Projet 6 - Classifiez automatiquement des biens de consommation

Classification des articles d'une marketplace en les catégorisant automatiquement, sur la base de la description et de l'image fournie.



Sommaire



La problématique & Le jeu de données



Prétraitement, extraction de Features sur la partie texte



Prétraitement, extraction de Features & Etude de faisabilité sur la partie image



Classification supervisée image



Test de l'API



RGPD





Contexte: "Place de marché" souhaite lancer une marketplace sur laquelle les vendeurs peuvent proposer des produits à la vente en fournissant une photo et une description. Les vendeurs choisissent manuellement la catégorie du produit. Pour augmenter en fiabilité Linda, Lead DS, me confie la mission d'étudier la possibilité d'automatiser cette tâche d'attribution de catégories.

Objectif : Etudier la faisabilité d'un moteur de classification des articles en différentes catégories, avec un niveau de précision suffisant. Dans l'objectif de gagner en fiabilité et en temps.

place de marché

Moyens: Une sélection d'articles déjà catégorisées.



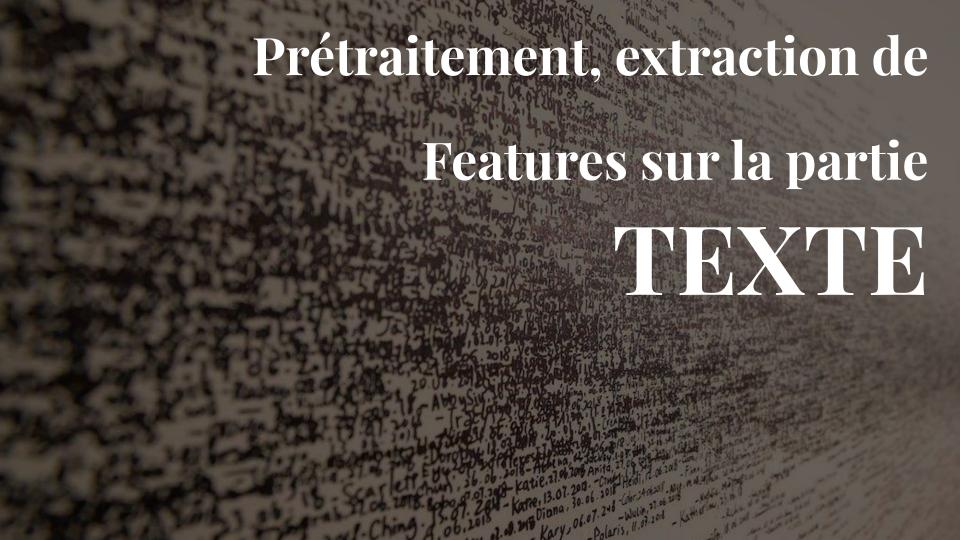
Les données - un fichier CSV & images associées

Catégories	Quantité
Home Furnishing	150
Baby Care	150
Watches	150
Home Decor & Festive Needs	150
Kitchen & Dining	150
Beauty and Personal Care	150
Computers	150
Total	1050

Les variables à utiliser :

- → product_name
- → product_category_tree ⇒ category
- → image ⇒ image_path (des images fournies)
- → description

Taux de remplissage de 100% pour ces variables!





Prétraitement du TEXTE

Texte = product_name + description

Elegance Polyester Multicolor Abstract Eyelet Door Curtain Key Features of Elegance Polyester Multicolor Abstract Eyelet Door Curtain Floral Curtain, Elegance Polyester Multicolor Abstract Eyelet Door Curtain (213 cm in Height, Pack of 2) Price: Rs. 899 This curtain enhances the look of the interior...

- Fonction de preprocessing :
- Conversion en minuscules & Suppression de la ponctuation
- Tokenisation des mots
- Suppression des stop words & des mots de moins de 3 lettres

- Stemming
- Suppression des caractères non alphabétiques
- Suppression des chaînes de texte vides

['eleg','polyest','multicolor','abstract','eyelet','door','curtain','key','featur','eleg','polyest','multicolor','abstract','eyelet','door','curtain','height','pack','pric','curtain','enhanc','look','interiors th',...

 Suppression du vocabulaire des mots qui apparaissent moins de 6 fois et ceux qui apparaissent trop souvent



Les méthodes testées pour représenter le texte

USE

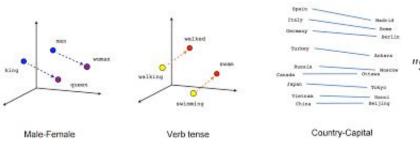
DEEP LEARNING

MÉTHODES PLUS TRADITIONNELLES

Bag of words et TF-IDF (fréquence, stats.)

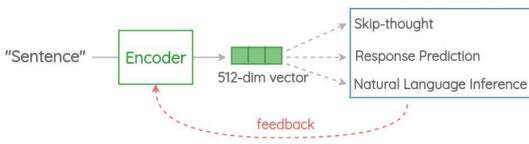
	1 This	2 movie	3 is	4 very	5 scary	6 and	7 long	8 not	9 slow	10 spooky	11 good	Length of the review(in words)
Review 1	1	1	1	1	1	1	1	0	0	0	0	7
Review 2	1	1	2		$w = tf \times log(\frac{N}{N})$					0	8	
Review 3	1	1	1	'	$W_{x,y} = tt_{x,y} \times$				$y \times \log \left(df_{x} \right)$			6
	L	l			TF-IDI		tf _{x,y} = frequency df _x = numbe N = total number	r of docum	ents containi	ng x		

Word vectors (doc2vec, word2vec) (peu profonds



MNLI NER SQUAD Mask LM Start/End Span C T, T_N T_(SEP) BERT BERT Question Answer Pair Fine-Tuning Pre-training

Multi-task Learning



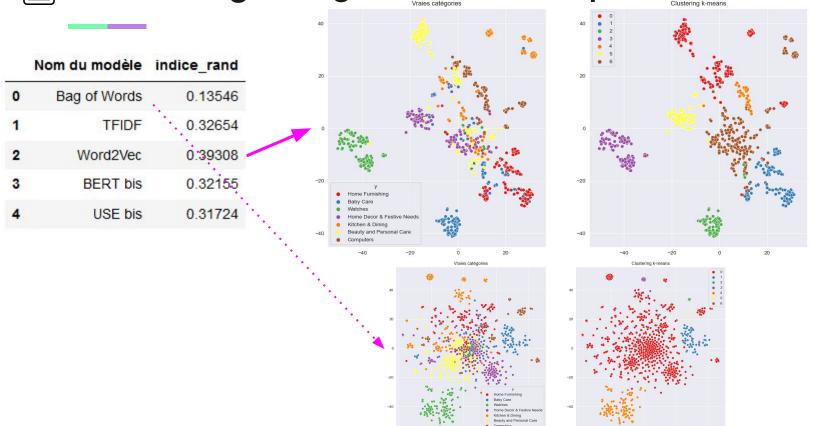


Méthodologie: Segmentation & Comparaison des résultats

- Utilisation du feature engineering présenté
- K-means avec 7 clusters sur les features adaptés pour chaque modèle
- Réprésentation visuelle via TSNE à 2 dimensions
- Comparaison vraies catégories et segmentation
- Indice de Rand pour comparer les résultats des différentes méthodes

T

Méthodologie : Segmentation & Comparaison des résultats

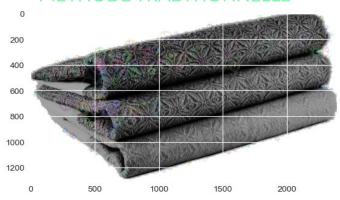






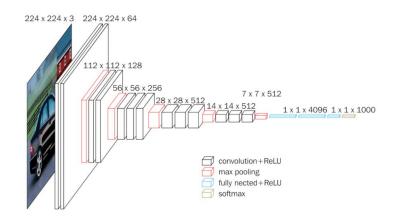
2 Méthodes, modèles et preprocessing associé

MÉTHODE TRADITIONNELLE



- Resize
- Nuances de gris
- Egalisation des histogrammes
- descripteurs d'image SIFT
- créations de clusters de descripteurs (625)
- Création histogramme par image

VGG16 PRE ENTRAÎNE



- Resize
- Turn into array
- **Expand dimension**
- Fonction de preprocessing dédiée à VGG16



Méthodologie: Segmentation & Comparaison des résultats



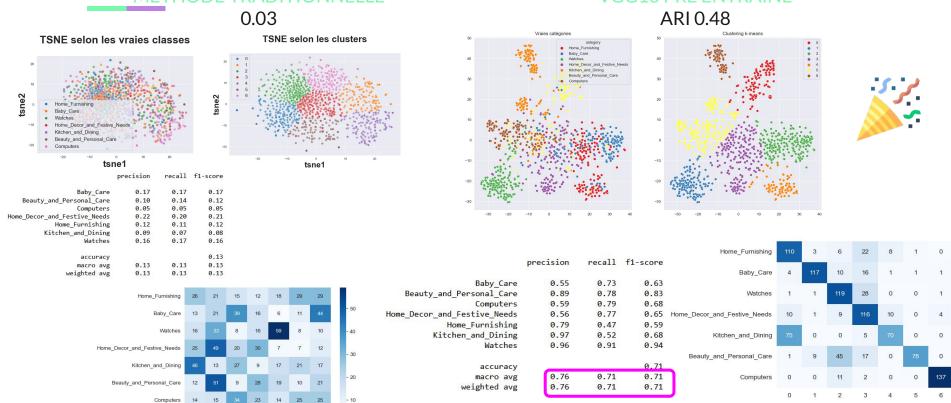
- Preprocessing des images selon méthode
- Utilisation des features dégagées selon chacune des méthodes
- Réduction de dimension PCQ n_components=0.99
- T-SNE à n = 2
- K-means avec 7 clusters sur les features le T-SNE
- Représentation visuelle via TSNE à 2 dimensions
- Comparaison vraies catégories et segmentation
- Indice de Rand, Accuracy, Matrice de confusion



Méthodologie : Segmentation & Comparaison des résultats



VGG16 PRE ENTRAÎNE

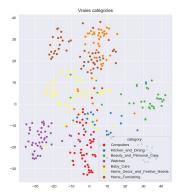


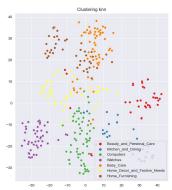


3 méthodes testées

Modèle SANS data augmentation sans entraînement

- Preprocessing VGG16 cf faisabilité
- PCA n=0,99
- TSNE n=2
- Train & Test Split (0.3)
- Knn plus proches voisins
- Accuracy = 0.80





Modèle SANS data augmentation avec entraînement

- Train, Test, Validation Split (0.7, 0.15, 0.15)
- Preprocessing VGG16 cf faisabilité
- Label encoding
- Chargement de VGG16 et conservation des top layers pour entraînement des couches de classification
- Accuracy = 0.78



Modèle AVEC data augmentation avec entraînement

- Train, Test, Validation Split (0.7, 0.15, 0.15)
- Création de 2 datagen pour train_generator
 & valid_generator (label encoding & preprocessing inclus)
- Chargement de VGG16 et conservation des top layers pour entraînement des couches de classification
- Accuracy = 0.63





Quelle méthode choisir?



Modèle SANS data augmentation sans entraînement

- Preprocessing VGG16 cf faisabilité
- PCA n=0,99
- TSNE n=2
- Train & Test Split (0.3)
- Knn plus proches voisins
- Accuracy = 0.80

Problème de dimensionnalité, selon le nombre de données qui alimenteraient le modèle

Beaucoup d'étapes

Modèle SANS data augmentation avec entraînement

- Train, Test, Validation Split (0.7, 0.15, 0.15)
- Preprocessing VGG16 cf faisabilité
- Label encoding
- Chargement de VGG16 et conservation des top layers pour entraînement des couches de classification
- Accuracy = 0.78



L'entraînement n'améliore pas les résultats

Bien que bons, ils sont peut être liés à un mauvais choix des données

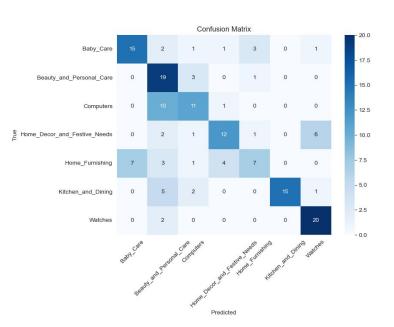
Modèle AVEC data augmentation avec entraînement

- Train, Test, Validation Split (0.7, 0.15, 0.15)
- Création de 2 datagen pour train_generator
 & valid_generator (label encoding & preprocessing inclus)
- Chargement de VGG16 et conservation des top layers pour entraînement des couches de classification
- Accuracy = 0.63



L'entraînement améliore les résultats

Rapport sur la méthode choisie



Classification Report:

	precision	recall	f1-score	support
Baby_Care	0.68	0.65	0.67	23
Beauty_and_Personal_Care	0.44	0.83	0.58	23
Computers	0.58	0.50	0.54	22
Home_Decor_and_Festive_Needs	0.67	0.55	0.60	22
Home_Furnishing	0.58	0.32	0.41	22
Kitchen_and_Dining	1.00	0.65	0.79	23
Watches	0.71	0.91	0.80	22
accuracy			0.63	157
macro avg	0.67	0.63	0.63	157
weighted avg	0.67	0.63	0.63	157

Exemples d'erreurs:



Prédit Home_Furnishing à la place de Baby_Care



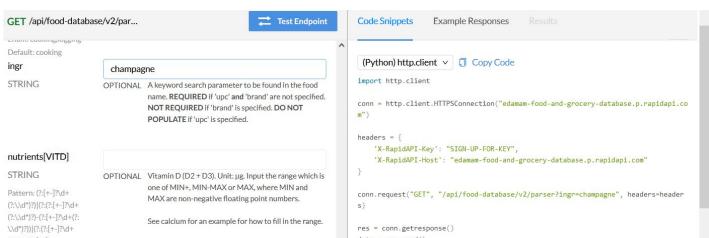




Contexte : Elargir notre gamme de produits, en particulier dans l'épicerie fine.

Objectif: Tester la collecte de produits à base de "champagne" et fournir une extraction des 10 premiers produits dans un fichier ".csv", contenant pour chaque produit les données suivantes : foodld, label, category, foodContentsLabel, image.

Moyens: via l'API EDAMAM.



Le fichier CSV

Α	В	С	D	E	F	G	Н
	foodId	label	category	foodContentsLabel	image		
0	food_a656mk2a5dmqb2adia	Champagne	Generic foods	None	https://www.edamam.com		
1	food_a656mk2a5dmqb2adia	Champagne	Generic foods	None	https://w	n.com/fo	
2	food_b753ithamdb8psbt0w	Champagne Vinaigrette, Cham	Packaged foods	OLIVE OIL; BALSAMIC VINEG	None		
3	food_b3dyababjo54xobm6r	Champagne Vinaigrette, Cham	Packaged foods	INGREDIENTS: WATER; CANO	https://w	ww.edaman	n.com/fo
4	food_a9e0ghsamvoc45bwa	Champagne Vinaigrette, Cham	Packaged foods	CANOLA AND SOYBEAN OIL;	None		
5	food_an4jjueaucpus2a3u1n	Champagne Vinaigrette, Cham	Packaged foods	WATER; CANOLA AND SOYBE	None		
6	food_bmu5dmkazwuvpaa5p	Champagne Dressing, Champa	Packaged foods	SOYBEAN OIL; WHITE WINE	https://w	ww.edaman	n.com/fo
7	food_alpl44taoyv11ra0lic1q	Champagne Buttercream	Generic meals	sugar; butter; shortening; var	None		
8	food_byap67hab6evc3a0f9v	Champagne Sorbet	Generic meals	Sugar; Lemon juice; brandy; (None		
9	food_am5egz6aq3fpjlaf8xpl	Champagne Truffles	Generic meals	butter; cocoa; sweetened co	None		





Les grands principes :

- Licéité, loyauté et transparence : données collectées de manière légale, transparente et avec le consentement des utilisateurs.
- Limitation des finalités
- Minimisation des données
- Exactitude des données : maintenues à jour et exactes.
- Limitation de la conservation

Mail de Linda concernant la PI:

Linda

PS : J'ai bien vérifié qu'il n'y avait aucune contrainte de propriété intellectuelle sur les données et les images.

L'application au projet :



Ici: photos & descriptions sur une marketplace

Seulement des données non personnellement identifiables ?





- consentement des utilisateurs sur l'utilisation des données.
- consigne de ne pas mettre de photos identifiables (ex visage, identité ...)
- ni de description avec des données personnelles (ex numéro de téléphone, adresse, identité...)
- suppression des images et texte une fois que les modèles sont entraînés (indépendamment de leur présence sur le site)

Conclusions & Remarques

Ce que nous soumettons :

- La méthode de classification supervisée avec data augmentation avec entraînement du modèle semble être la plus prometteuse. Il faudrait étendre le test à un plus grand nombre de e données pour être sûr.

Si c'était à refaire:

- Pour le test split, j'utiliserais des images que j'ai extraites à l'aide de web scraping depuis Google Images, afin d'obtenir de nouvelles données et d'utiliser la totalité de mon jeu de données, qui est de taille très réduite.

Point de vigilance:

- Il faut effectuer de la veille régulièrement sur les nouvelles méthodes d'analyse d'image

=> Slide: Veille et Proof of Concept d'une récente technique?

Le plus important : j'ai gagné en compétences

