Analyse
- Nettoyage
- Modèle
- Teste



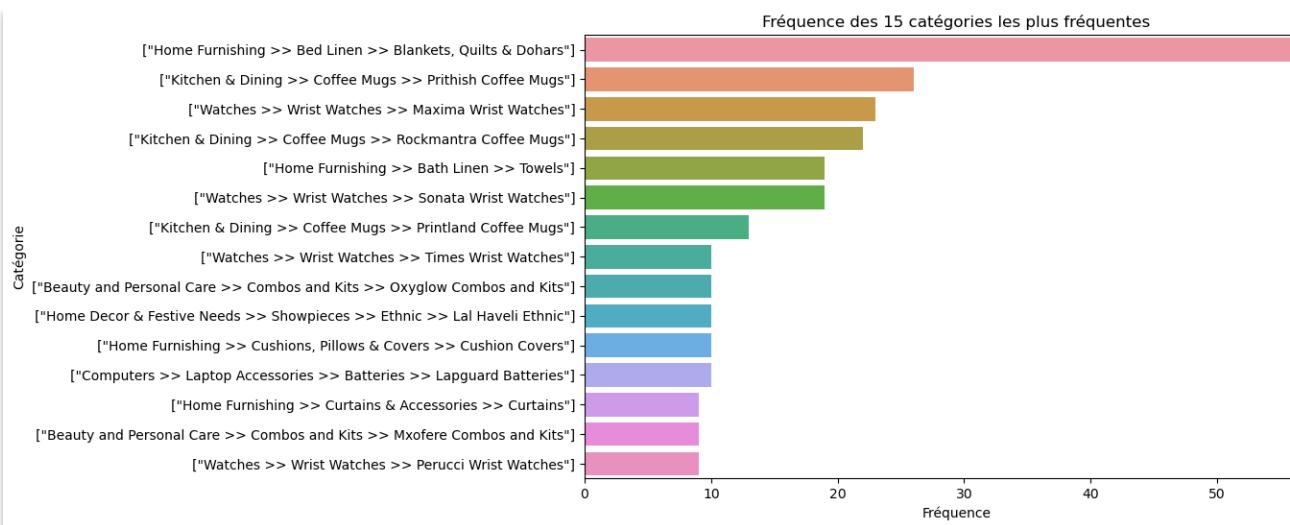
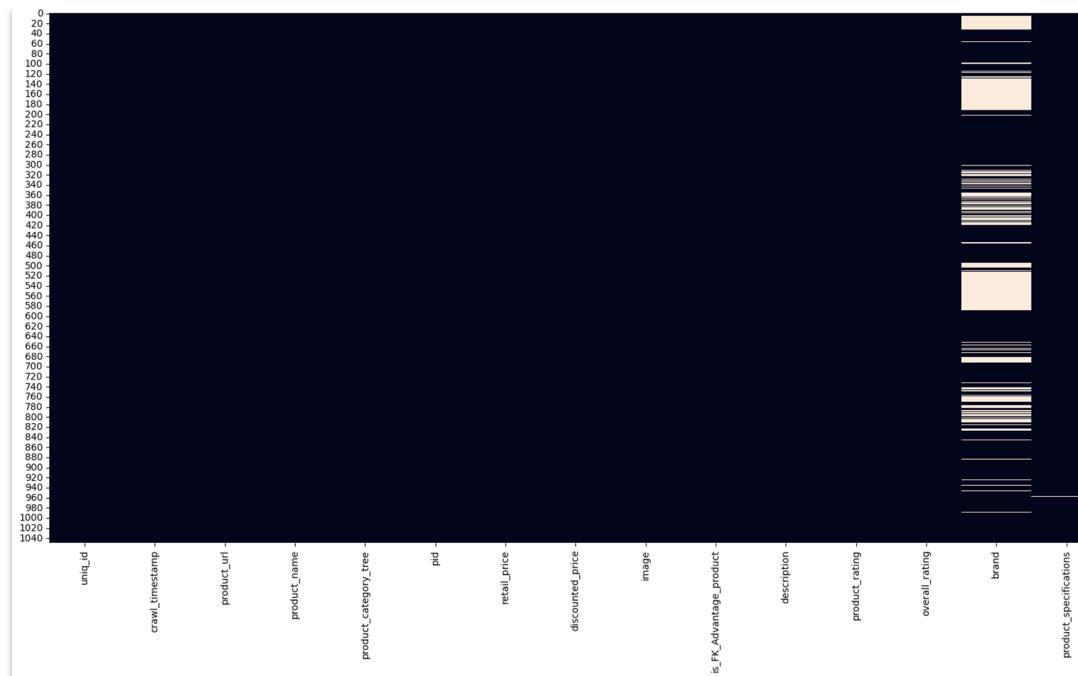
DESCRIPTION DU JEU DE DONNÉES

- **Contexte** : informations sur les produits disponibles sur le site e-commerce Flipkart
- **Nombre de produits** : 1050
- **Caractéristiques des produits** : des produits un identifiant unique, le nom du produit, la catégorie, la marque, le prix de détail, le prix réduit, etc.
- **Descriptif des colonnes** : 15 colonnes
- **Échantillon de données** :

	uniq_id	crawl_timestamp	product_url	product_name	product_category_tree	pid	retail_pric
0	55b85ea15a1536d46b7190ad6fff8ce7	2016-04-30 03:22:56 +0000	http://www.flipkart.com/elegance-polyester-mul...	Elegance Polyester Multicolor Abstract Eyelet ...	["Home Furnishing >> Curtains & Accessories >>..."]	CRNEG7BKMFFYHQ8Z	1899
1	7b72c92c2f8c40268628ec5f14c6d590	2016-04-30 03:22:56 +0000	http://www.flipkart.com/sathiyas-cotton-bath-t...l	Sathiyas Cotton Bath Towel	["Baby Care >> Baby Bath & Skin >> Baby Bath T..."]	BTWEGFZHGBXPHZUH	600
2	64d5d4a258243731dc7bbb1eeef49ad74	2016-04-30 03:22:56 +0000	http://www.flipkart.com/eurospa-cotton-terry-f...	Eurospa Cotton Terry Face Towel Set	["Baby Care >> Baby Bath & Skin >> Baby Bath T..."]	BTWEG6SHXTDB2A2Y	Na
3	d4884d0dc759dd90df41504698d737d8	2016-08-20 08:49:52 +0000	http://www.flipkart.com/santosh-royal-fashion-...	SANTOSH ROYAL FASHION Cotton Printed King size...	["Home Furnishing >> Bed Linen >> Bedsheets >>..."]	BOSEJT0UQWHDUBH4	2699
4	6325b6870c54cd47be8ebfbfa620ec7	2016-08-20 08:49:52 +0000	http://www.flipkart.com/jaipur-print-cotton-f...	Jaipur Print Cotton Floral King sized Double B...	["Home Furnishing >> Bed Linen >> Bedsheets >>..."]	BOSEJTHNGWVWWQU	2599

ANALYSE RAPIDE DU JEU DE DONNÉES

	colonne	valeurs manquantes	pourcentage
13	brand	338	32.190476
6	retail_price	1	0.095238
7	discounted_price	1	0.095238
14	product_specifications	1	0.095238
0	uniq_id	0	0.000000
1	crawl_timestamp	0	0.000000
2	product_url	0	0.000000
3	product_name	0	0.000000
4	product_category_tree	0	0.000000
5	pid	0	0.000000
8	image	0	0.000000
9	is_FK_Advantage_product	0	0.000000
10	description	0	0.000000
11	product_rating	0	0.000000
12	overall_rating	0	0.000000



Base de données au niveau de texte :

- Beaucoup de données manquantes au niveau de la variable « Brand » et peu sur product_specifications
- Les catégories les plus présentes dans notre base sont les produits de maison et literie.



Watches



Home Decor & Festive Needs



Home Furnishing



Baby Care



Kitchen & Dining



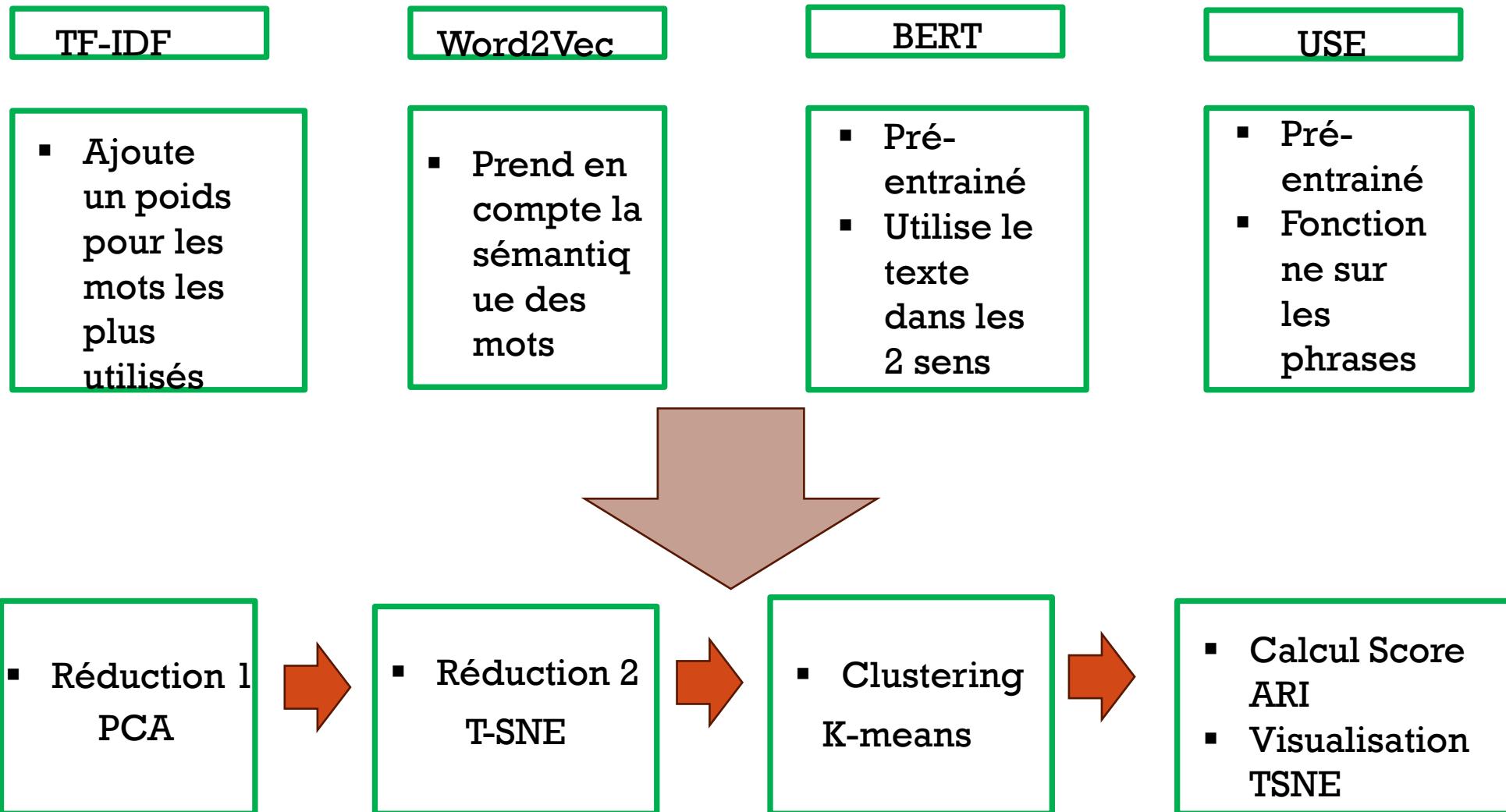
Computers



Beauty and Personal Care



CLASSIFICATION DESCRIPTION : FEATURES EXTRACTION



3. FAISABILITÉ DE CLASSIFICATION – TEXTE

NETTOYAGE DES FEATURES TEXTUELLES

Tokenisation :

```
# Exemple d'utilisation
example_sentence = "Le chat dort."
tokenized_sentence = tokenize_sentence(example_sentence)
tokenized_sentence

['Le', 'chat', 'dort', '.']
```

Stemming / Lemmatisation:

```
# Exemple d'utilisation
example_sentence = "The cats are running and jumping."
stemmed_sentence = stem_sentence(example_sentence)
stemmed_sentence

'the cat are run and jump .'
```

Suppression ponctuation / conversion minuscule :

```
# Application de la fonction de nettoyage à la colonne de description
df['cleaned_description'] = df['description'].apply(clean_text)
# Application de la fonction de nettoyage à la colonne de nom du produit
df['cleaned_product_name'] = df['product_name'].apply(clean_text)

# Affichage des résultats
df[['description', 'cleaned_description', 'product_name', 'cleaned_product_name']].head()
```

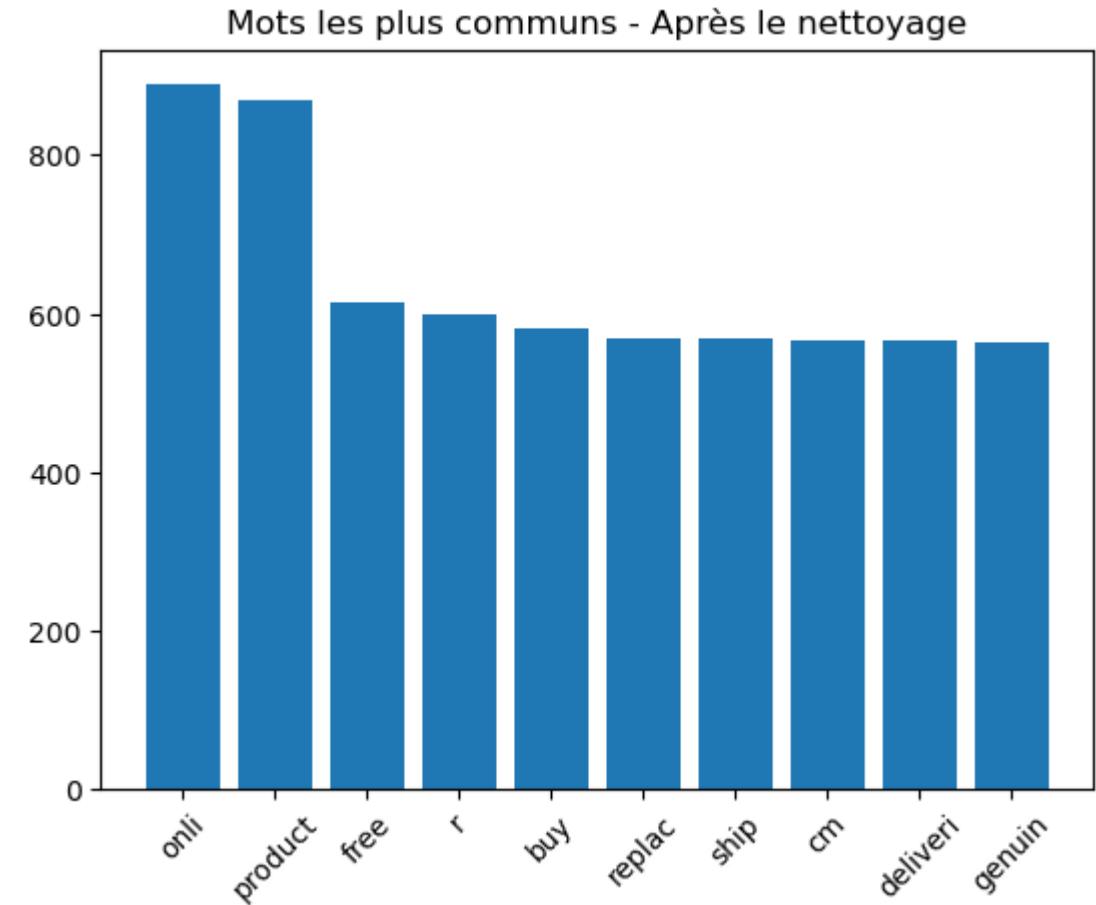
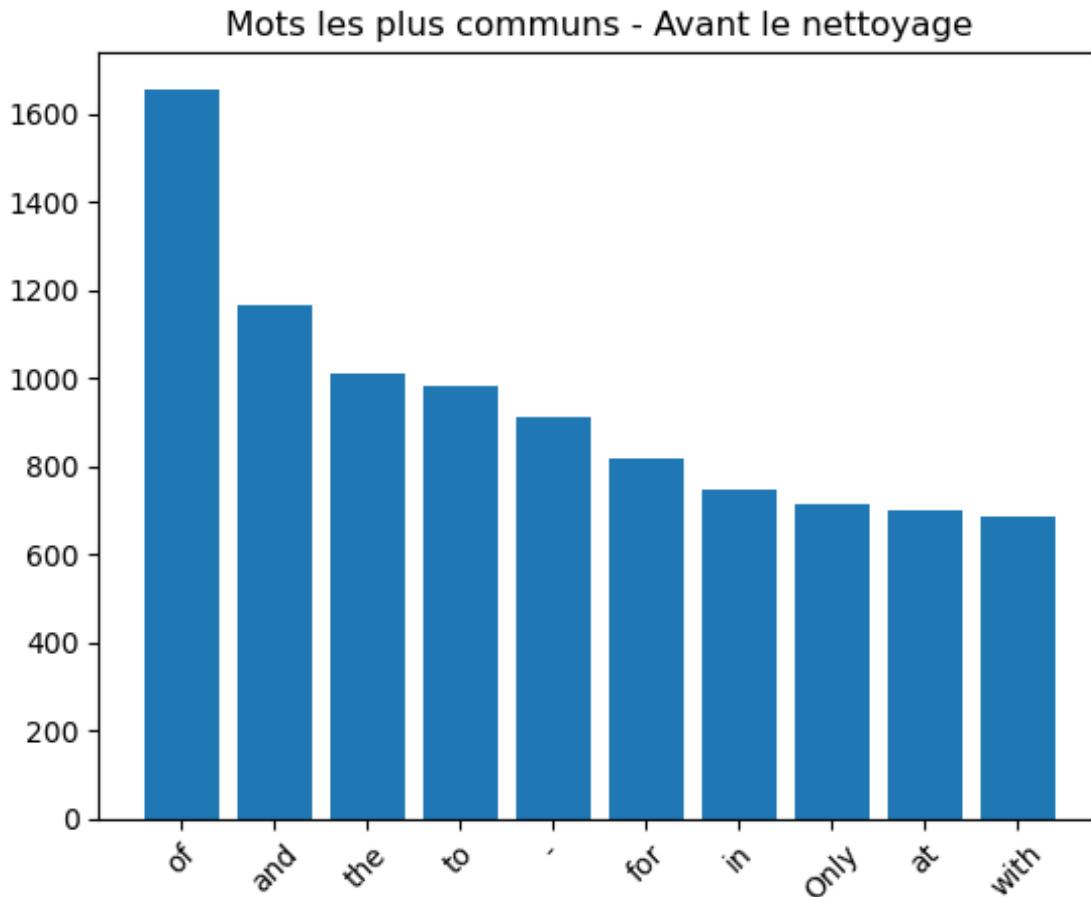


	description	cleaned_description
0	Key Features of Elegance Polyester Multicolor ...	key featur eleg polyest multicolor abstract ey...
1	Specifications of Sathiya Cotton Bath Towel (...	specif sathiya cotton bath towel 3 bath towel...
2	Key Features of Eurospa Cotton Terry Face Towel...	key featur eurospa cotton terri face towel set...
3	Key Features of SANTOSH ROYAL FASHION Cotton P...	key featur santosh royal fashion cotton print ...
4	Key Features of Jaipur Print Cotton Floral King ...	key featur jaipur print cotton floral king siz...



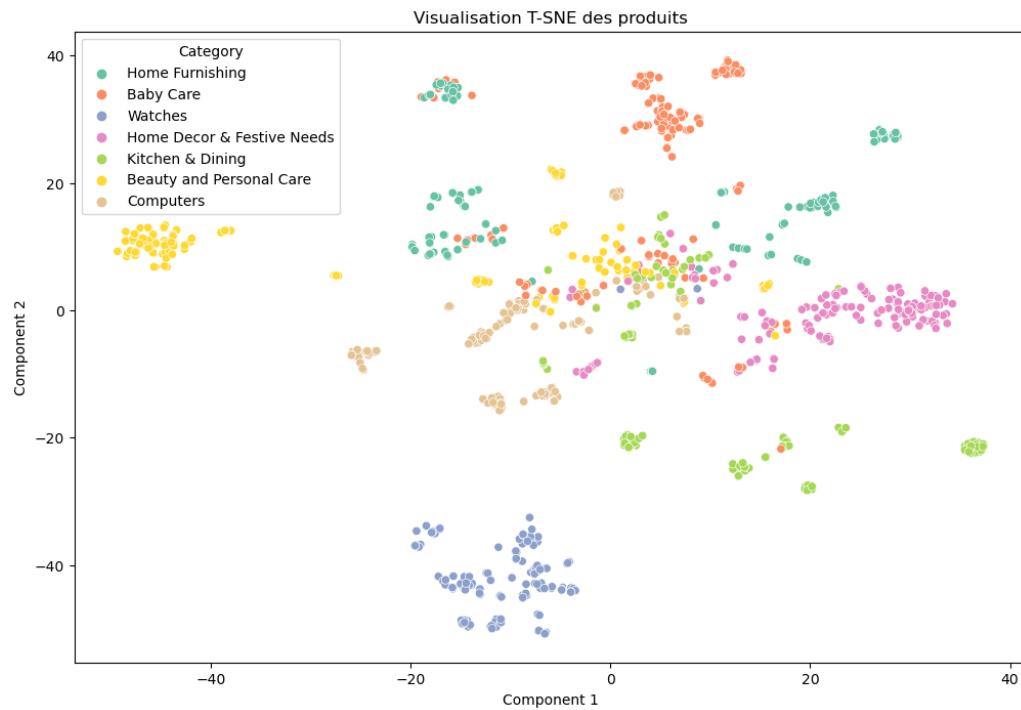
AVANT / APRÈS NETTOYAGE

Il n'y a plus de stopwords, mais uniquement des mots liés aux descriptions des biens



Faisabilité de classification – Texte Bag-of-Words

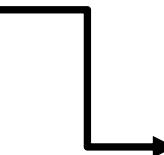
□ Fit et Transform sur « product_name » + « description » TfidfVectorizer



Adjusted Rand Score: 0.3988821944927195

MÉTHODOLOGIE ADOPTÉE POUR LA CLASSIFICATION NON SUPERVISÉE DE DONNÉES TEXTUELLES

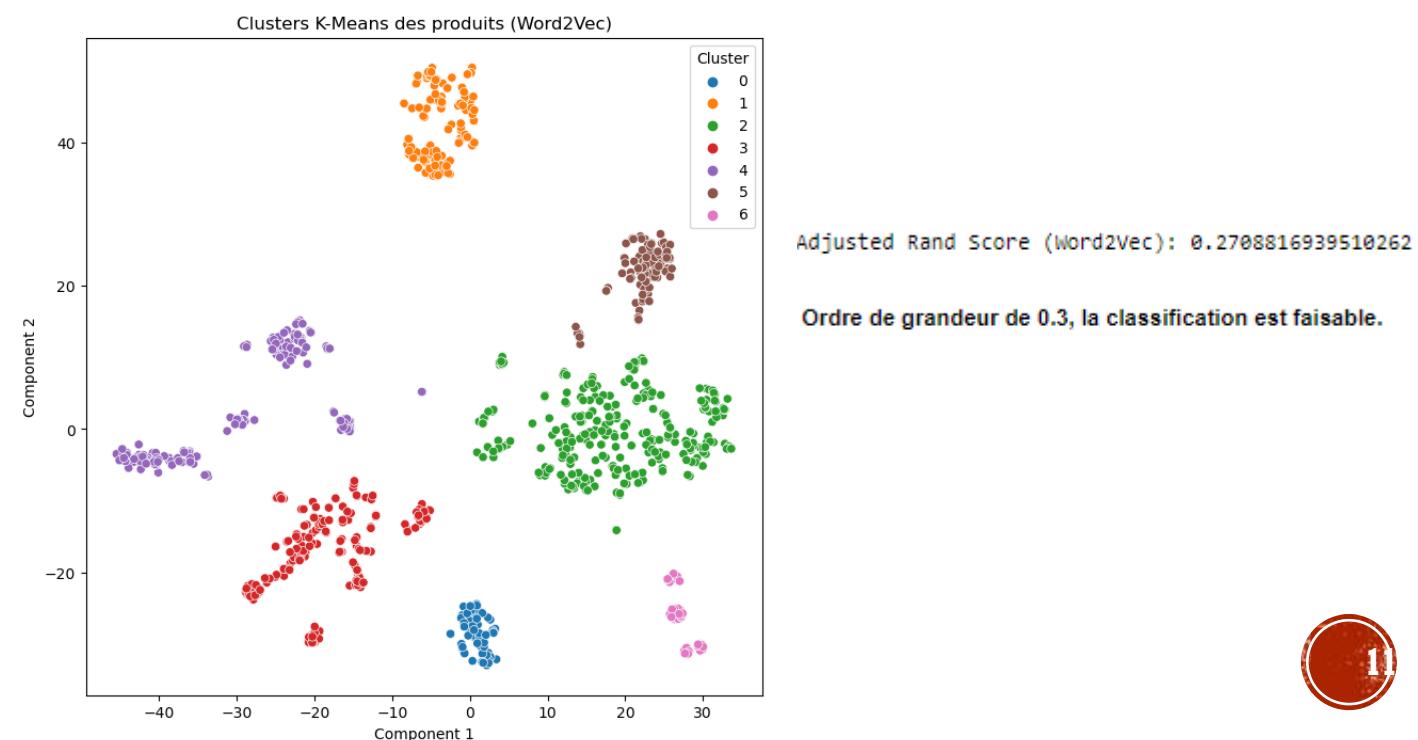
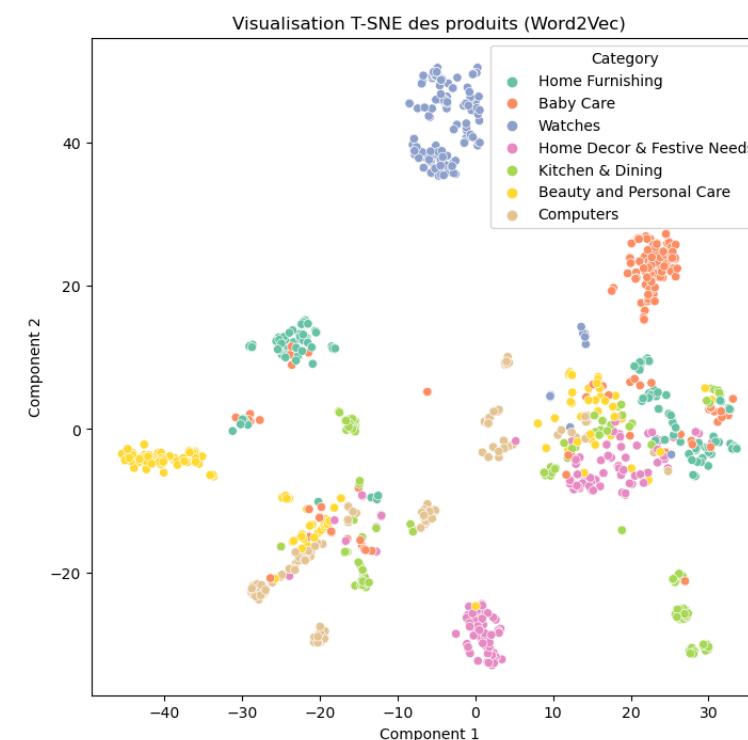
TEXT EMBEDDING



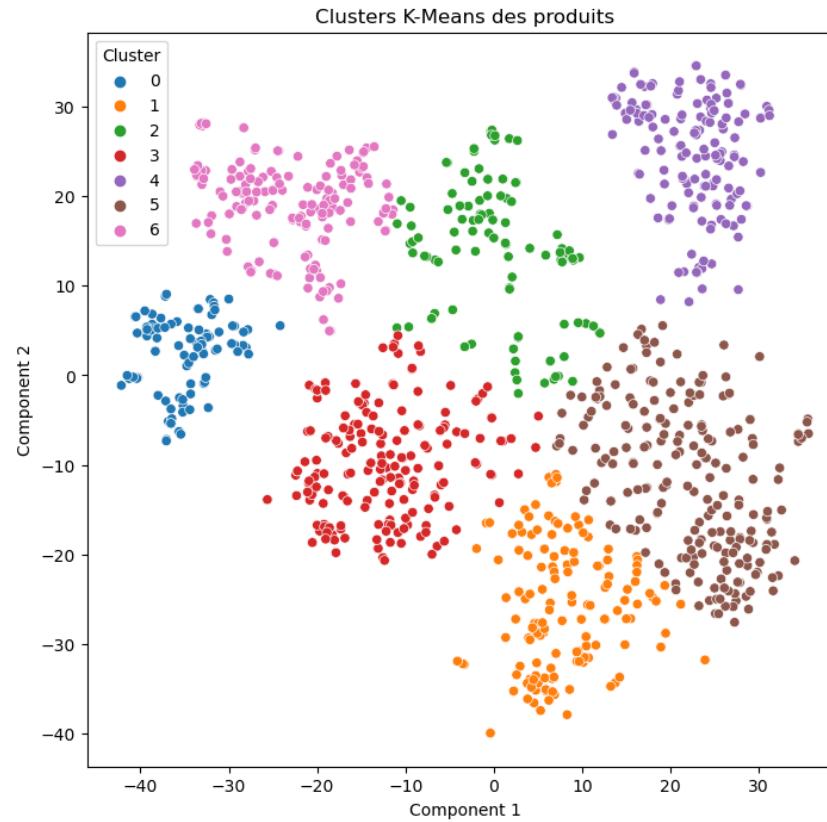
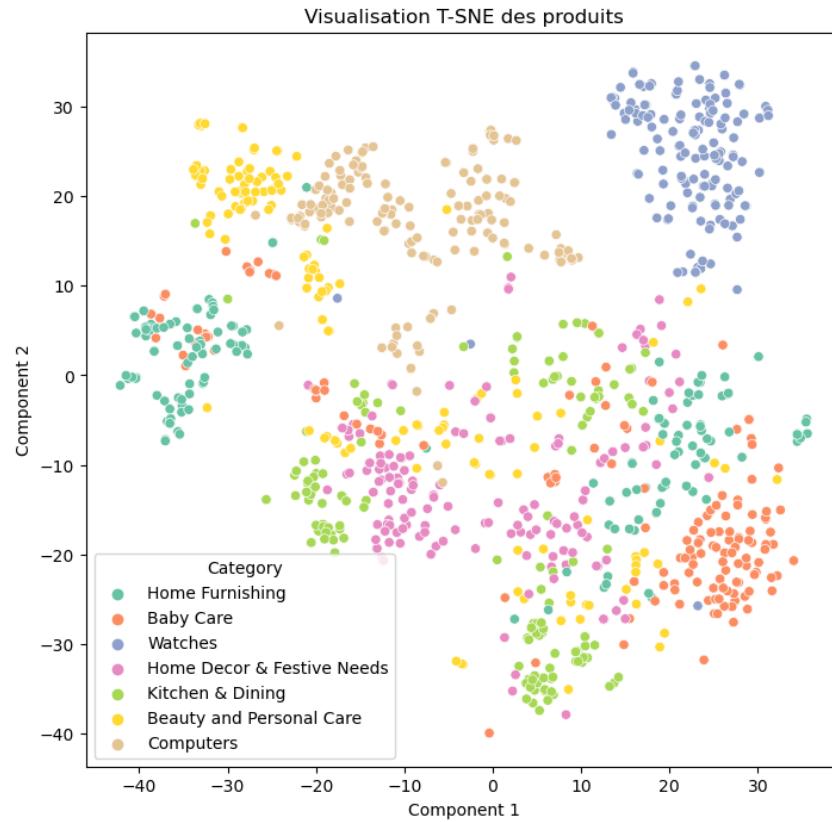
transformer du text brut non structuré et difficile à traiter par les modèles d'IA, en représentations structurées

3 outils : Word2Vec / BERT / USE

Word2Vec



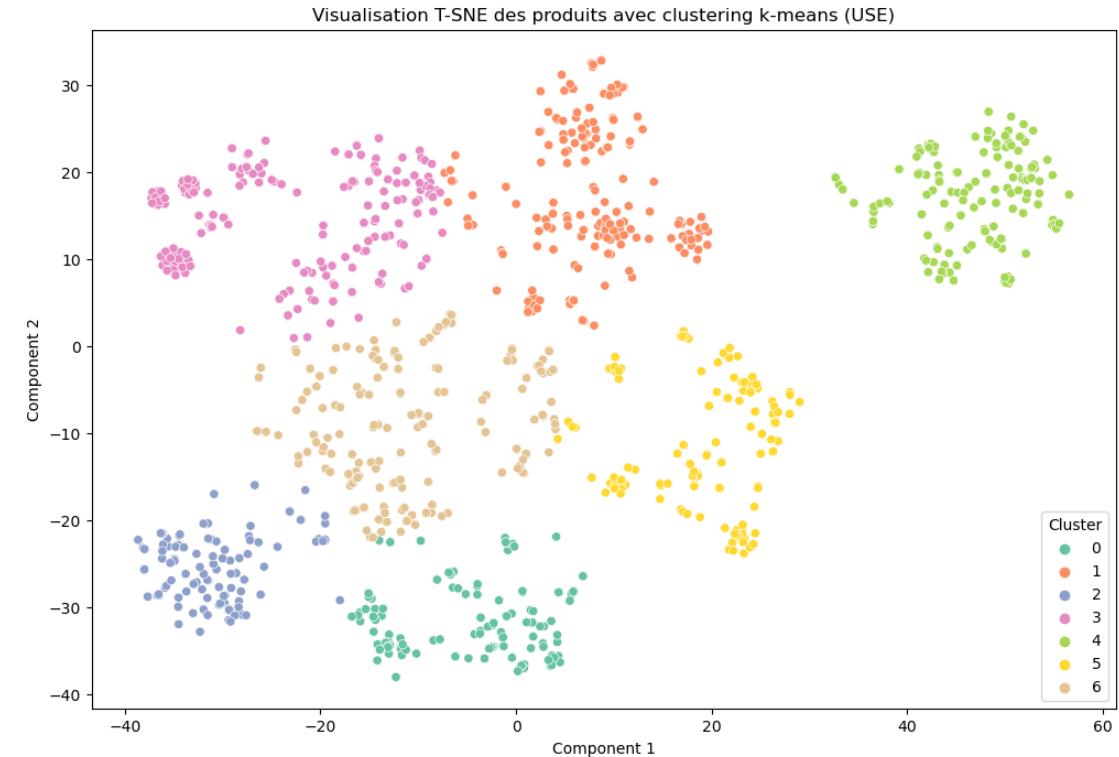
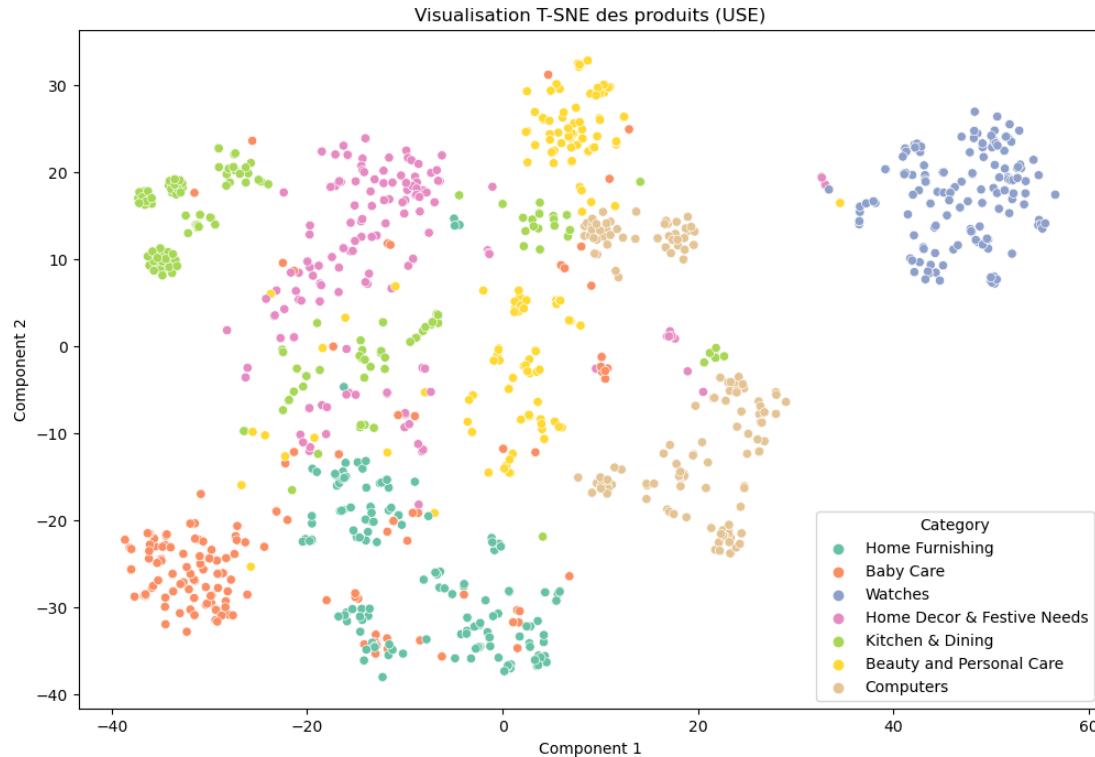
BERT



Adjusted Rand Score: 0.36237013040045846

Ordre de grandeur de 0.3, la classification est faisable.

USE



Adjusted Rand Score (USE): 0.43914335331065046

Ordre de grandeur de 0.4, la classification est faisable.

4. MÉTHODOLOGIE ADOPTÉE POUR LA CLASSIFICATION NON SUPERVISÉE DE DONNÉES IMAGES

ANALYSE D'UNE IMAGE ET DIFFÉRENTES APPROCHES DE TRANSFORMATION

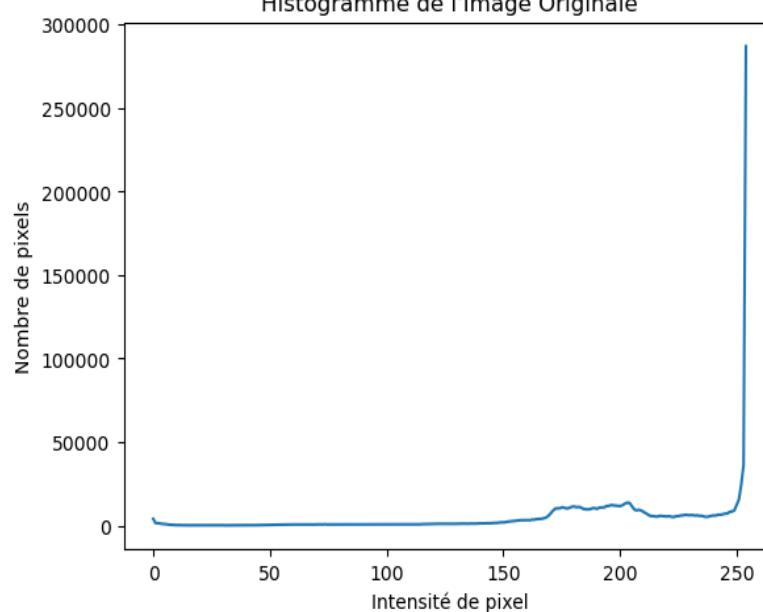
Image Originale



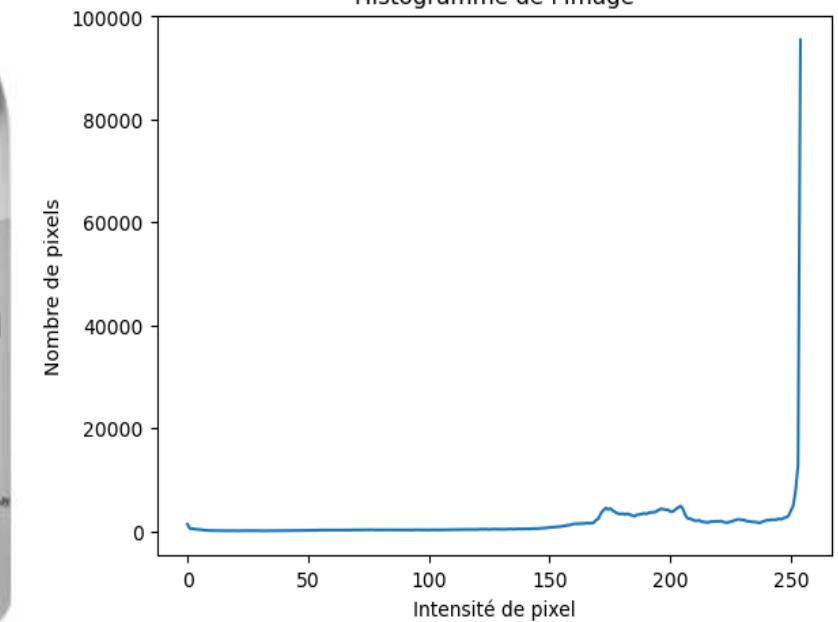
Image en Niveaux de Gris



Histogramme de l'Image Originale



Histogramme de l'Image



FAISABILITÉ DE CLASSIFICATION AUTOMATIQUE D'IMAGES

ANALYSE D'UNE IMAGE ET DIFFÉRENTES APPROCHES DE TRANSFORMATION

Image avec Égalisation d'Histogramme

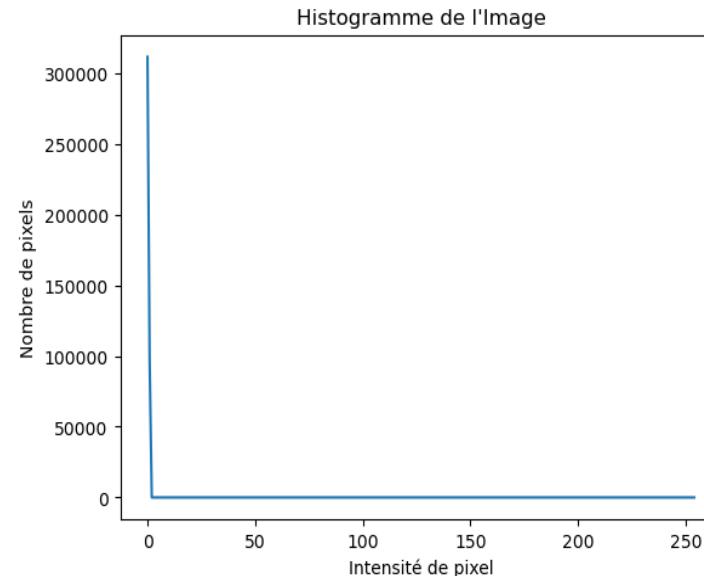
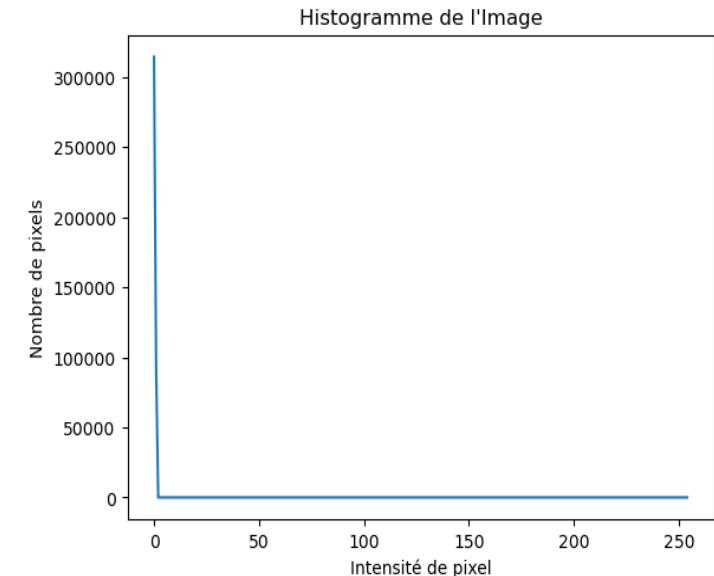


Image avec Filtrage du Bruit



FAISABILITÉ DE CLASSIFICATION AUTOMATIQUE D'IMAGES

ANALYSE D'UNE IMAGE ET DIFFÉRENTES APPROCHES DE TRANSFORMATION

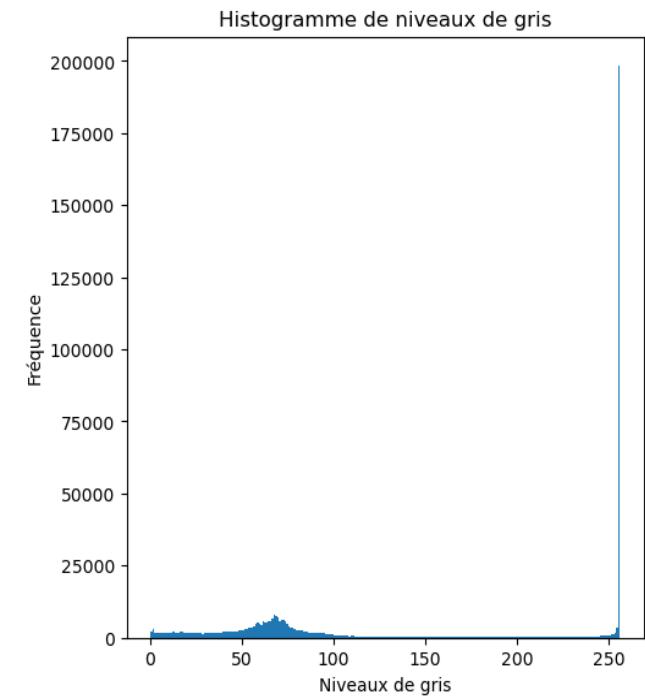
Image Originale



Image avec Contraste Amélioré

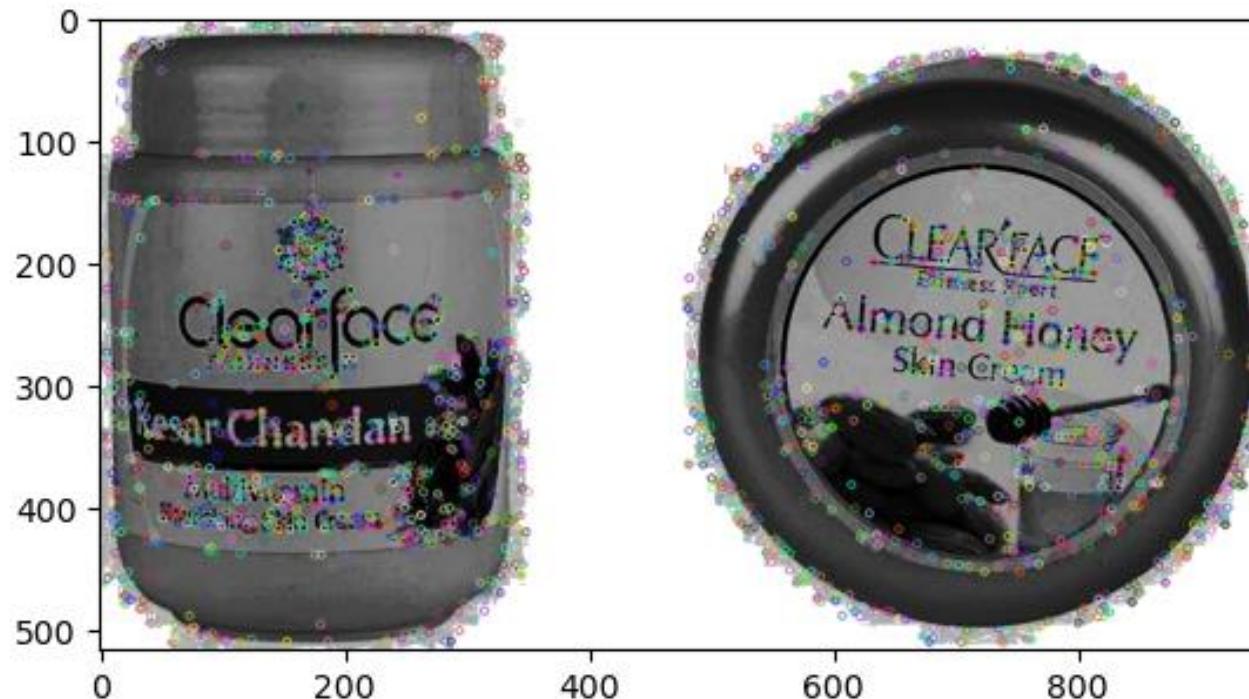


Image originale



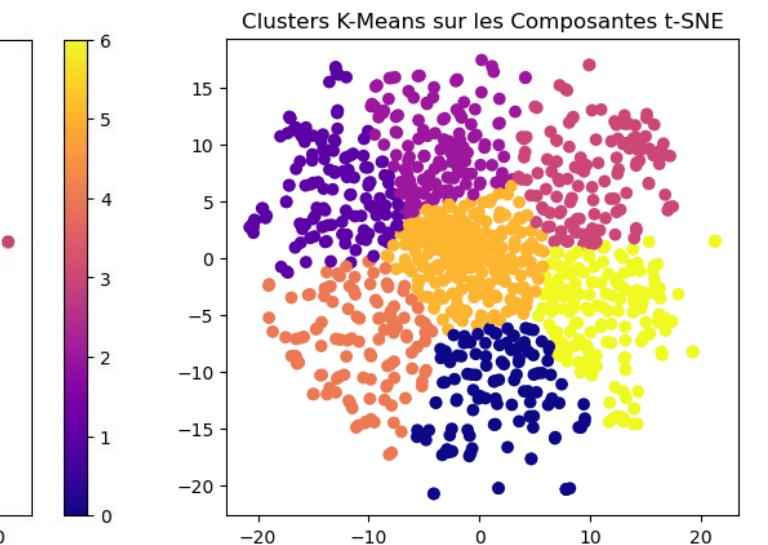
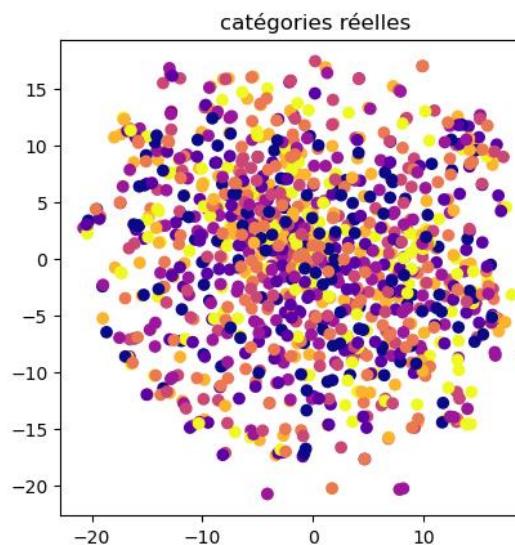
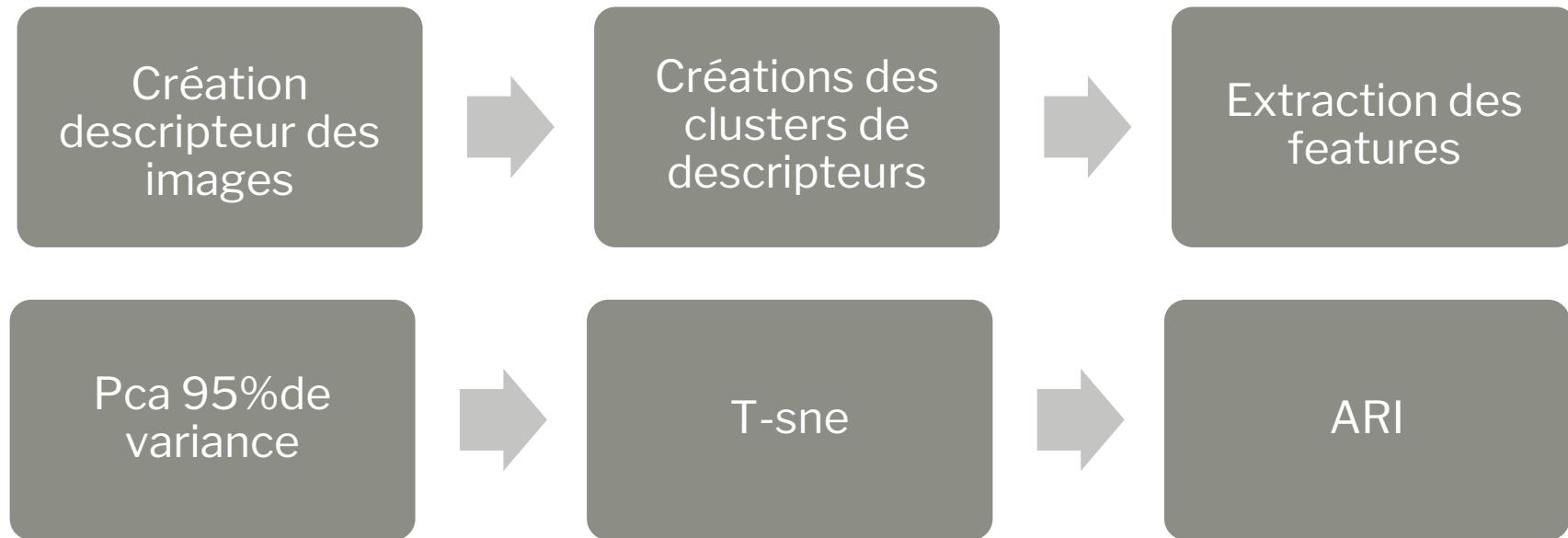
FAISABILITÉ DE CLASSIFICATION AUTOMATIQUE D'IMAGES

Le Sift détecte des points d'intérêt, construis de l'échelle de l'image, détecte de l'orientation, décris des caractéristiques.



Descripteurs : (2387, 128)

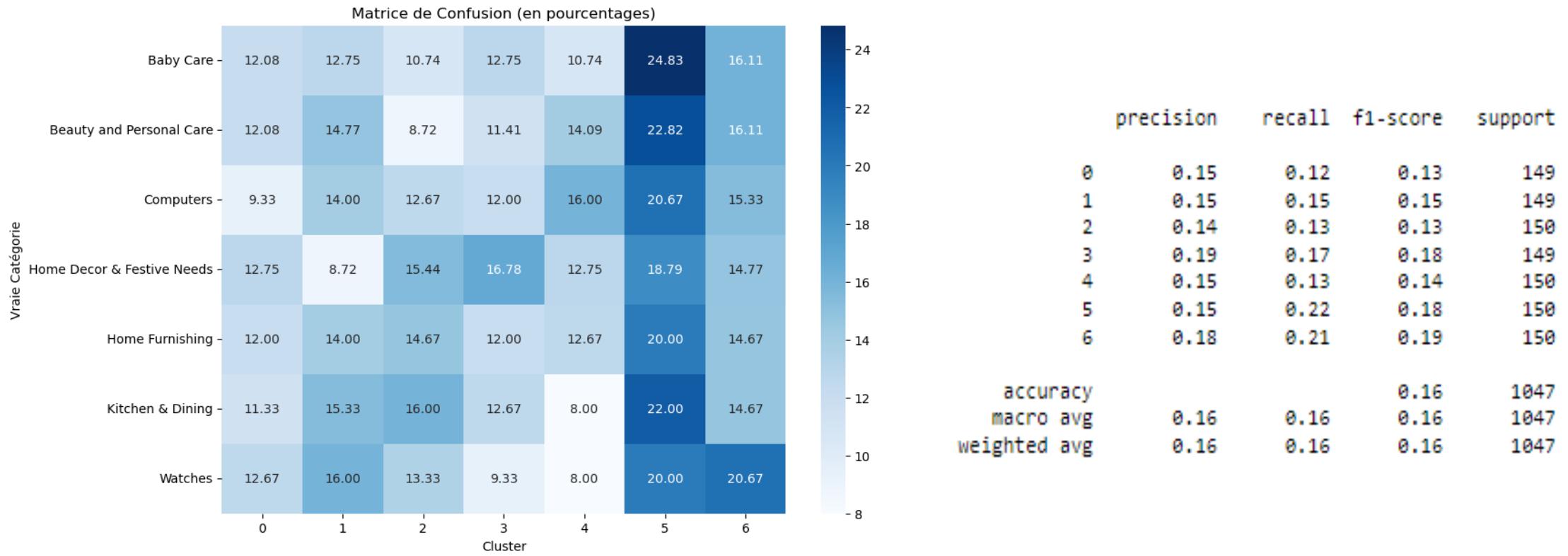
```
[[124.  47.  26. ...  0.  0.  0.]  
 [ 1.   2.  0. ...  0.  0.  3.]  
 [ 0.   0.  0. ... 20.  1.  7.]  
 ...  
 [ 0.   0.  0. ...  2.  1.  2.]  
 [ 0.   0.  2. ... 51.  2.  0.]  
 [ 0.   0.  0. ... 70.  3.  1.]]
```



Adjusted Rand Score (ARI): -0.00197064706630822

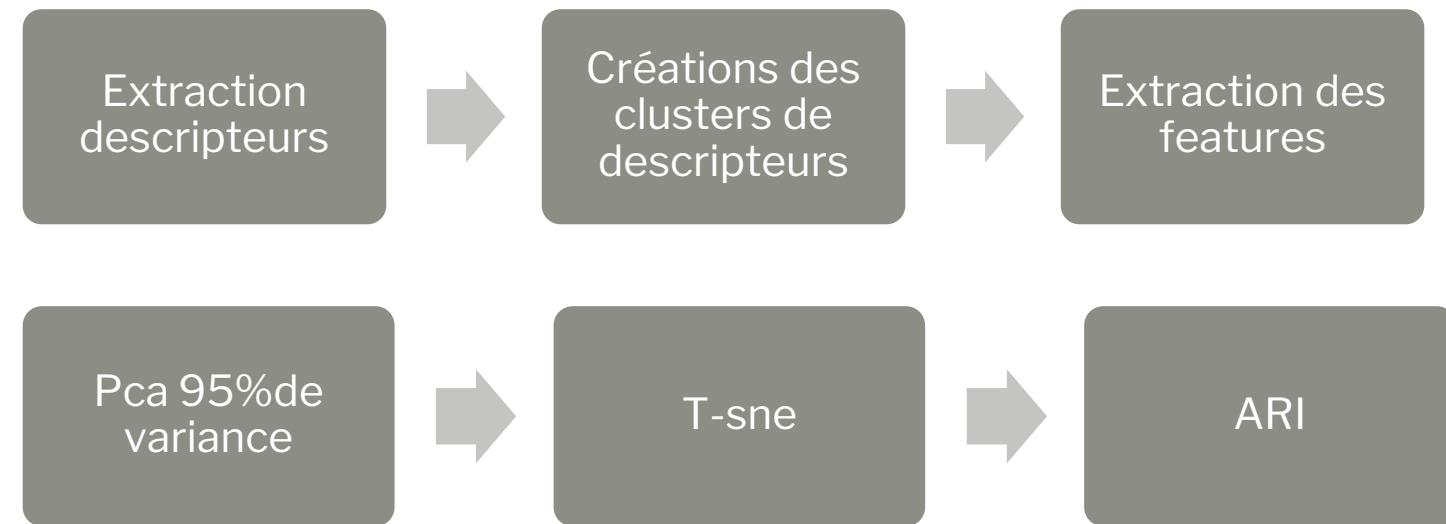
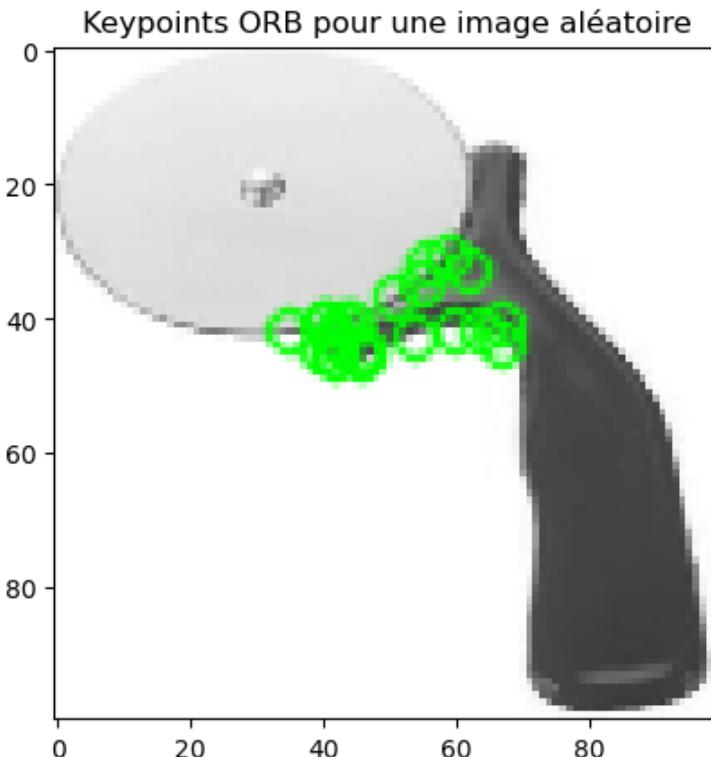
La valeur, de l'ordre de 0.00 à 0.1, confirme le visuel.

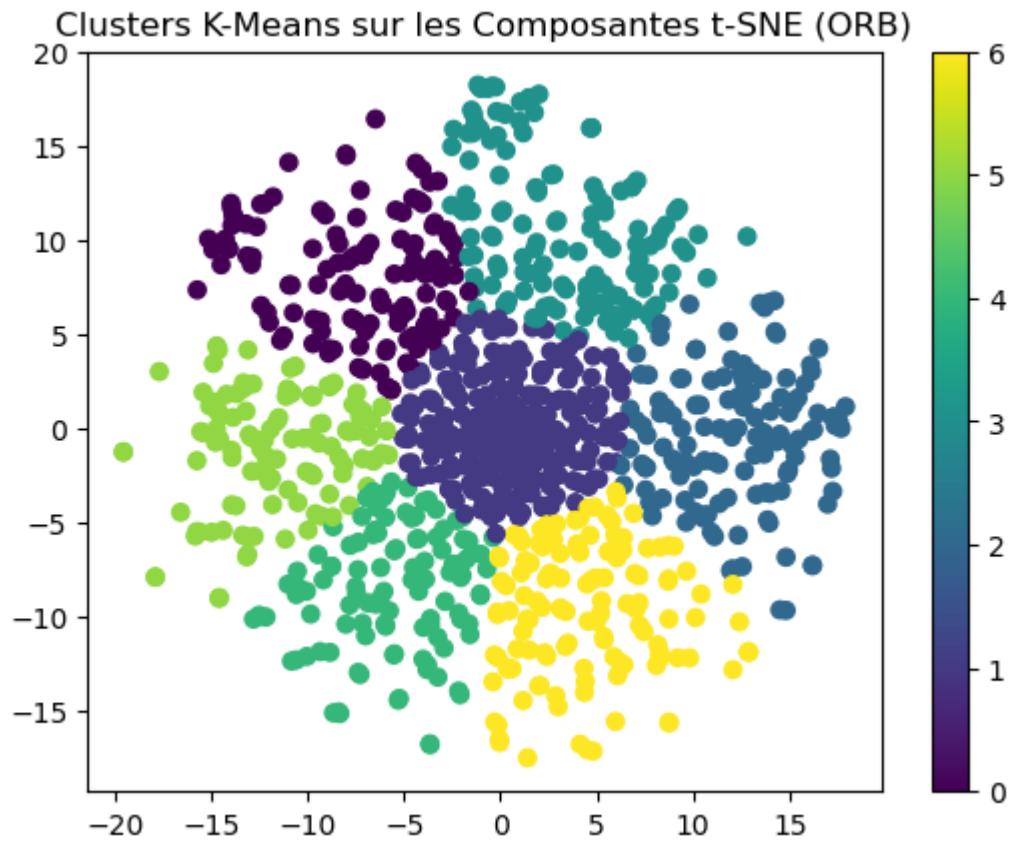
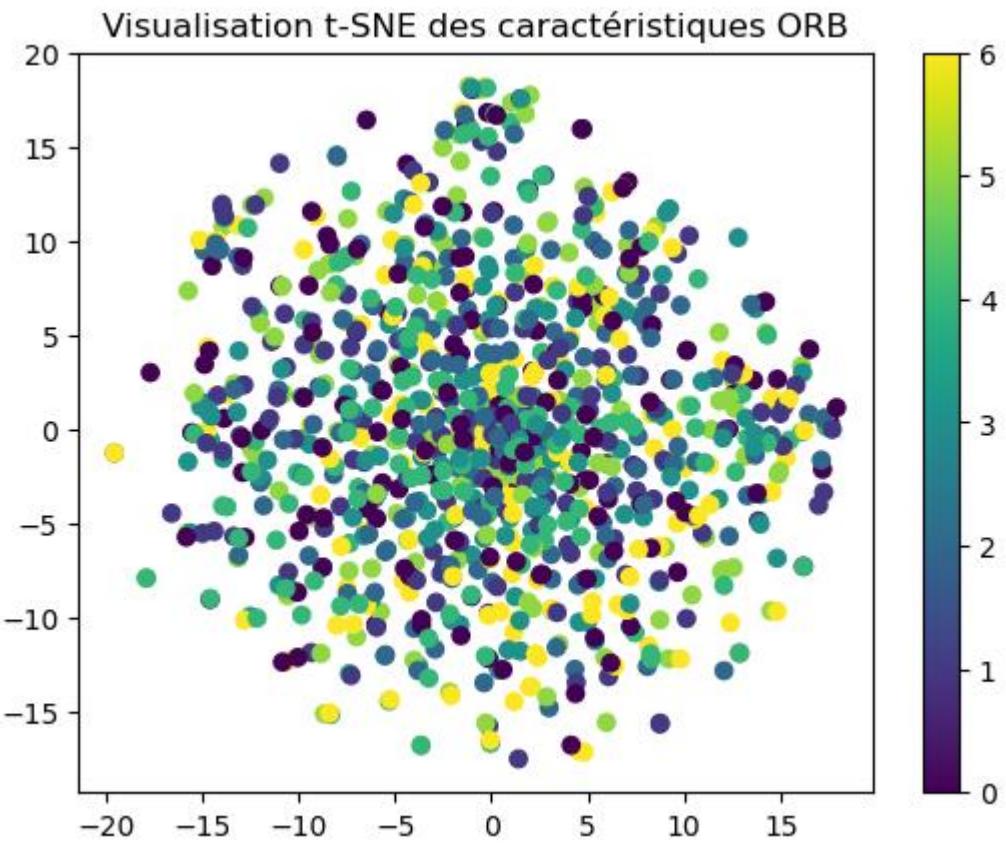
ANALYSE PAR CLASSE



FAISABILITÉ DE CLASSIFICATION AUTOMATIQUE D'IMAGES

ORB



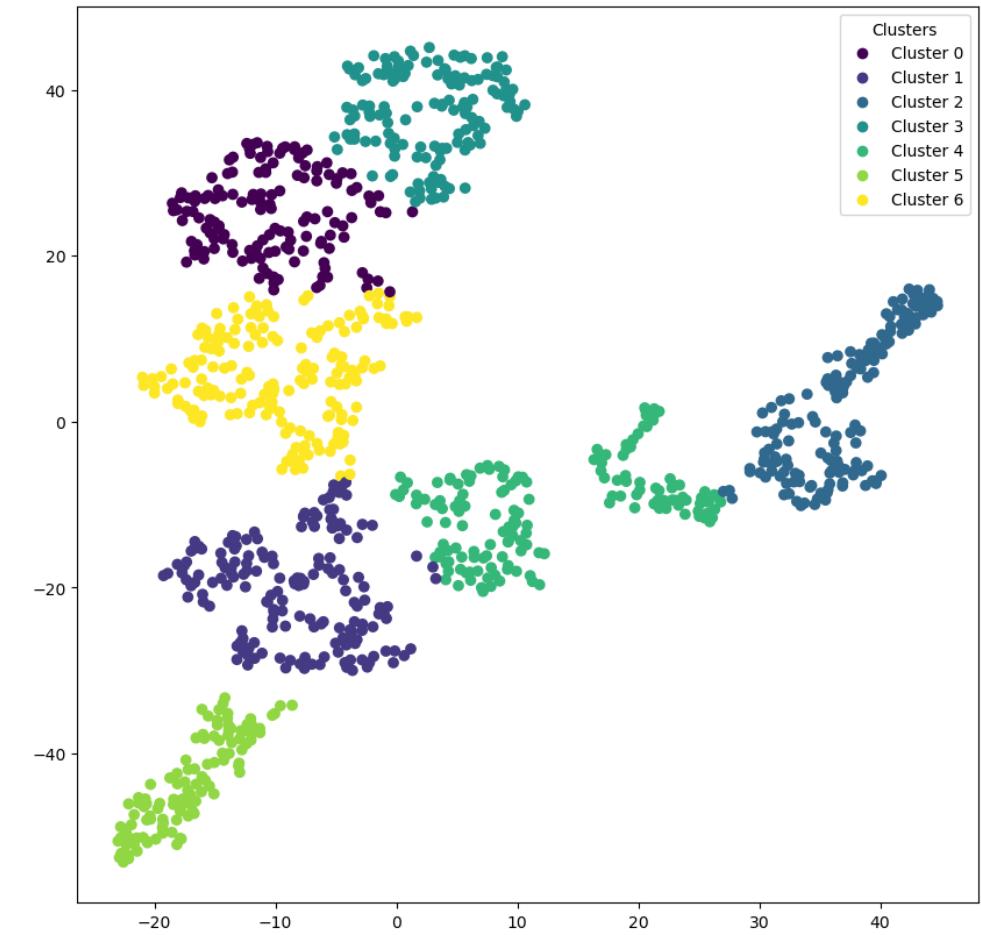
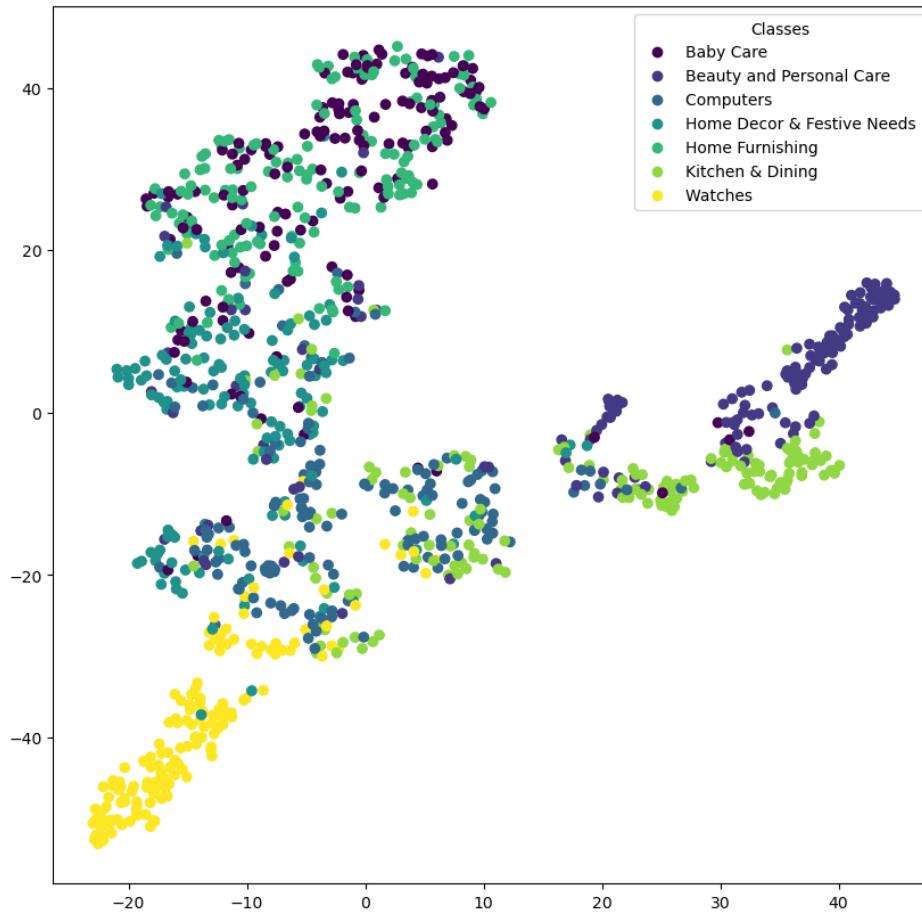


Adjusted Rand Score (ARI) avec ORB: 0.00039670233175964935

Le résultat est similaire au sift

Adjusted Rand Score (ARI): -0.0012080496299577332

FAISABILITÉ DE CLASSIFICATION AUTOMATIQUE D'IMAGES VIA CNN TRANSFERLEARNING



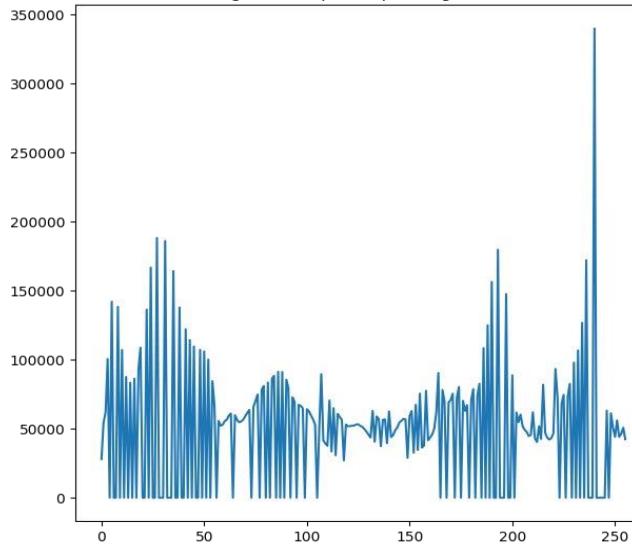
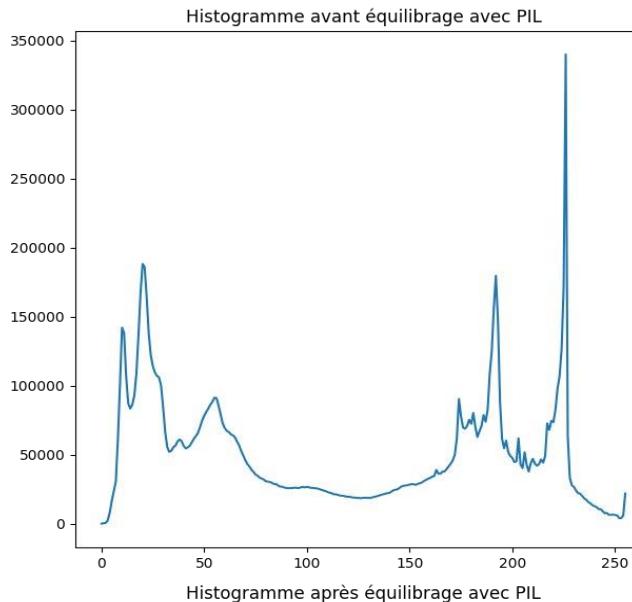
Score Ari :

ARI entre t-SNE et K-Means: 0.30636696012242487

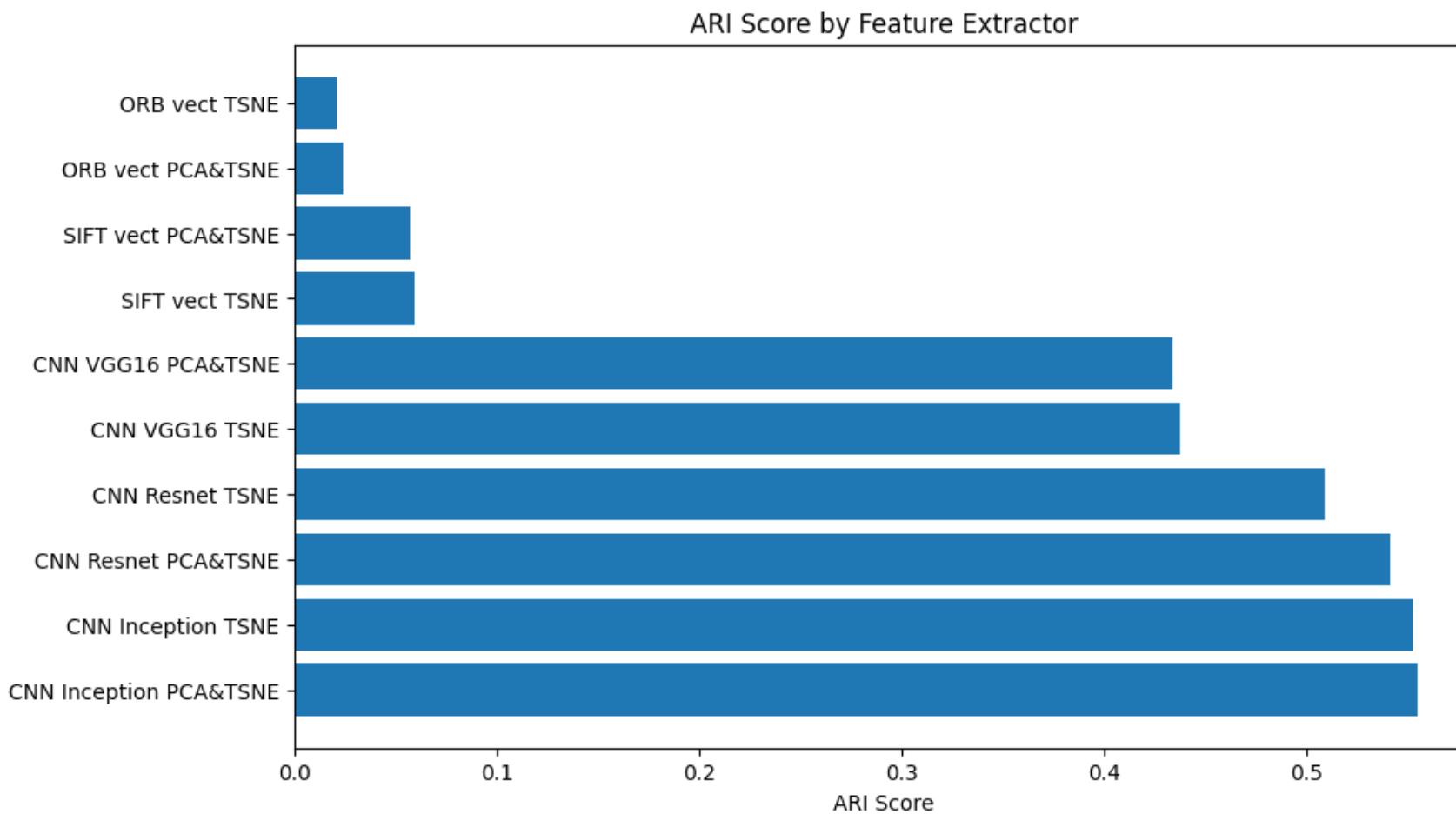
5. CLASSIFICATION SUPERVISÉ

- **Classification Supervisé**
- **Classification Supervisé avec Data Augmentation avant le modèle**

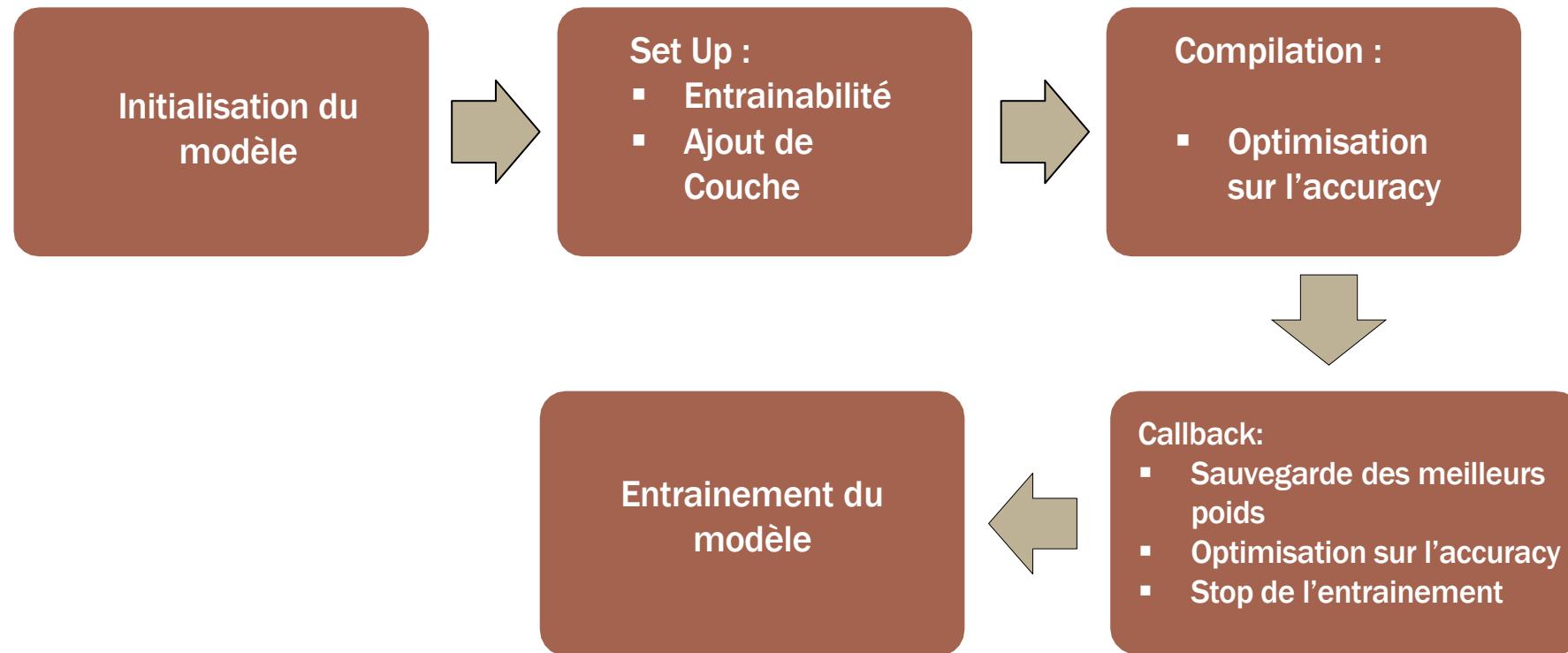
CLASSIFICATION IMAGE : PRE-PROCESSING



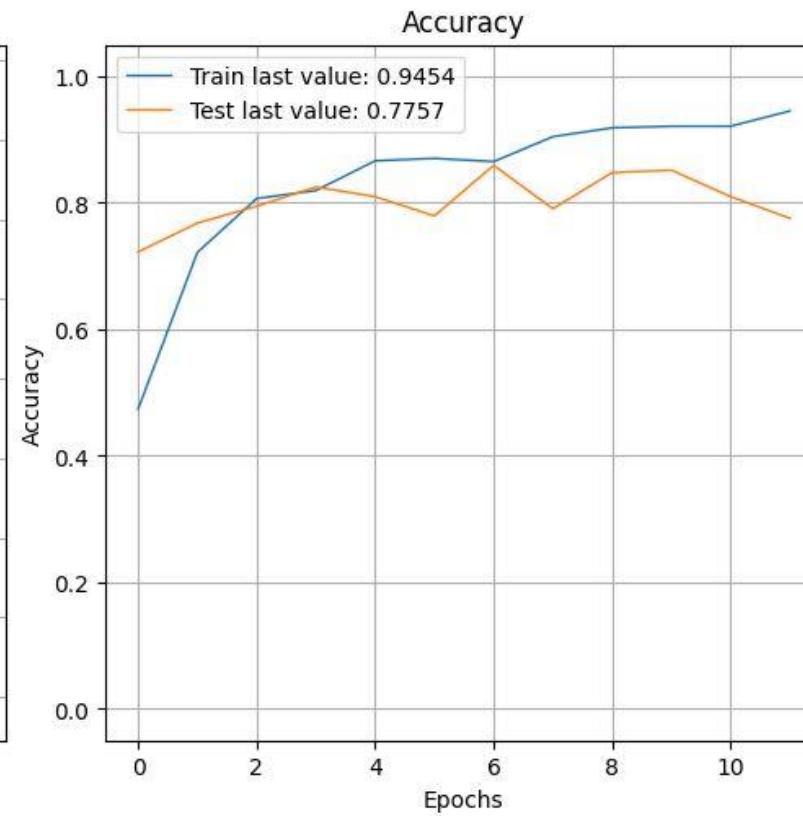
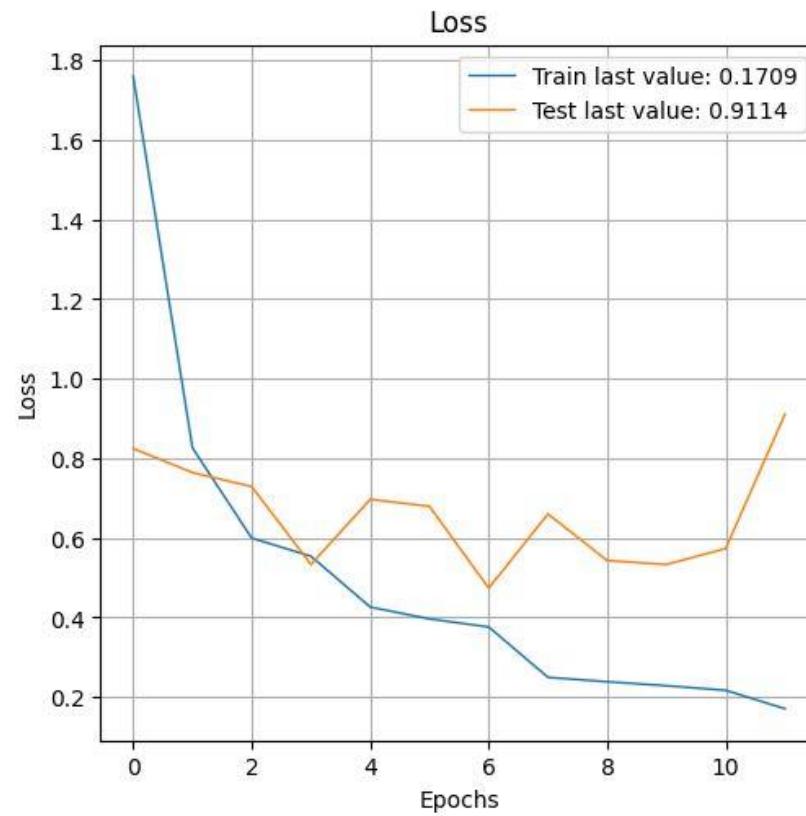
CLASSIFICATION IMAGE: RÉDUCTION DE DIMENSION ET EFFICACITÉ



CLASSIFICATION SUPERVISÉE D'IMAGES

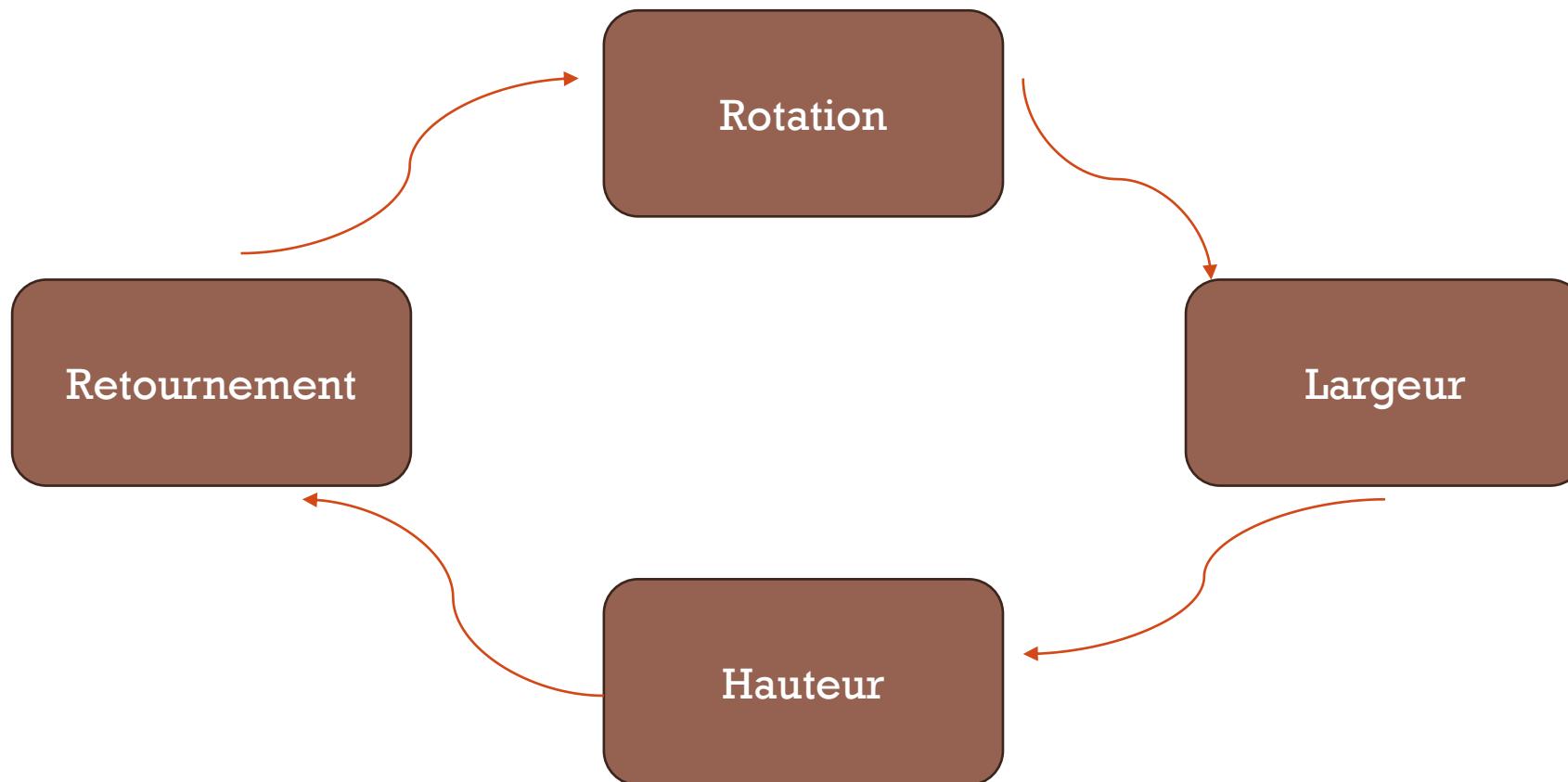


CLASSIFICATION IMAGE SUPERVISÉE



CLASSIFICATION IMAGE : SUPERVISÉE AVEC DATA AUGMENTATION AVANT LE MODÈLE

❖ DATA AUGMENTATION AVANT LE MODEL



Classification Image: Supervisee, Accuracy

Sans data augmentation

	0	1	2	3	4	5	6
watches -	33	0	0	0	2	0	0
kitchen & dining -	3	27	4	4	0	0	0
home furnishing -	1	0	28	0	0	0	0
beauty and personal care -	6	1	0	32	2	3	0
computers -	10	0	0	0	22	0	0
home decor & festive needs -	1	2	4	1	0	34	0
baby care -	0	0	1	0	0	0	41

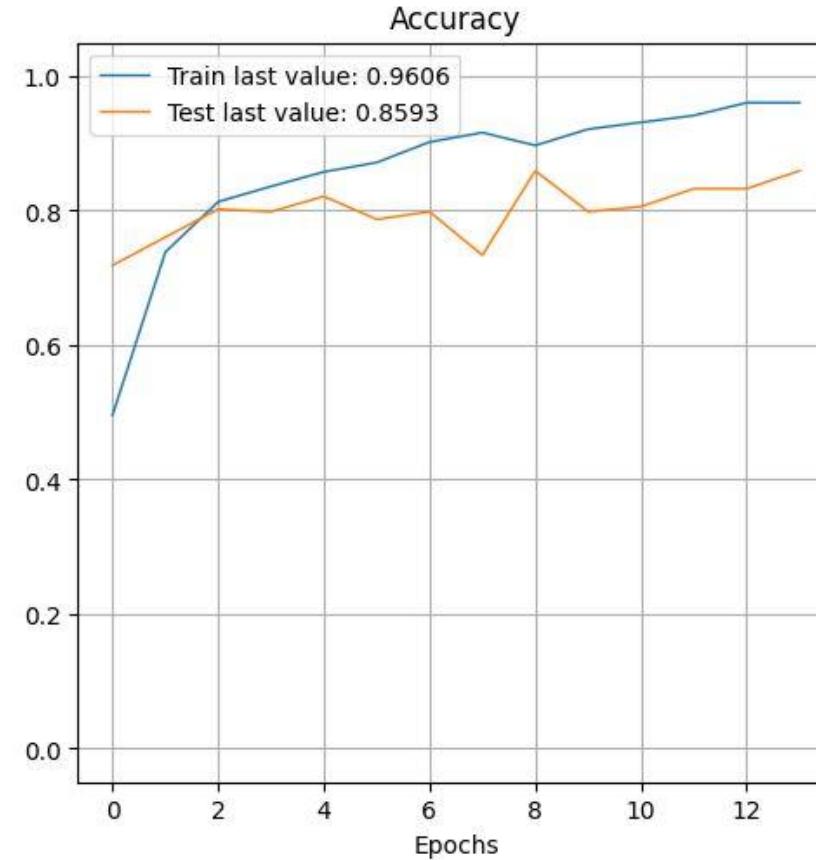
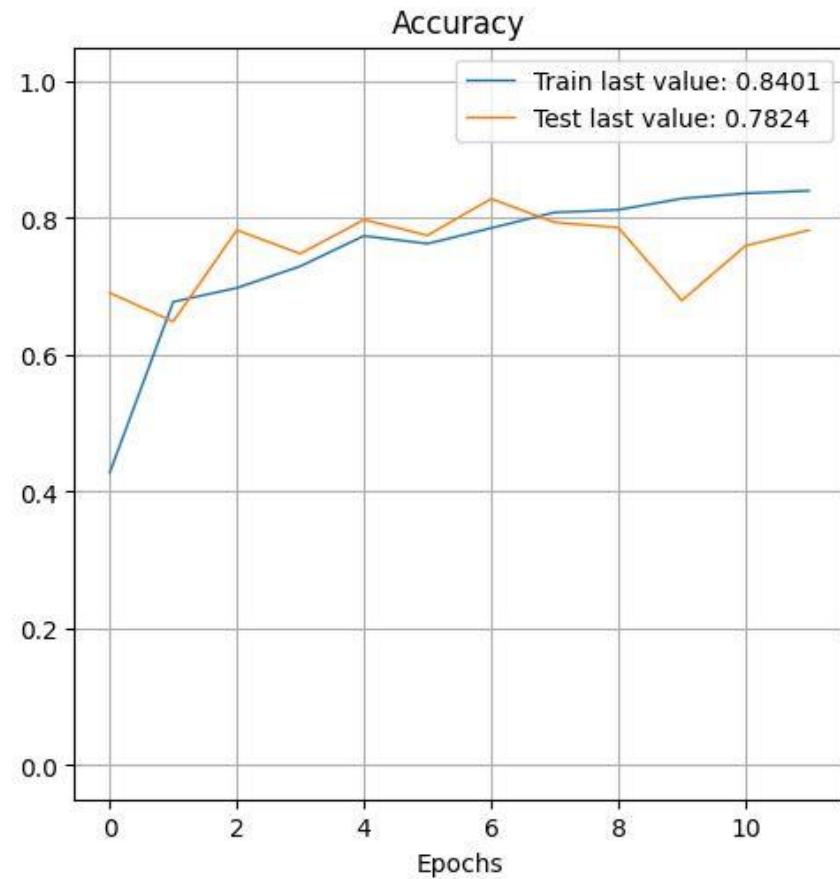
Accuracy

	Watches	1,00	0,95
Kitchen	0,92	0,94	
Home furnish	0,85	0,92	
Beauty	0,90	0,86	
Computer	0,76	0,82	
Home decor	0,86	0,65	
Baby care	0,61	0,87	
Avg	0,84	0,86	

Avec data augmentation

	0	1	2	3	4	5	6
watches -	26	0	0	8	1	0	0
kitchen & dining -	0	30	1	5	0	0	2
home furnishing -	0	0	28	0	0	1	0
beauty and personal care -	0	1	2	39	1	1	0
computers -	4	0	0	5	23	0	0
home decor & festive needs -	0	4	2	3	0	33	0
baby care -	0	0	1	0	0	0	41

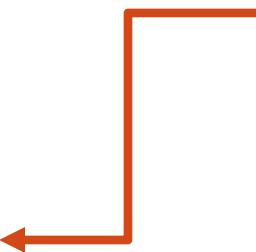
Classification Image : Supervisée, sur-apprentissage



6 -PRÉSENTATION DU TEST DE L'API



RESPECT DES NORMES RGPD



```
querystring = {"ingr":"champagne"}  
['foodId', 'label', 'category', 'foodContentsLabel', 'image']
```

	foodId	label	category	foodContentsLabel	image
0	food_a650mk2a5dmqb2adiamu6bei/duu	Champagne	Generic foods	NaN	https://www.edamam.com/food-img/a71/a718cf3c52...
1	food_b753ithmdb8psbt0w2k9aquo0c	Champagne Vinaigrette, Champagne	Packaged foods	OLIVE OIL; BALSAMIC VINEGAR; CHAMPAGNE VINEGAR...	NaN
2	food_b3dyababj054xobm0r8jzbghjgqe	Champagne Vinaigrette, Champagne	Packaged foods	INGREDIENTS: WATER; CANOLA OIL; CHAMPAGNE VINE...	https://www.edamam.com/food-img/d88/d88b64d973...
3	food_a9e0ghsamvoc45bw2ybsa3gken9	Champagne Vinaigrette, Champagne	Packaged foods	CANOLA AND SOYBEAN OIL; WHITE WINE (CONTAINS S...	NaN
4	food_an4jjeaucpus2a3u1ni8auhe7q9	Champagne Vinaigrette, Champagne	Packaged foods	WATER; CANOLA AND SOYBEAN OIL; WHITE WINE (CON...	NaN
5	food_bmu5dmkazwuvpa5prh1daa8jxs0	Champagne Dressing, Champagne	Packaged foods	SOYBEAN OIL; WHITE WINE (PRESERVED WITH SULFI...	https://www.edamam.com/food-img/ab2/ab245fc2a...
6	food_alpl44taoyv11ra0lic1qa8xculi	Champagne Buttercream	Generic meals	sugar; butter; shortening; vanilla; champagne;...	NaN
7	food_byap67hab6evc3a0f9w1oag3s0qf	Champagne Sorbet	Generic meals	Sugar; Lemon juice; brandy; Champagne; Peach	NaN
8	food_am5egz6aq3fpjlaf8xpkdbc2asis	Champagne Truffles	Generic meals	butter; cocoa; sweetened condensed milk; vanil...	NaN
9	food_bcz8rhiajk1fuva0vkfmekabouc0	Champagne Vinaigrette	Generic meals	champagne vinegar; olive oil; Dijon mustard; s...	NaN

CONCLUSION

- Possibilité de Classification Automatique
- S'orienter vers les images

MERCI POUR VOTRE ÉCOUTE