

Pour info, sur les slides 5, 9, 10, 11, 13 et 14 il y a des animations qui cachent le texte en dessous donc si possible plutôt regarder le ppt en mode diaporama pour voir le texte sur ces slides.



CLASSIFIEZ AUTOMATIQUEMENT DES BIENS DE CONSOMMATION



SOMMAIRE

1. Rappel de la problématique
2. Présentation du jeu de données
3. Approches de modélisation du texte
4. Approches de modélisation des images
5. Association du texte et de l'image
6. Conclusion



RAPPEL DE LA PROBLÉMATIQUE

- ▶ L'entreprise "Place de marché" est une marketplace qui propose des articles à des acheteurs en postant une photo et une description.
- ▶ Besoin d'automatiser l'attribution d'une catégorie à un article.
- ▶ Etudier la faisabilité d'un moteur de classification avec un niveau de précision suffisant.



PRÉSENTATION

- ▶ Jeu de données composé de 50 articles.
- ▶ Les catégories sont hiérarchisées
 - Exemple : ["Waterproof Baby Wipes"]
- ▶ Manipulations pour créer un article.
- ▶ Certaines sous-catégories ont plus de données donc décisions de faisabilité d'un modèle.

cat	sub_cat	uniq_id
Baby Care	Baby & Kids Gifts	15
	Baby Bath & Skin	14
	Baby Bedding	15
	Baby Grooming	4
	Diapering & Potty Training	7
	Feeding & Nursing	8
	Furniture & Furnishings	2
	Infant Wear	84
Beauty and Personal Care	Strollers & Activity Gear	1
	Bath and Spa	7
	Beauty Accessories	1
	Body and Skin Care	15
	Combos and Kits	24
	Eye Care	2
	Fragrances	65
	Hair Care	9
	Health Care	7
	Makeup	18
	Men's Grooming	1
	Women's Hygiene	1

50 articles.

ne :

gories par

de
tude de

place de marché

PRÉSENTATION DU JEU DE DONNÉES

► 7 catégories principales parfaitement équilibrées :

- ❑ Home Furnishing : 150 articles
- ❑ Baby Care : 150 articles
- ❑ Watches : 150 articles
- ❑ Home Decor & Festive Needs : 150 articles
- ❑ Kitchen & Dining : 150 articles
- ❑ Beauty and Personal Care : 150 articles
- ❑ Computers : 150 articles



place de marché

APPROCHES DE MODÉLISATION DU TEXTE

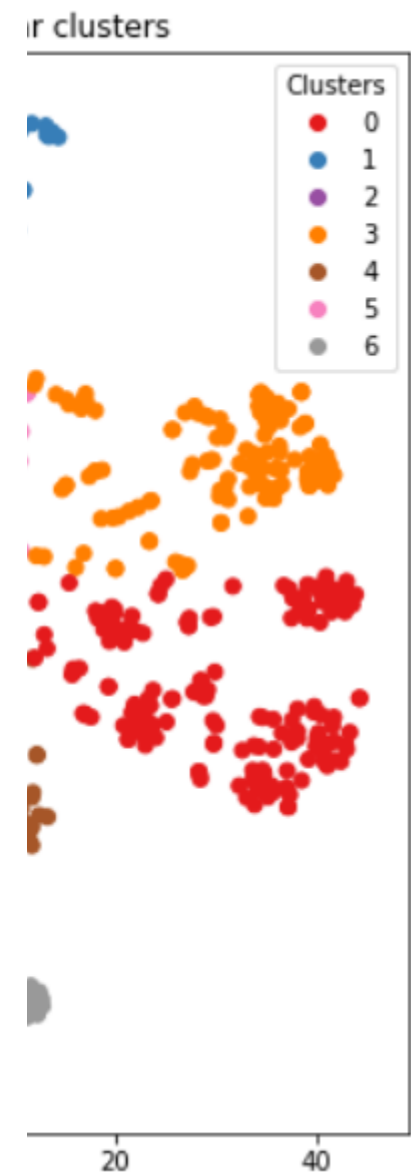
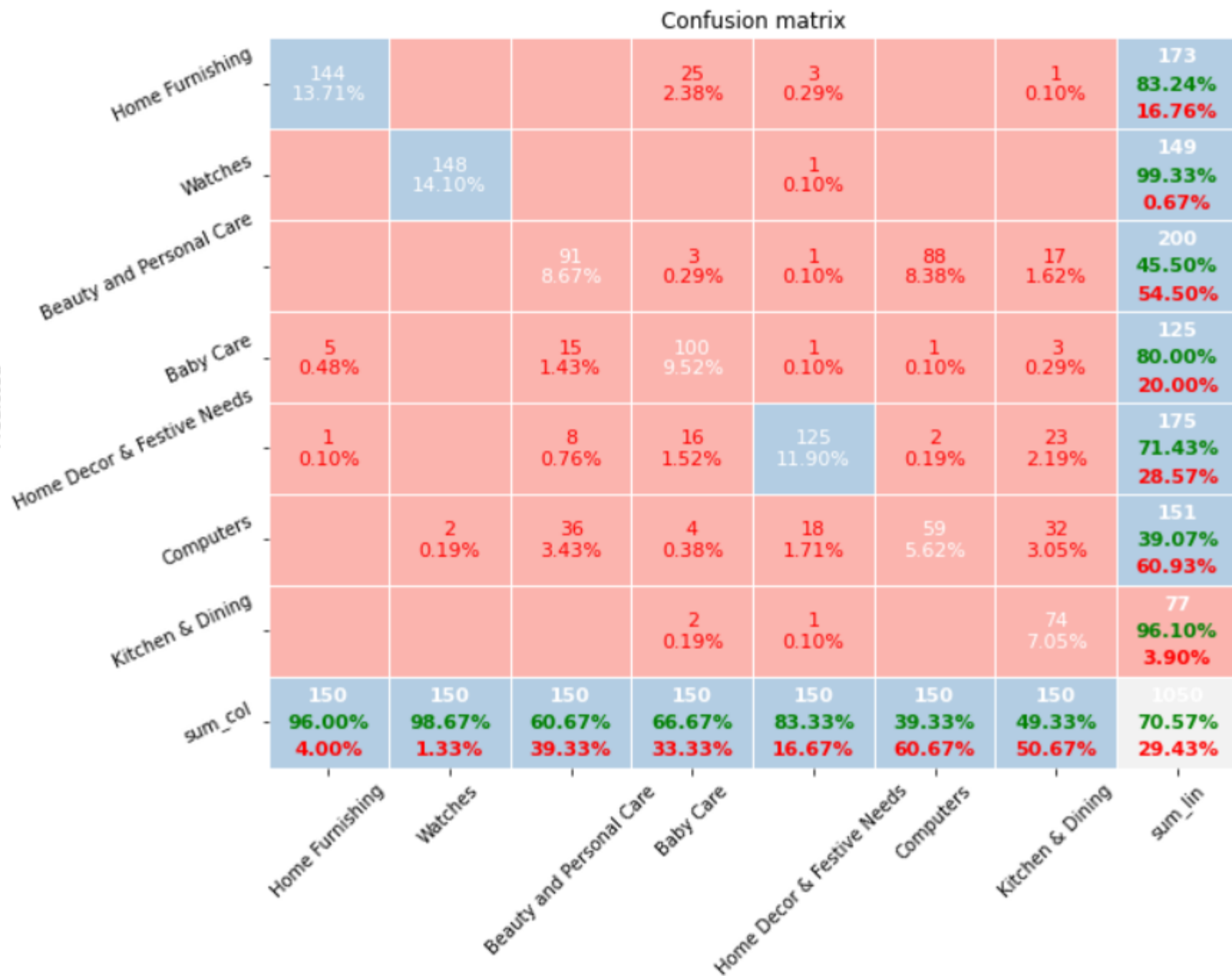
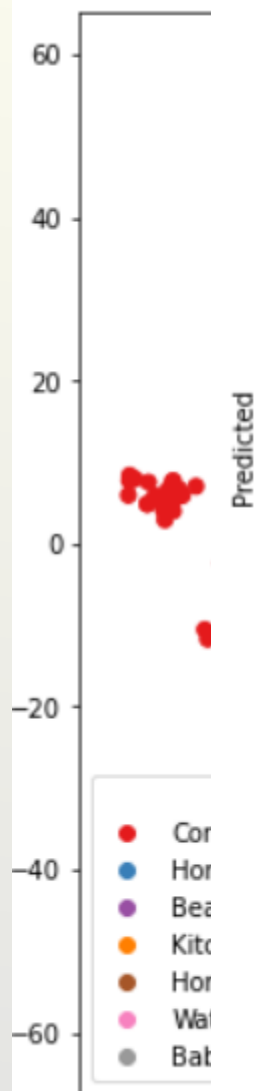
- ▶ Extraction des features texte via différentes approches :
 - ❑ De type 'bag of words' avec un comptage simple de mots et avec Tf-idf
 - ❑ De type word embedding avec Word2Vec et Doc2Vec
 - ❑ De type word embedding avec BERT (Bidirectional Encoder Representations from Transformers)
 - ❑ De type sentence embedding avec USE (Universal Sentence Encoder)
- ▶ Réduction de dimension des features obtenus grâce à un t-sne (2 composantes).
- ▶ Clustering (k-means) sur les données réduites puis calcul de l'Adjusted Rand Index (ARI) par rapport aux vrais catégories.



APPROCHES DE MODÉLISATION DU TEXTE

- ▶ Prétraitements réalisés pour l'approche 'bag of words' et word embedding avec Word2Vec et Doc2Vec :
 - ❑ Mise du texte en minuscule
 - ❑ Suppression des ponctuations
 - ❑ Suppression des nombres
 - ❑ Tokenisation
 - ❑ Suppression des 'stopwords' anglais
 - ❑ Lemmatisation
 - ❑ Suppression des mots inférieurs ou égaux à 2 caractères





Descriptions préprocessées de certains articles 'Computers' clusterisés en 'Beauty and Personal Care' :

EXTE

```
'buy smart router genuine product day replacement guarantee free shipping cash delivery',  
'buy genuine product day replacement guarantee free shipping cash delivery',  
'buy wireless dual band router genuine product day replacement guarantee free shipping cash delivery',  
'buy apple genuine product day replacement guarantee free shipping cash delivery',  
'buy genuine product day replacement guarantee free shipping cash delivery',  
'buy wireless range extender genuine product day replacement guarantee free shipping cash delivery',  
'buy tew genuine product day replacement guarantee free shipping cash delivery',
```

Descriptions préprocessées de certains articles 'Beauty and Personal Care' bien clusterisés en 'Beauty and Personal Care' :

```
and link box extender pl  
array(['buy industry bangle four roll ring earring box vanity pouch industry bangle four roll ring earring box vanity pouch best  
price free shipping cash delivery genuine product day replacement guarantee',  
      'buy vanity pouch vanity pouch best price free shipping cash delivery genuine product day replacement guarantee',  
      'buy sally cross festival kit price genuine product day replacement guarantee free shipping cash delivery',  
      'buy equinox body fat analyzer genuine product day replacement guarantee free shipping cash delivery',  
      'buy wild stone red juice set genuine product day replacement guarantee free shipping cash delivery',  
      'buy ice drive dynamic pulse set genuine product day replacement guarantee free shipping cash delivery',  
      'vincent valentine set dark fire dark fire majesty deodorant set set price dark fire sparkle enchanting top note blend b  
lack currant star anise distinctive spiciness galanga root depth green pepper extract warmth base rich woody balsam cedar vanil  
la bean passionate royal majesty latest beautiful perfume anna young extravagant woman neo classical touch dark fire sparkle en  
chanting top note blend black currant star anise distinctive spiciness galanga root depth green pepper extract warmth base rich  
woody balsam cedar vanilla bean passionate royal majesty latest beautiful perfume anna young extravagant woman neo classical to  
uch',  
      'buy black code gift set set genuine product day replacement guarantee free shipping cash delivery',  
      'buy wild stone smoke juice set genuine product day replacement guarantee free shipping cash delivery',  
      'buy set genuine product day replacement guarantee free shipping cash delivery'],
```

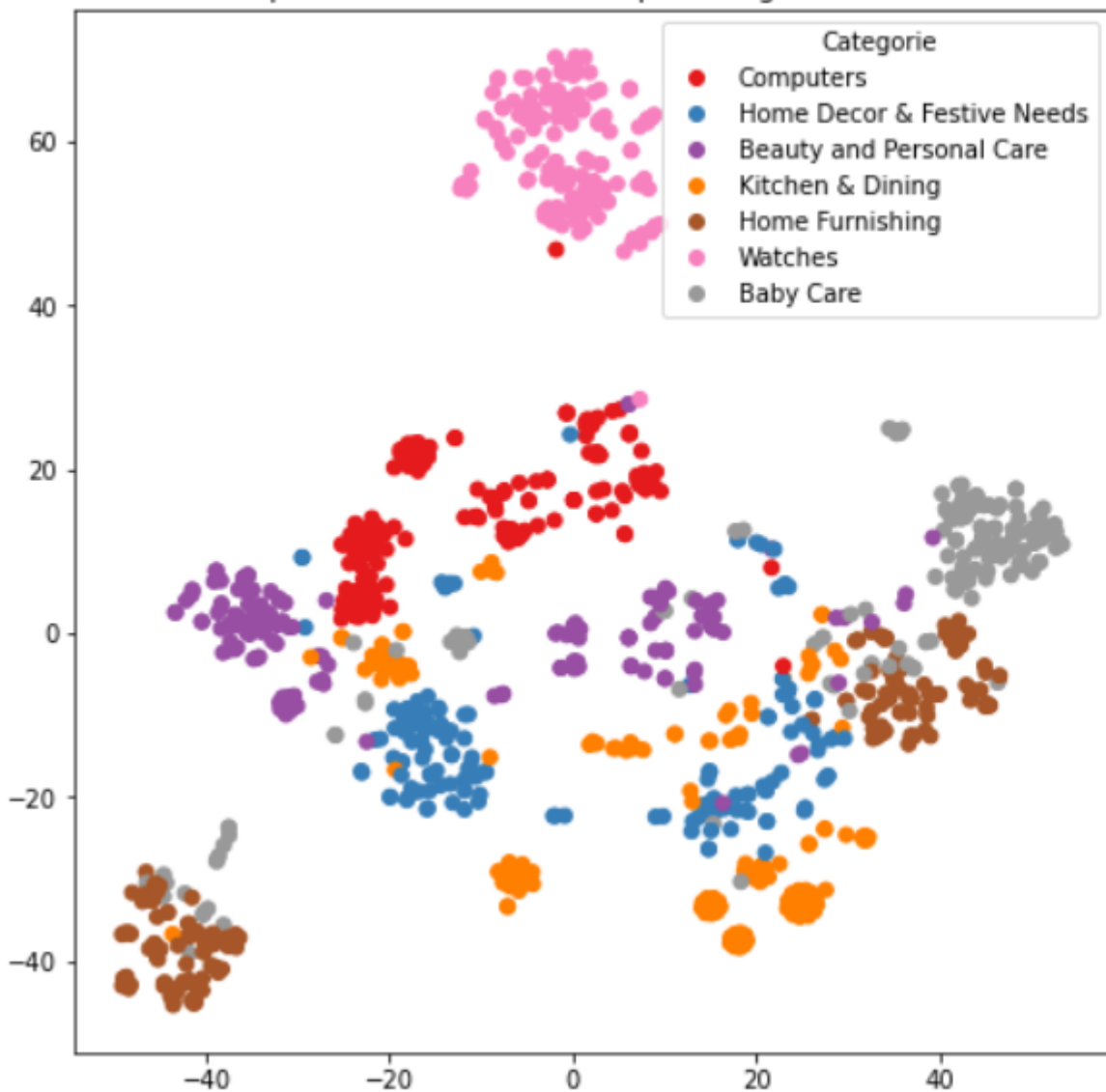
arché

'word2v

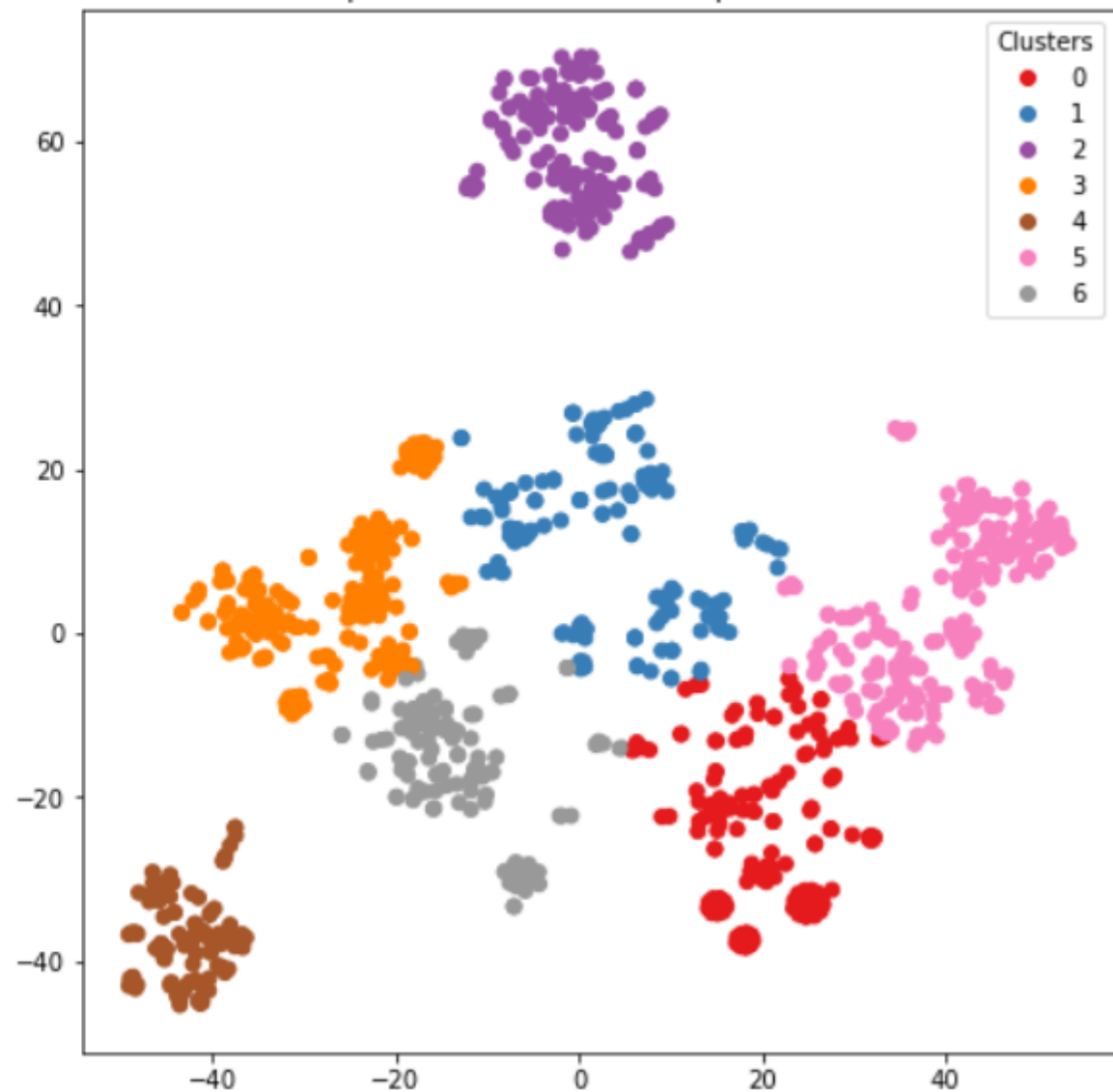
USE

entraîné

Représentation des articles par catégories réelles



Représentation des articles par clusters



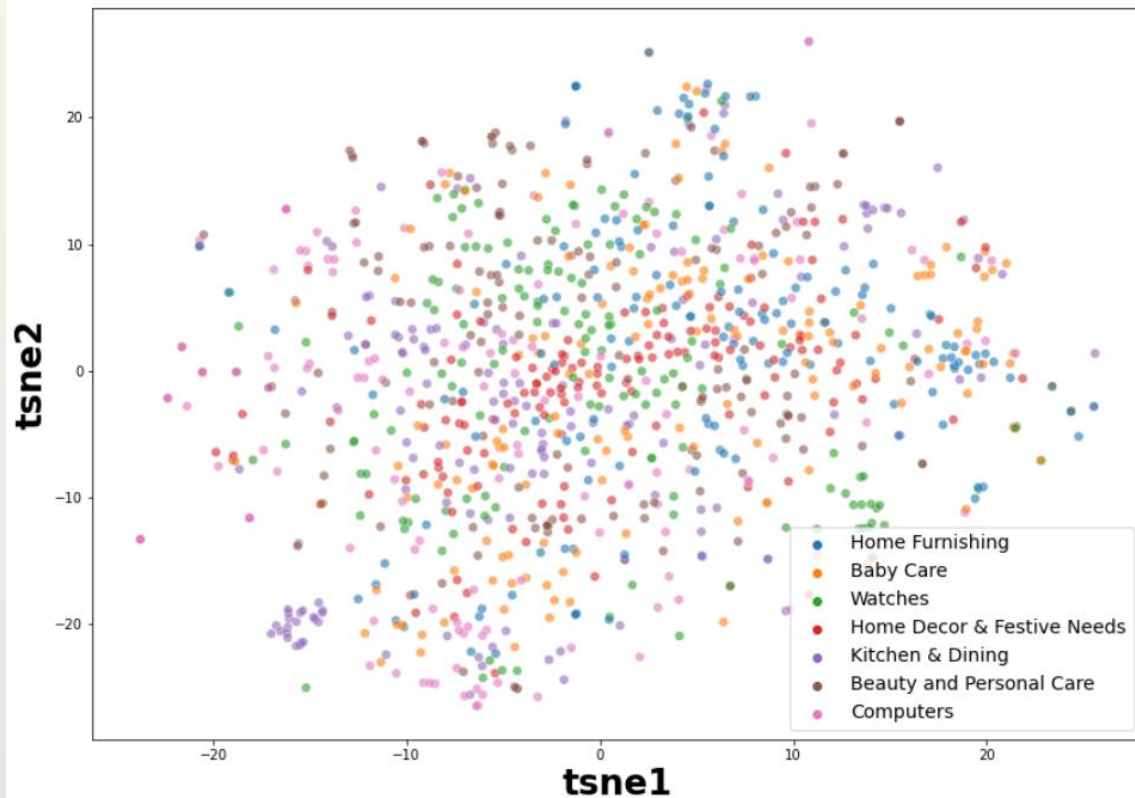
APPROCHES DE MODÉLISATION DES IMAGES

- ▶ Extraction des features image via différentes approches :
 - ❑ Un algorithme de type SIFT
 - ❑ Des algorithmes de type CNN Transfer Learning
 - ❖ VGG16
 - ❖ EfficientNet
 - ❖ ResNet50
- ▶ Réduction de dimension des features obtenus grâce à un t-sne (2 composantes).
- ▶ Clustering (k-means) sur les données réduites puis calcul de l'Adjusted Rand Index (ARI) par rapport aux vraies catégories.

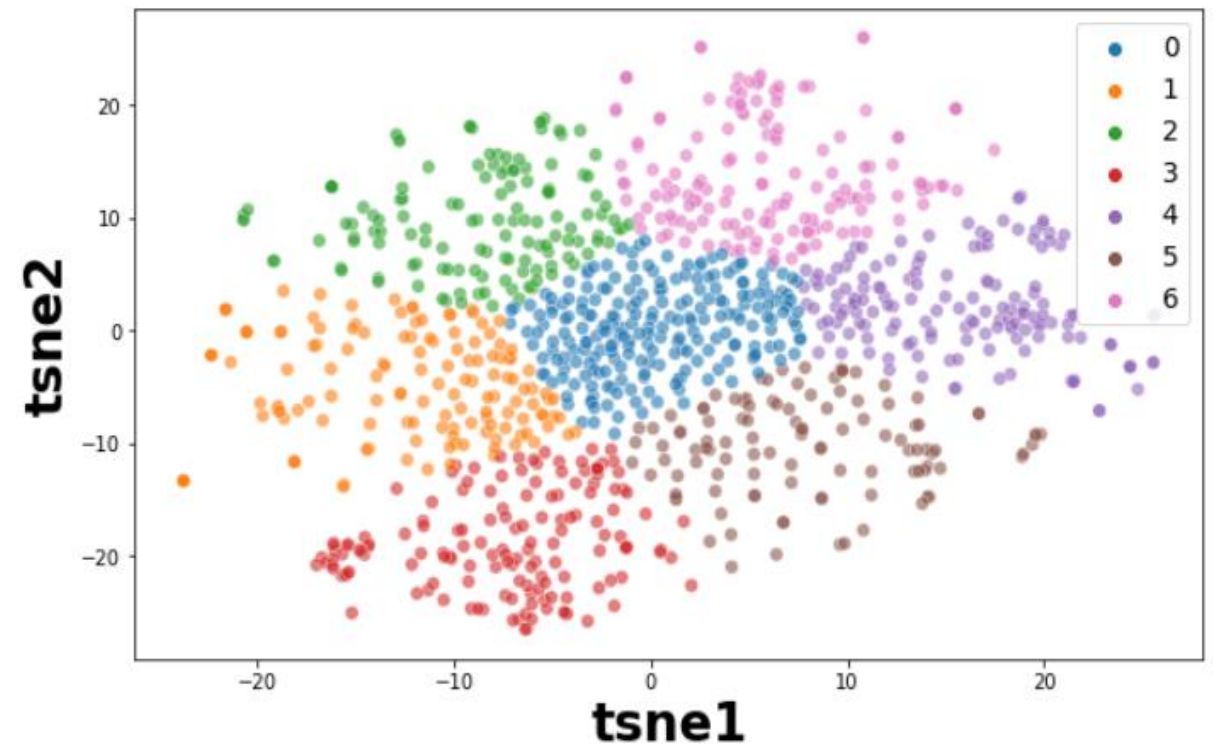


APPROCHES DE MODÉLISATION DES IMAGES

TSNE selon les vraies classes



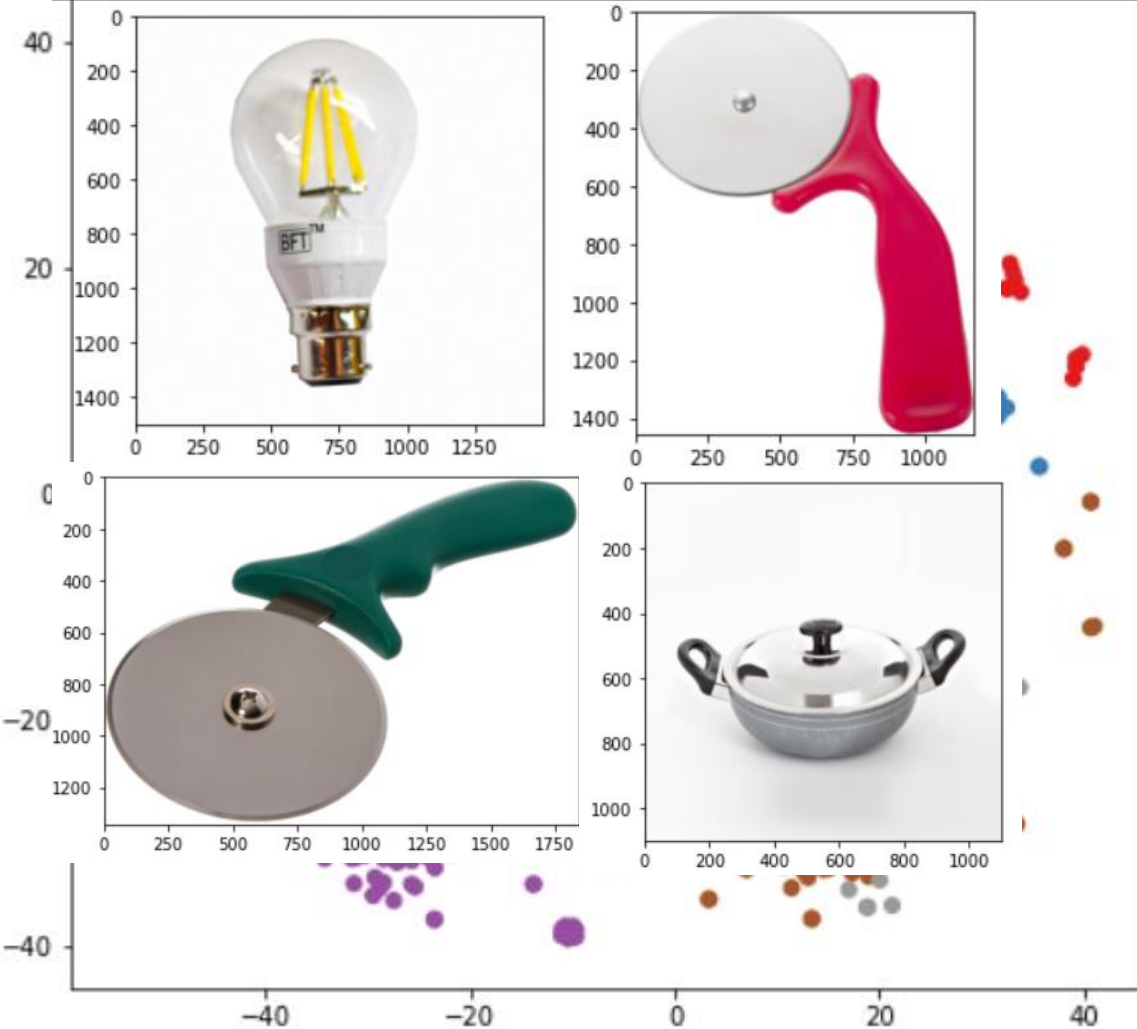
TSNE selon les clusters



ResNet50

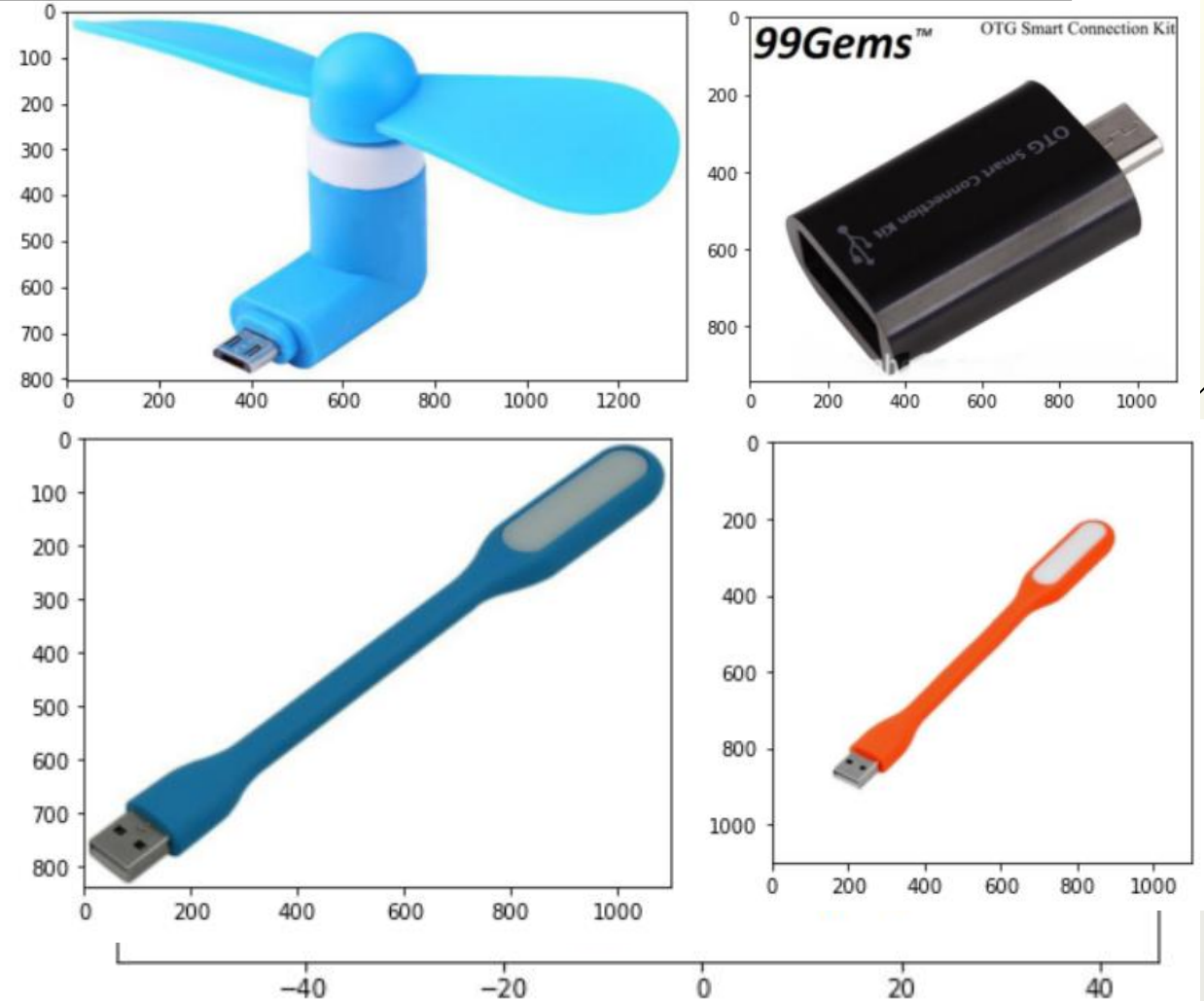
Représentation des articles par catégories réelles

Images d'articles qui ont été clusturisés en tant que 'Computers' alors que ce sont des 'Kitchen & Dining' :



Représentation des articles par clusters

Images d'articles qui ont été clusturisés en tant que 'Computers' et qui sont bien des 'Computers' :



ASSOCIATION DU TEXTE ET DE L'IMAGE

- ▶ Association des approches ayant obtenus le meilleur score ARI pour le texte avec celle ayant obtenus le meilleur score ARI pour les images.
- ▶ Deux méthodes testées :
 - ❑ Agrégation des 2 composantes du t-sne obtenus pour le texte et les images puis réalisation d'un clustering (k-means)
 - ❑ Agrégation directement des features extraites pour le texte et les images puis réalisation d'un t-sne et d'un clustering (k-means)
- ▶ Meilleurs scores ARI obtenus pour le texte avec l'approche Tf-idf et USE.
- ▶ Meilleurs scores ARI obtenus pour les images avec les algorithmes EfficientNet et ResNet50.



ASSOCIATION DU TEXTE ET DE L'IMAGE

► Rappel des scores ARI obtenus :

- Tf-idf = 0,53 ; USE = 0,43
- EfficientNet = 0,36 ; ResNet50 = 0,43

► Résultats des associations :

t-sne	EfficientNet	ResNet50
Tf-idf	0,53	0,55
USE	0,46	0,53

Features puis t-sne	EfficientNet	ResNet50
Tf-idf	0,36	0,42
USE	0,35	0,42



CONCLUSION

- ▶ Les caractéristiques extraites pour le texte et les images permettent de regrouper des produits de même catégorie.
- ▶ Certaines catégories sont très bien clusterisées ('Watches', 'Home Furnishing'), d'autres sont plus compliquées à différencier ('Computers').
- ▶ Cependant en appliquant un modèle de classification aux images par exemple, on obtient tout de même de bons résultats sur l'ensemble du jeu de test (86% d'accuracy pour EfficientNet par exemple, non présenté dans ces slides mais présenté dans le notebook).