

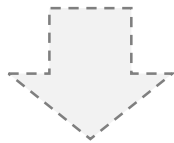
Formation Data Scientist

**Soutenance Projet 6 :
Classifiez automatiquement des biens de consommation**

-- Mélanie WARY --

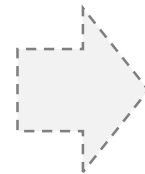
PROBLÉMATIQUE

L'entreprise "Place de marché", qui lance une marketplace e-commerce et souhaite **automatiser l'attribution de la catégorie de chaque article.**



MISSIONS

- (1) Réaliser une première **étude de faisabilité** d'un moteur de **classification** d'articles en différentes catégories, **basée sur l'image et la description** fournies pour chaque article, avec un **niveau de précision suffisant**.
- (2) Présenter les résultats de la classification sous la forme d'une **représentation en 2D**



CONTRAINTES & INTERPRÉTATION :

- **Pré-traitement et extraction de features** des données
 - **textuelles (description)**
 - **visuelles (photo)** – SIFT / ORB / SURF
- **Réduction(s) dimensionnelle(s)** – 2D
- **Clustering non supervisé** – nombre de clusters = nombre de « catégories vraies »
- **Optimisation** et mesure du niveau de **précision** par comparaison à la « catégorisation vraie ».

JEU DE DONNEES

→ 1050 produits

→ 3 variables d'intérêt:

CATEGORIES

Arbre de catégories

Min. 2 et max. 7 catégories

Format =

["Categ1 >> Categ2"]

["Categ1 >> Categ2 >>
Categ3 >> ... >> Categ7 "]

IMAGES

1 par produit

En couleur

Différentes tailles

Généralement sur fond blanc

DESCRIPTIONS

1 par produit

En anglais

De 13 à 572 mots

MÉTHODOLOGIE GÉNÉRALE

IMAGES

Pré-traitement ★

Extraction de features ★

Réduction dimensionnelle 1 --- PCA ★

DESCRIPTIONS

Pré-traitement

Extraction de features ★

Réduction dimensionnelle 1 --- PCA ★

merge

Dataset features
images réduitDataset features images +
descriptions réduitDataset features
descriptions réduit

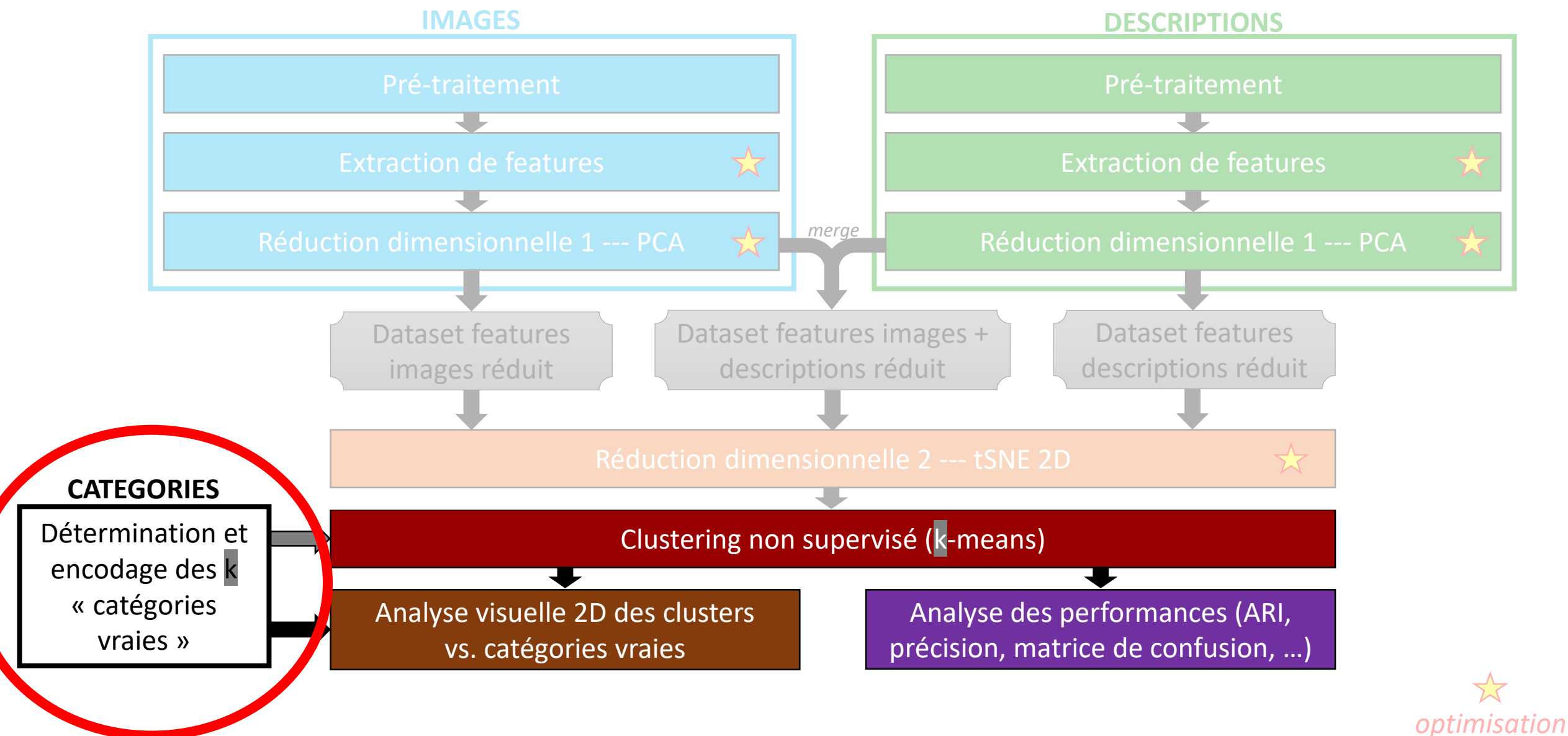
Réduction dimensionnelle 2 --- tSNE 2D ★

CATEGORIES

Détermination et
encodage des k
« catégories
vraies »

Clustering non supervisé (k-means)

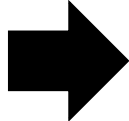
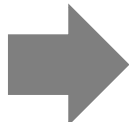
Analyse visuelle 2D des clusters
vs. catégories vraiesAnalyse des performances (ARI,
précision, matrice de confusion, ...)



Détermination et encodage des « catégories vraies »

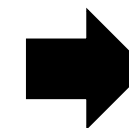
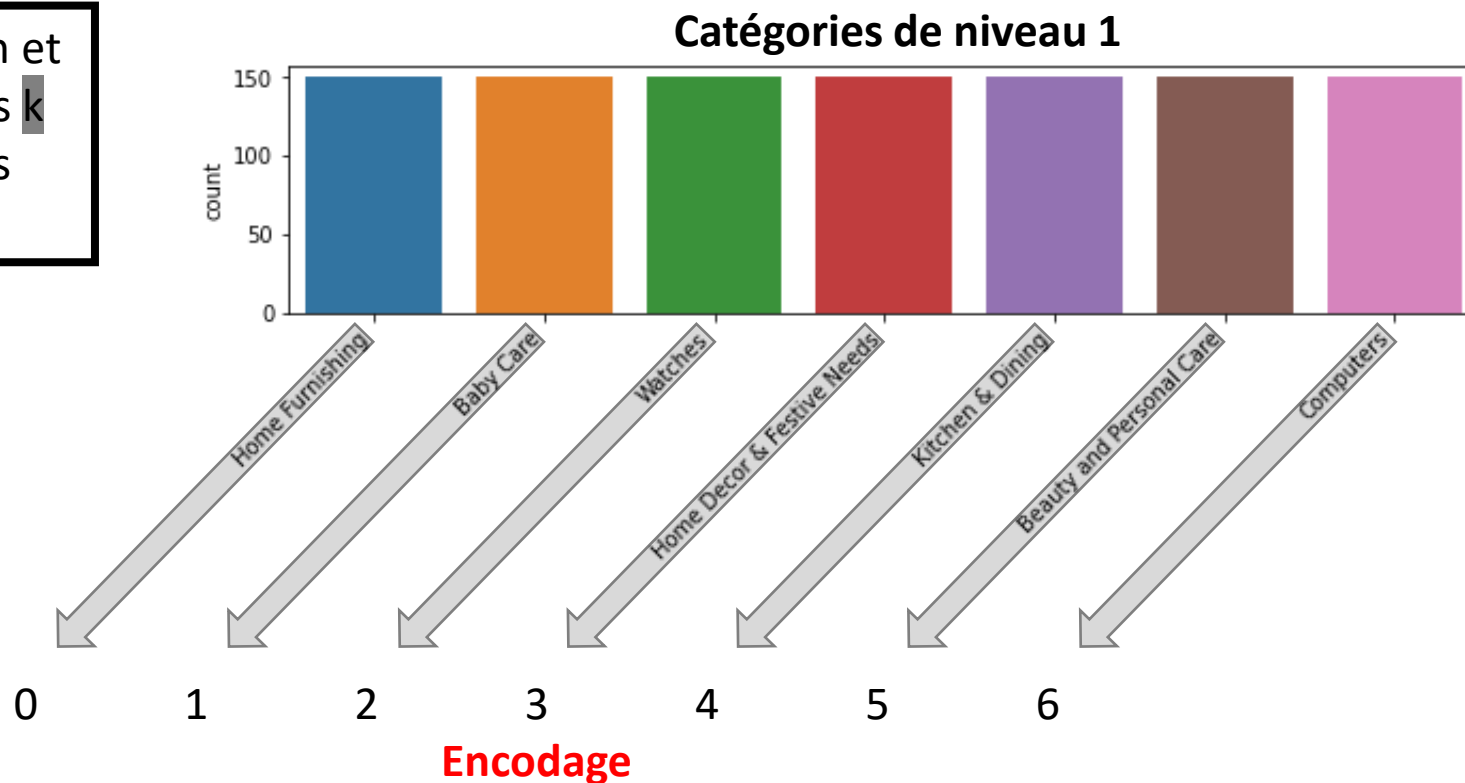


Category	Count
Home Furnishing	150
Baby Care	150
Watches	150
Home Decor & Festive Needs	150
Kitchen & Dining	150
Beauty and Personal Care	150
Computers	150

[illegible]

CATEGORIES

Détermination et
encodage des k
« catégories
vraies »



Sélection des $k=7$
catégories de niveau 1
comme « catégories
vraies » car
équilibrées

IMAGES

Pré-traitement



Extraction de features



Réduction dimensionnelle 1 --- PCA

Dataset features
images réduit

Réduction dimensionnelle 2 --- tSNE 2D



Clustering non supervisé (k-means)

Analyse visuelle 2D des clusters
vs. catégories vraiesAnalyse des performances (ARI,
précision, matrice de confusion, ...)

CATEGORIES

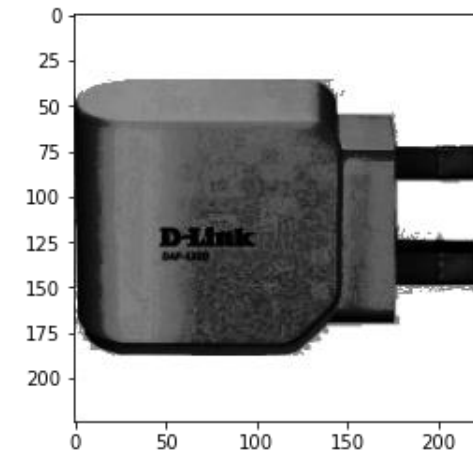
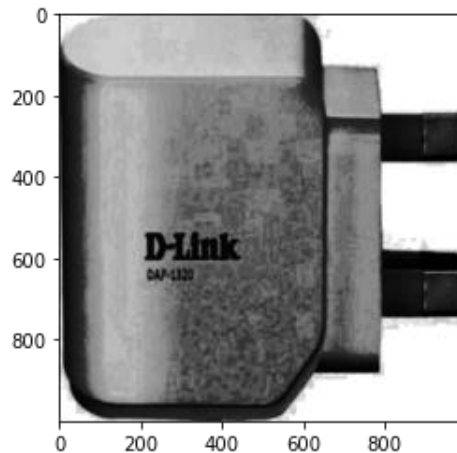
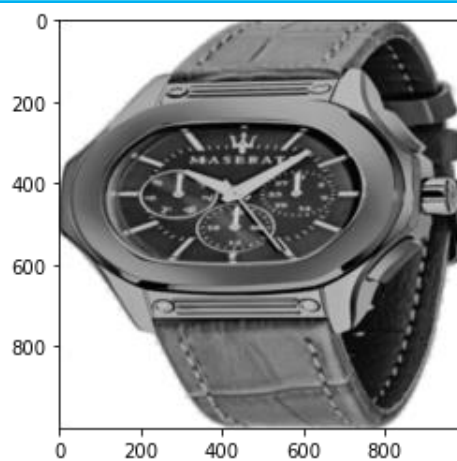
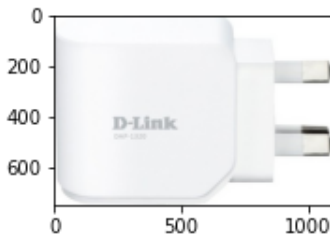
Détermination et
encodage des k
« catégories
vraies »

Pré-traitement

Redimensionnement (1000,1000)
Conversion en niveaux de gris
Egalisation histogramme (contraste)
Bruit filtré

Redim. (224,224) + padding blanc
Conversion en niveaux de gris
Egalisation histogramme (contraste)
Bruit NON filtré

Images originales



Introduction

Catégories

Images

Descriptions

Img + descript.

Conclusions

Pré-traitement

Size1000 + noise

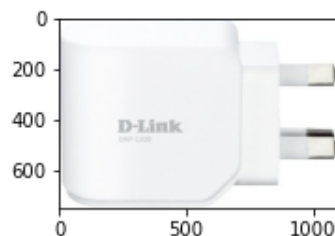
Size224 + padding

SIFT

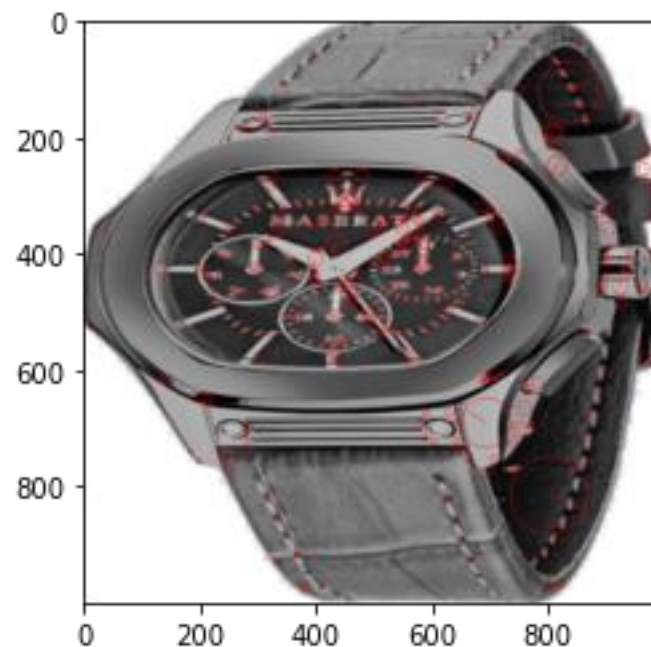
Extraction de features

Nb max descripteurs : [25,100,500,1000,5000,10000]

Images originales

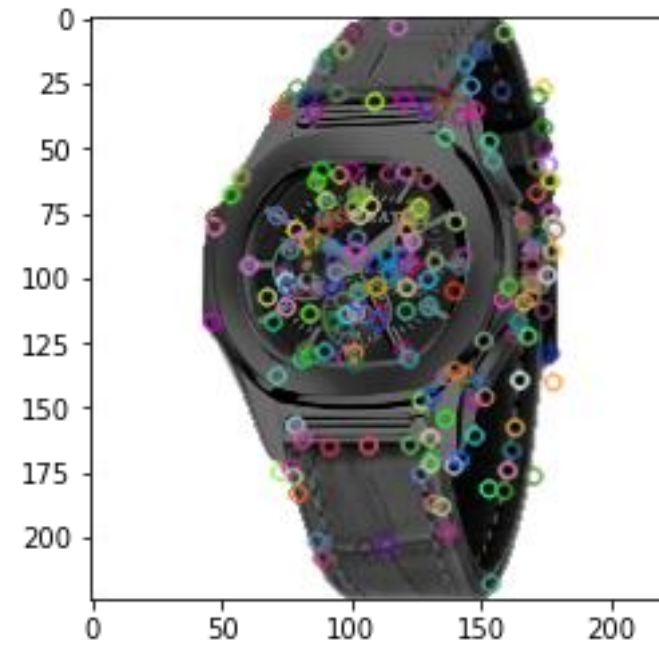


Ex: nb max descripteurs = 500



Size1000 + noise

500 descripteurs



Size224 + padding

244 descripteurs

Pré-traitement

Size1000 + noise

Size224 + padding

Extraction de features

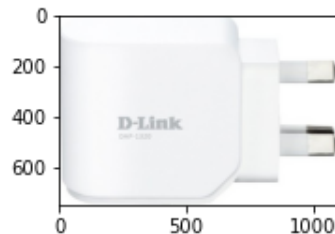
Nb max descripteurs : [25,100,500,1000,5000,10000]

Nb visual words: [k_{min} , k_{mean} , k_{max}]

SIFT

MiniBatchKMeans

Images originales



$$k_{min} = 10 \times nb \text{ catégories} = 70$$

$$k_{max} = \lfloor \sqrt{nb \text{ total de descripteurs de l'ensemble des images}} \rfloor$$

$$k_{mean} = \left\lfloor \frac{k_{min} + k_{max}}{2} \right\rfloor$$

Introduction

Catégories

Images

Descriptions

Img + descript.

Conclusions

Pré-traitement

Size1000 + noise

Size224 + padding

Extraction de features

Nb max descripteurs : [25,100,500,1000,5000,10000]

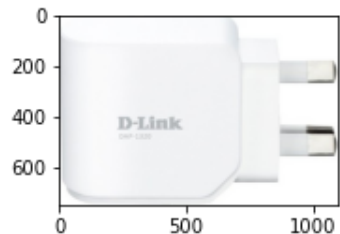
Nb visual words: [k_{min} , k_{mean} , k_{max}]

Vectorisation – Bag of visual words

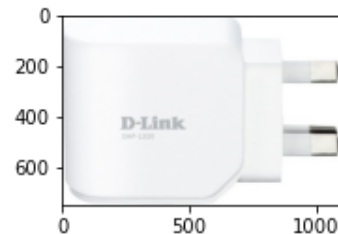
SIFT

MiniBatchKMeans

Images originales



Images originales



SIFT

MiniBatchKMeans

Pré-traitement

Size1000 + noise

Size224 + padding

Extraction de features

Nb max descripteurs : [25,100,500,1000,5000,10000]

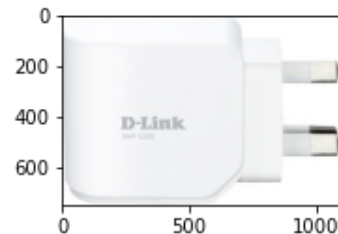
Nb visual words: [k_{min} , k_{mean} , k_{max}]

Vectorisation – Bag of visual words

Réduction dimensionnelle 1 --- PCA

Variance expliquée: [80%, 90%, 99%]

Images originales



SIFT

MiniBatchKMeans

Pré-traitement

Size1000 + noise

Size224 + padding

Extraction de features

Nb max descripteurs : [25,100,500,1000,5000,10000]

Nb visual words: [k_{min} , k_{mean} , k_{max}]

Vectorisation – Bag of visual words

Réduction dimensionnelle 1 --- PCA

Variance expliquée: [80%, 90%, 99%]

Réduction dimensionnelle 2 --- tSNE 2D

Perplexity: [2,5,10,25,50,75,100,500]

Images originales



SIFT

MiniBatchKMeans

Pré-traitement

Size1000 + noise

Size224 + padding

Extraction de features

Nb max descripteurs : [25,100,500,1000,5000,10000]

Nb visual words: [k_{min} , k_{mean} , k_{max}]

Vectorisation – Bag of visual words

Réduction dimensionnelle 1 --- PCA

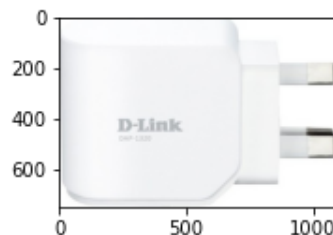
Variance expliquée: [80%, 90%, 99%]

Réduction dimensionnelle 2 --- tSNE 2D

Perplexity: [2,5,10,25,50,75,100,500]

Clustering non supervisé (k-means à k=7)

Images originales



SIFT

MiniBatchKMeans

Pré-traitement

Size1000 + noise

Size224 + padding

Extraction de features

Nb max descripteurs : [25,100,500,1000,5000,10000]

Nb visual words: [k_{min} , k_{mean} , k_{max}] =610

Vectorisation – Bag of visual words

Réduction dimensionnelle 1 --- PCA

Variance expliquée: [80%, 90%, 99%]

Réduction dimensionnelle 2 --- tSNE 2D

Perplexity: [2,5,10,25,50,75,100,500]

Clustering non supervisé (k-means à k=7)

Analyse des performances (ARI)

Meilleur ARI ≈ 0.08

Introduction

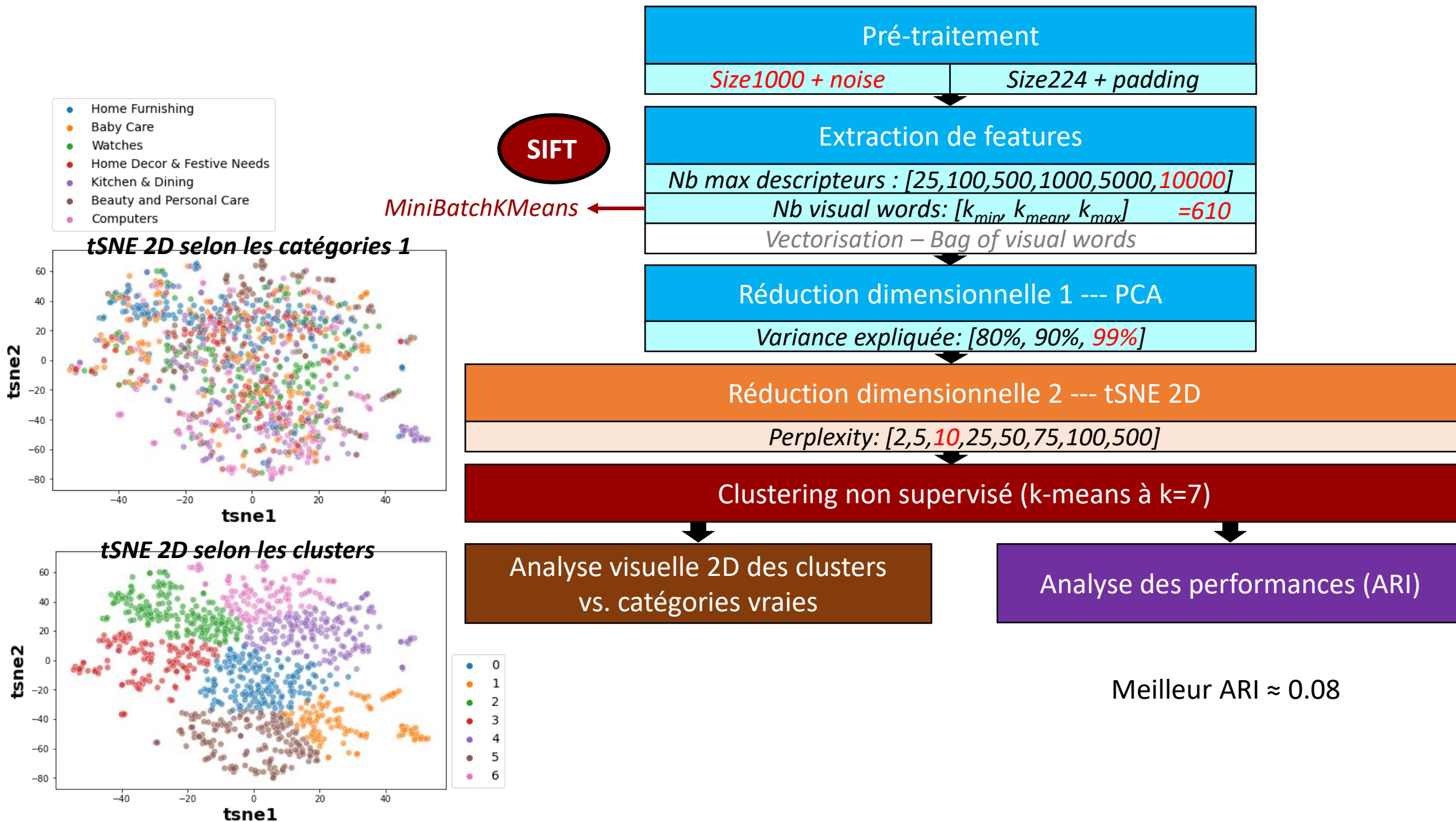
Catégories

Images

Descriptions

Img + descript.

Conclusions



Introduction

Catégories

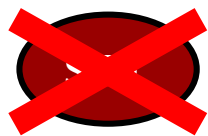
Images

Descriptions

Img + descript.

Conclusions

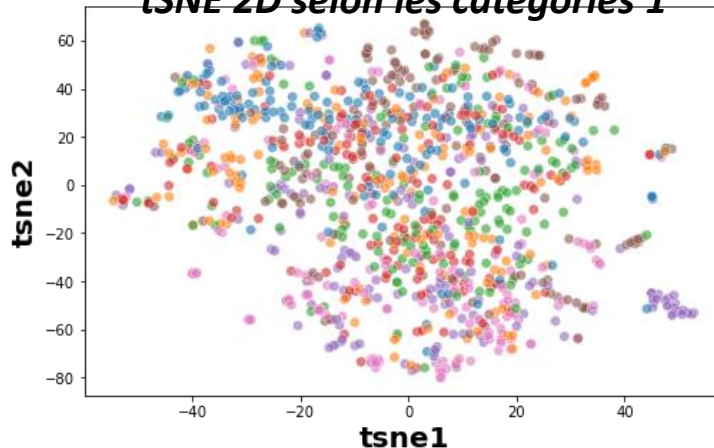
CNN (VGG16)
– Transfer Learning



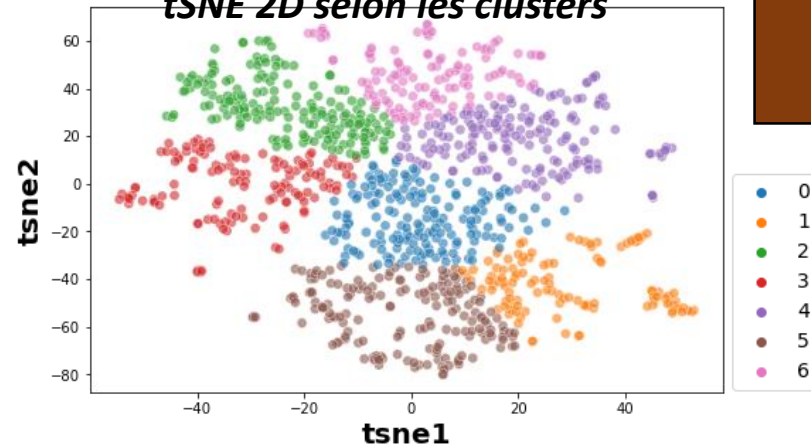
MiniBatchKMeans

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

tSNE 2D selon les catégories 1



tSNE 2D selon les clusters



Pré-traitement

Size1000 + noise

Size224 + padding

Extraction de features

Nb max descripteurs : [25,100,500,1000,5000,10000]

Nb visual words: [k_{min} , k_{mean} , k_{max}] =610

Vectorisation – Bag of visual words

Réduction dimensionnelle 1 --- PCA

Variance expliquée: [80%, 90%, 99%]

Réduction dimensionnelle 2 --- tSNE 2D

Perplexity: [2,5,10,25,50,75,100,500]

Clustering non supervisé (k-means à k=7)

Analyse visuelle 2D des clusters
vs. catégories vraies

Analyse des performances (ARI)

Meilleur ARI ≈ 0.08

Introduction

Catégories

Images

Descriptions

Img + descript.

Conclusions

VGG16

Pré-traitement

Prétraitement VGG16

Extraction de features

Base VGG16
+ Flatten

Vecteur output de taille 25088

Base VGG16 + Pooling
+ Flatten

Vecteur output de taille 4608

Réduction dimensionnelle 1 --- PCA

Variance expliquée: [80%, 90%, 99%]

Réduction dimensionnelle 2 --- tSNE 2D

Perplexity: [2,5,10,15,20,25,30,35,40,45,50,75,100,500,750] + essai 3D

Clustering non supervisé (k-means à k=7)

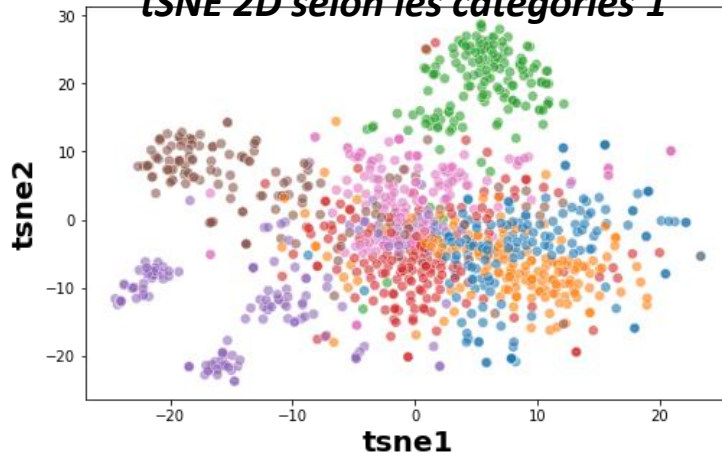
Analyse visuelle 2D des clusters
vs. catégories vraies

Analyse des performances (ARI,
précision, matrice de confusion, ...)

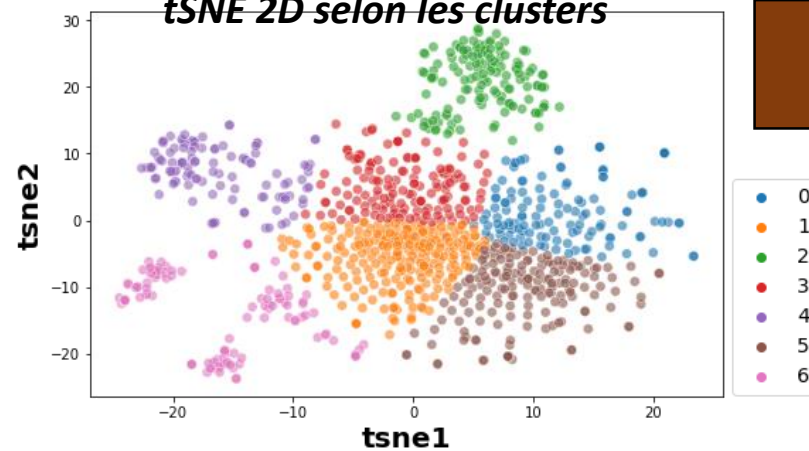
Meilleur ARI ≈ 0.42

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

tSNE 2D selon les catégories 1



tSNE 2D selon les clusters

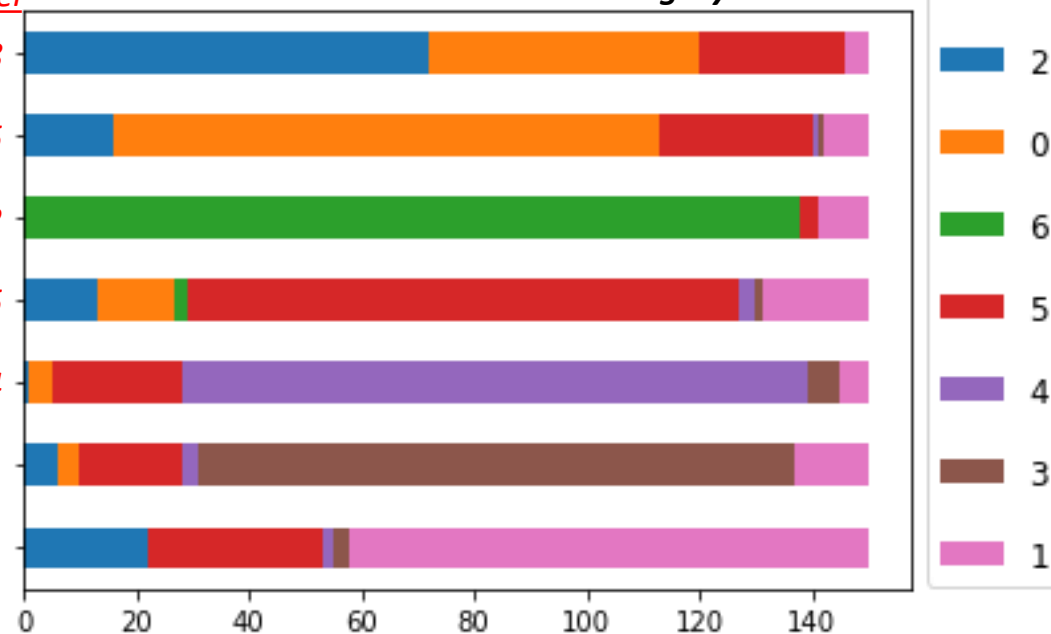


Product images clustering

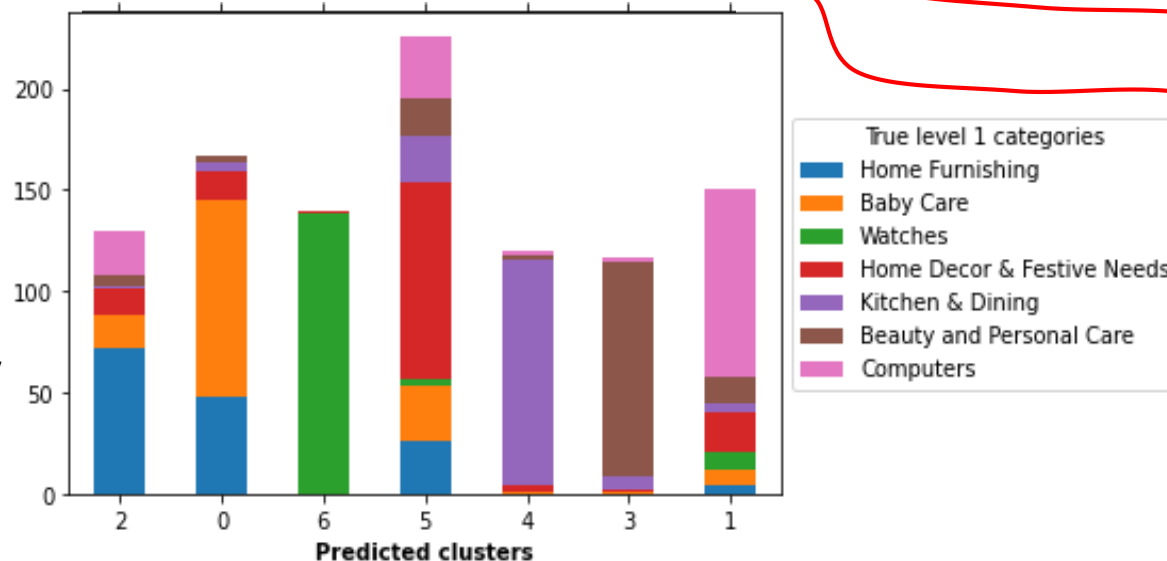
True level 1 categories

Home Furnishing	72	48	0	26	0	0	4	<i>Rappel</i> 0,48
Baby Care	16	97	0	27	1	1	8	0,65
Watches	0	0	138	3	0	0	9	0,92
Home Decor & Festive Needs	13	14	2	98	3	1	19	0,65
Kitchen & Dining	1	4	0	23	111	6	5	0,74
Beauty and Personal Care	6	4	0	18	3	106	13	0,71
Computers	22	0	0	31	2	3	92	0,61
<i>Précision</i>	0,55	0,58	0,99	0,43	0,93	0,91	0,61	

Distribution of predicted clusters within each true level 1 category



Distribution of true level 1 category within each predicted clusters



Rappel moyen = 0,68
 ➔ 68% des produits de chaque « catégorie vraie » correctement identifiés

Précision moyenne = 0,71
 ➔ 71% des produits de chaque cluster correctement classifiés

→ Bilan intermédiaire 1 (IMAGES):

- Extraction features via SIFT moins performante (et moins rapide) que via VGG16-Transfer Learning
- Meilleures performances : ARI = 0.42, rappel = 0.68, précision = 0.71

Améliorations possibles:

- Niveau dataset :
 - Regarder les produits mal classifiés pour identifier les raisons sous-jacentes et si possible les corriger avant le passage à l'échelle de la marketplace ?
- Niveau pré-traitement et extraction de features :
 - Tester ORB / SURF ?
 - Tester autre modèle CNN (+ Transfer Learning) avec images plus similaires ?
 - Ajouter une couche de pooling pour réduire encore la dimension des vecteurs features en sortie du CNN ?
- Niveau classification :
 - Autre approche Transfer Learning = ajouter un classifieur 7 classes (couche Dense, activation softmax) en fin de modèle → entraînement seulement du classifieur → prédiction des clusters ?

Introduction

Catégories

Images

Descriptions

Img + descript.

Conclusions

Pré-traitement

Tokenisation, suppr ponctuation, suppr caractères numériques, mise en minuscules, suppr English stopwords + mots fréquents, lemmatisation

Extraction de features

BoW

BoW tf-idf

bigrams tf-idf

Word2Vec

Réduction dimensionnelle 1 --- PCA

Variance expliquée: [80%, 90%, **99%**]

Nb vecteurs =
[50, 100,
200, 300, 500]

Réduction dimensionnelle 2 --- tSNE 2D

Perplexity: [2, 5, 10, 15, **20**, 25, 30, 35, 40, 45, 50, 75, 100, 500] (+ init. PCA pour W2V)

Clustering non supervisé (k-means à k=7)

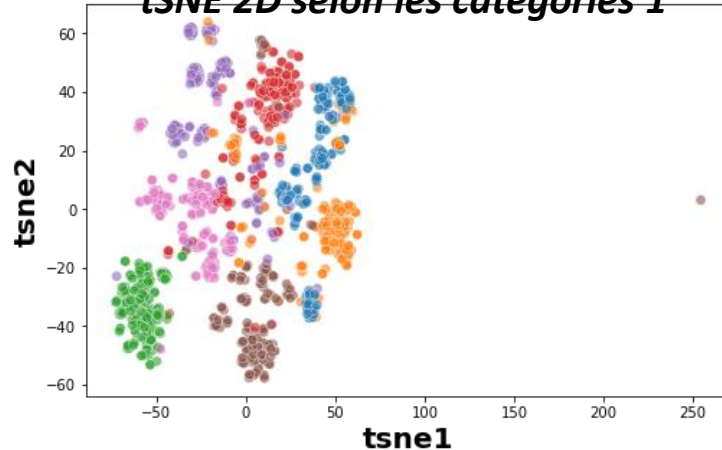
Analyse visuelle 2D des clusters
vs. catégories vraies

Analyse des performances (ARI,
précision, matrice de confusion, ...)

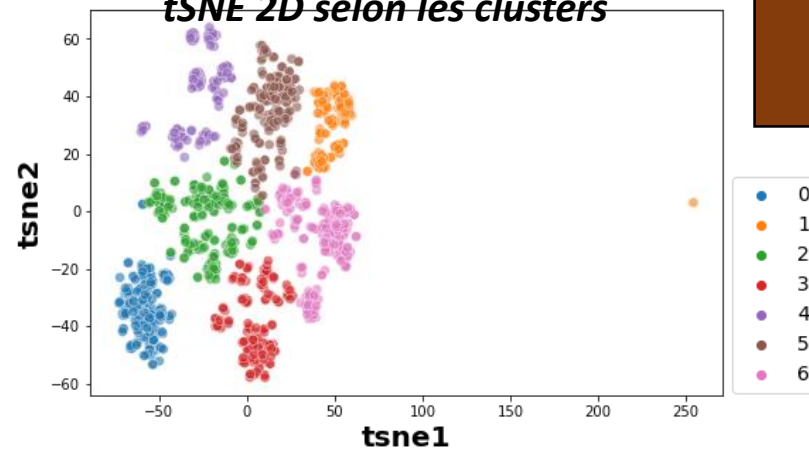
Meilleur ARI ≈ 0.63

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

tSNE 2D selon les catégories 1



tSNE 2D selon les clusters

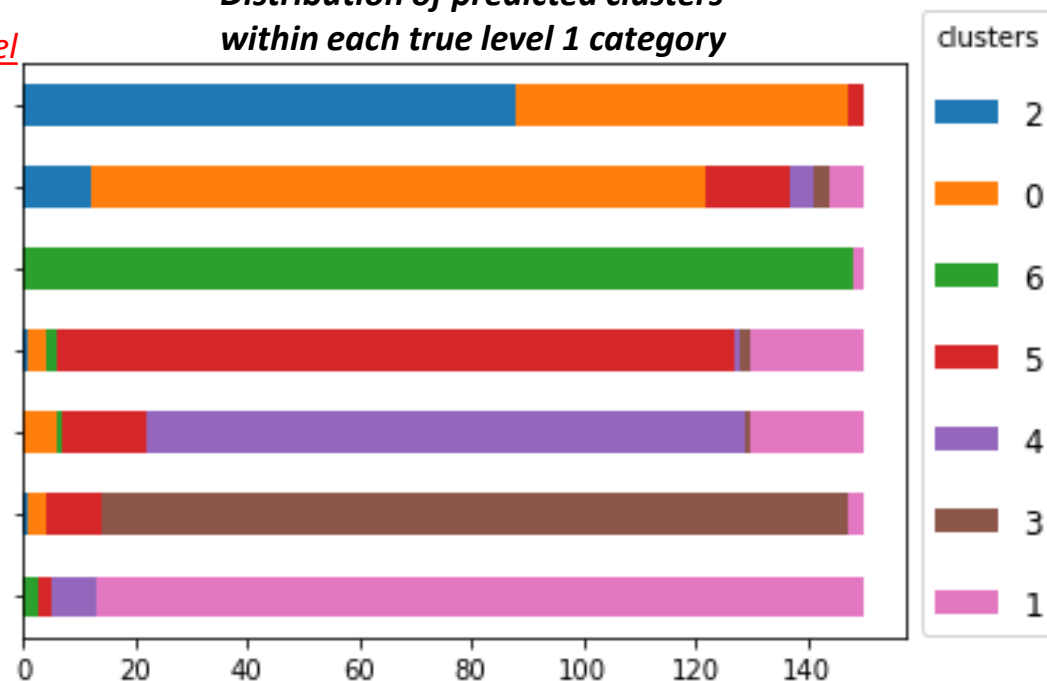


Product description clustering

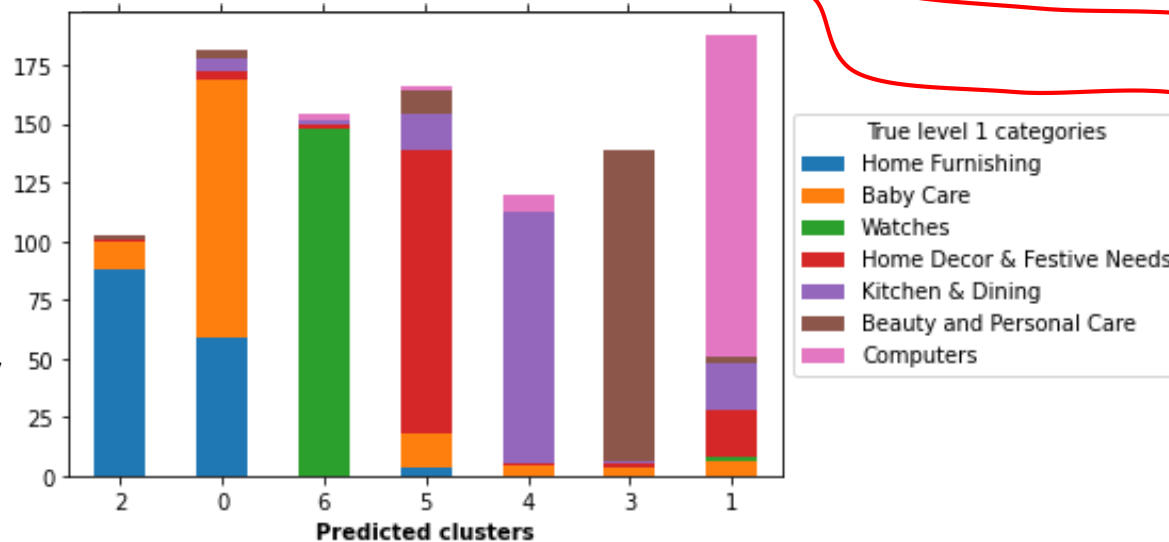
Distribution of predicted clusters within each true level 1 category

True level 1 categories

Home Furnishing	88	59	0	3	0	0	0	<i>Rappel</i> 0,59
Baby Care	12	110	0	15	4	3	6	0,73
Watches	0	0	148	0	0	0	2	0,99
Home Decor & Festive Needs	1	3	2	121	1	2	20	0,81
Kitchen & Dining	0	6	1	15	107	1	20	0,71
Beauty and Personal Care	1	3	0	10	0	133	3	0,89
Computers	0	0	3	2	8	0	137	0,91
	<i>Précision</i> 0,86	0,61	0,96	0,73	0,89	0,96	0,73	



Distribution of true level 1 category within each predicted clusters



Rappel moyen = 0,80
 ➔ 80% des produits de chaque « catégorie vraie » correctement identifiés

Précision moyenne = 0,82
 ➔ 82% des produits de chaque cluster correctement classifiés

→ Bilan intermédiaire 2 (DESCRIPTIONS):

- extraction features via Bag-of-words tf-idf plus performante que BoW classique, Bag-of-bigrams tf-idf et word embedding via Word2Vec
- Meilleures performances : ARI = 0.63, rappel = 0.80, précision = 0.82

Améliorations possibles:

- Niveau dataset :
 - Regarder les produits mal classifiés pour identifier les raisons sous-jacentes et si possible les corriger avant le passage à l'échelle de la marketplace ?
- Niveau pré-traitement et extraction de features :
 - Optimiser le pré-traitement de texte (e.g. lemmatisation vs racinisations)
 - Tester un autre algorithme de plongement de mots ?
 - Tester un modèle RNN (+ Transfer Learning) ?
- Niveau modélisation:
 - Tester clustering sur données réduites en 3D ?
 - Tester l'utilisation d'un RNN (+ Transfer Learning) pour la classification ?

IMAGES

Pré-traitement ★

Extraction de features ★

Réduction dimensionnelle 1 --- PCA ★

Dataset features
images réduit

DESCRIPTIONS

Pré-traitement

Extraction de features ★

Réduction dimensionnelle 1 --- PCA ★

Dataset features
descriptions réduit

Réduction dimensionnelle 2 --- tSNE 2D ★

Clustering non supervisé (k-means à k=7)

Analyse visuelle 2D des clusters
vs. catégories vraiesAnalyse des performances (ARI,
précision, matrice de confusion, ...)

Meilleurs ARIs = 0,42 | 0,63

CATEGORIES

Détermination et
encodage des 7
« catégories
vraies »

IMAGES

Pré-traitement ★

Extraction de features ★

Réduction dimensionnelle 1 --- PCA ★

DESCRIPTIONS

Pré-traitement

Extraction de features ★

Réduction dimensionnelle 1 --- PCA ★

merge

Dataset features
images réduitDataset features images +
descriptions réduitDataset features
descriptions réduit

Réduction dimensionnelle 2 --- tSNE 2D ★

Clustering non supervisé (k-means à k=7)

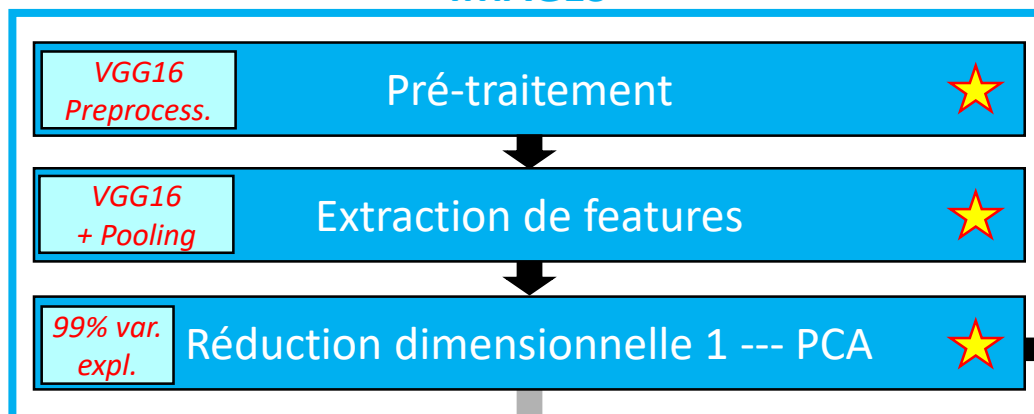
Analyse visuelle 2D des clusters
vs. catégories vraiesAnalyse des performances (ARI,
précision, matrice de confusion, ...)

Meilleurs ARIs = 0,42 | 0,63

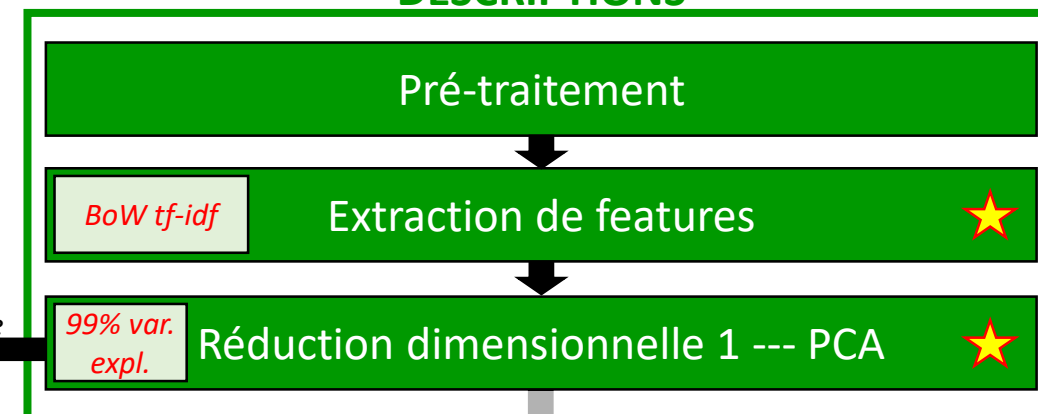
CATEGORIES

Détermination et
encodage des 7
« catégories
vraies »

IMAGES



DESCRIPTIONS



merge

Dataset features
images réduitDataset features images +
descriptions réduitDataset features
descriptions réduit

Réduction dimensionnelle 2 --- tSNE 2D ★

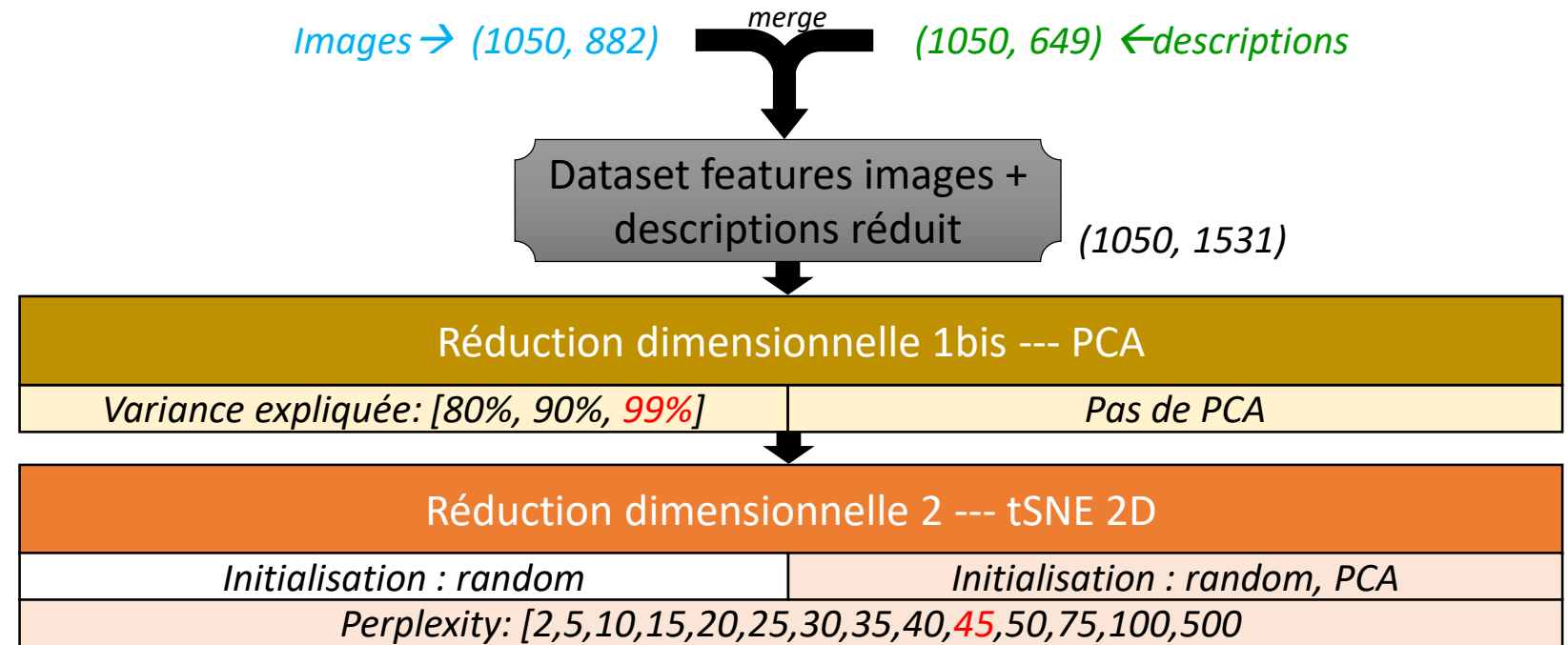
Clustering non supervisé (k-means à k=7)

Analyse visuelle 2D des clusters
vs. catégories vraiesAnalyse des performances (ARI,
précision, matrice de confusion, ...)

CATEGORIES

Détermination et
encodage des 7
« catégories
vraies »

Meilleurs ARIs = 0,42 | 0,63



Introduction

Catégories

Images

Descriptions

Img + descript.

Conclusions

Images $\rightarrow (1050, 882)$ $\xrightarrow{\text{merge}}$ $(1050, 649) \leftarrow$ descriptions

Dataset features images +
descriptions réduit

$(1050, 1531)$

Réduction dimensionnelle 1bis --- PCA

Variance expliquée: [80%, 90%, **99%**]

Pas de PCA

Réduction dimensionnelle 2 --- tSNE 2D

Initialisation : **random**

Initialisation : random, PCA

Perplexity: [2,5,10,15,20,25,30,35,40,**45**,50,75,100,500]

Clustering non supervisé (k-means à k=7)

Analyse visuelle 2D des clusters
vs. catégories vraies

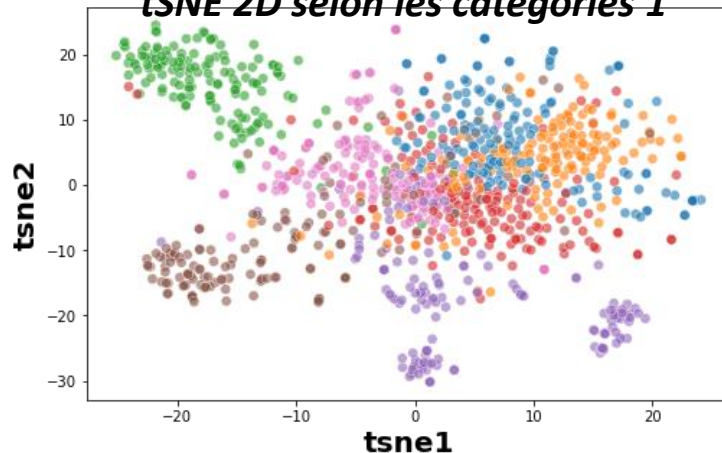
Analyse des performances (ARI,
précision, matrice de confusion, ...)

Meilleur ARI ≈ 0.43

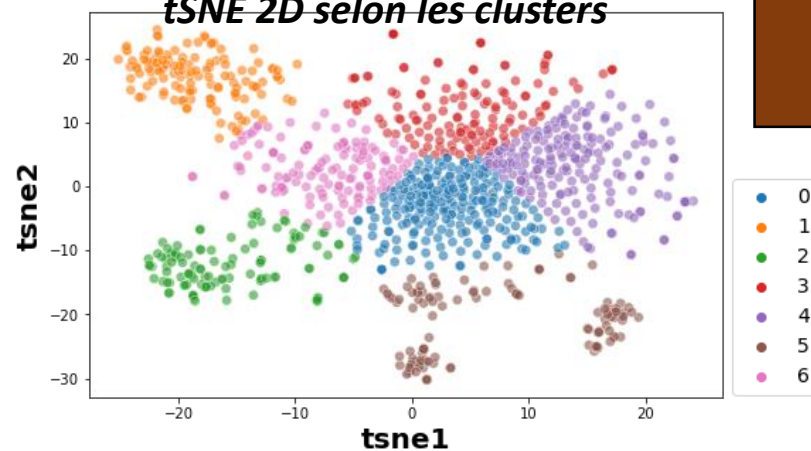
Meilleurs ARIs = **0,42** | **0,63**

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

tSNE 2D selon les catégories 1



tSNE 2D selon les clusters



Products description + image clustering

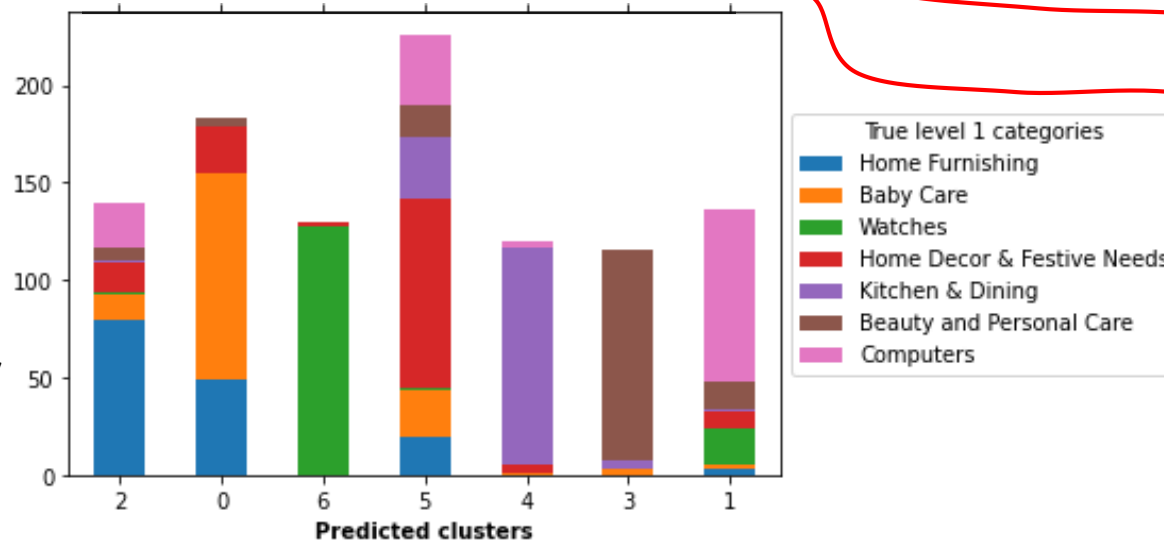
Rappel

Home Furnishing	79	49	0	19	0	0	3
Baby Care	14	106	0	24	1	3	2
Watches	1	0	128	2	0	0	19
Home Decor & Festive Needs	15	24	2	97	4	0	8
Kitchen & Dining	1	0	0	31	112	4	2
Beauty and Personal Care	7	4	0	17	0	108	14
Computers	22	0	0	36	3	1	88

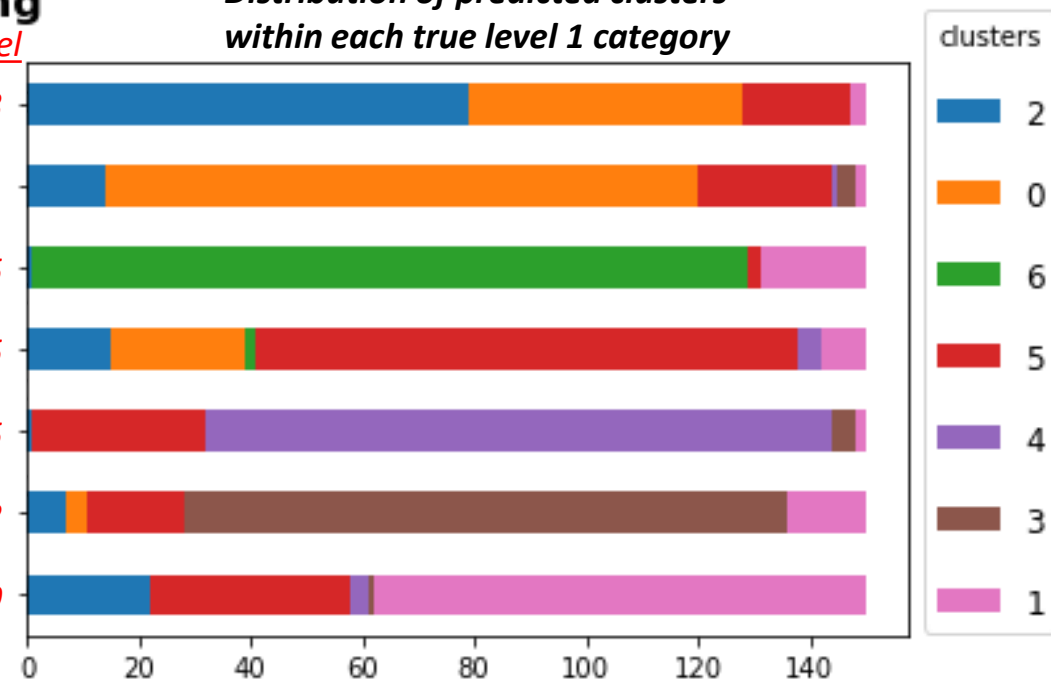
Précision

0,57 0,58 0,98 0,43 0,93 0,93 0,65

Distribution of
true level 1 category
within each
predicted clusters



Distribution of predicted clusters
within each true level 1 category



vs. 0,68 | 0,80

Rappel moyen = 0,68
→ 68% des produits de
chaque « catégorie
vraie » correctement
identifiés

Précision moyenne = 0,72 vs. 0,71 | 0,82
→ 72% des produits de chaque
cluster correctement classifiés

→ Bilan intermédiaire 3 (IMAGES + DESCRIPTIONS):

- Performances non améliorées, équivalentes à celles obtenues sur la base des images seules

Données	Meilleur ARI	Meilleur rappel	Meilleure précision
Images	0,42	0,68	0,71
Descriptions	0,63	0,80	0,82
Images + descriptions	0,43	0,68	0,72

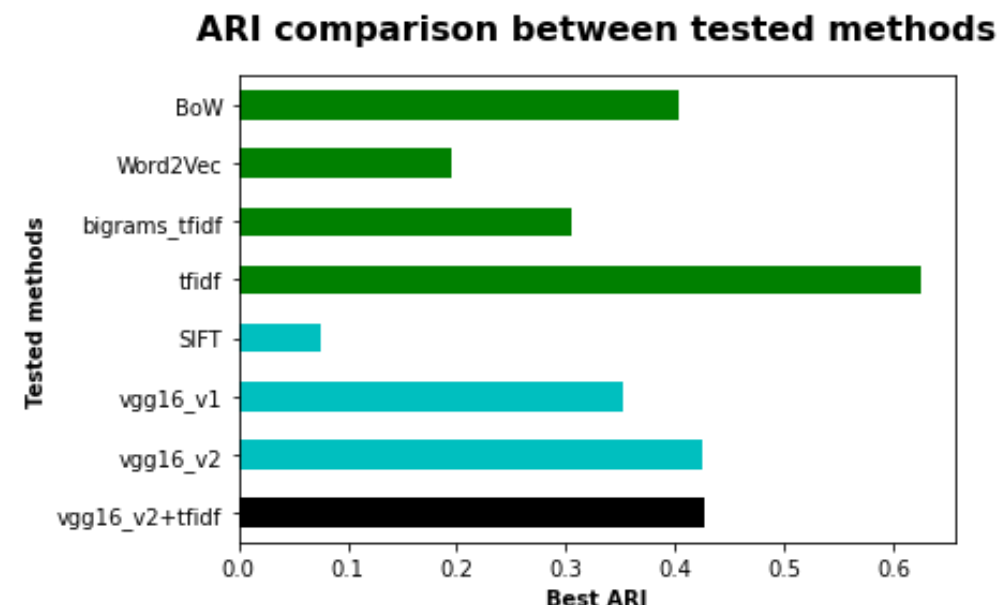
Améliorations possibles, en plus de celles spécifiques aux images et aux descriptions déjà mentionnées:

- Niveau pré-traitement et extraction de features :
 - Tester d'autres combinaisons de méthodes (pas uniquement les plus performantes séparément) ?
- Niveau modélisation:
 - Tester clustering sur données réduites en 3D ?

→ Conclusions générales

- Meilleures performances obtenues sur classification basée sur descriptions seules

Données	Meilleur ARI	Meilleur rappel	Meilleure précision
Images	0,42	0,68	0,71
Descriptions	0,63	0,80	0,82
Images + descriptions	0,43	0,68	0,72



- Suffisamment bonnes pour attester de la faisabilité (avec une bonne précision) :
 - d'une automatisation non supervisée (après tests supplémentaires et améliorations)
 - d'une automatisation **supervisée**, avec un **jeu de données plus conséquent** (après passage à l'échelle de la marketplace), de l'attribution de catégories aux produits.

Améliorations possibles, en plus de celles déjà mentionnées:

- Vérifier la validité des « catégories vraies » actuelles / détecter l'occurrence éventuelle de catégories plus adaptées ?
- Classification « manuelle » des descriptions via l'identification de mots-clés spécifiques pour chaque catégorie ?

Baby Care



Home Furnishing



Watches



Home Decor & Festive Needs



Kitchen & Dining



Beauty and Personal Care



Computers

