

P6 : CLASSIFIEZ AUTOMATIQUEMENT DES BIENS DE CONSOMMATION

October 5, 2021

Carlos SACRISTAN

OpenClassrooms

MISSION OBJECTIVE

place de marché

Context :

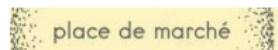
- 1 L'entreprise "Place de marché" souhaite lancer une marketplace e-commerce
- 2 Des vendeurs proposent des articles à des acheteurs en postant une photo et une description

Objective :

- 1 A partir de la photo et de la description, réaliser une première étude de faisabilité d'un moteur de **classification**

classification

TABLE DES MATIÈRES



- 1 Mission objective
- 2 Data preparation
 - Data cleaning
- 3 NLP (Text)
 - Pre-processing
 - Modeling
- 4 Machine Learning (CV)
 - Pre-processing
 - Feature extraction
- 5 Deep Learning (CV)
 - VGG16
- 6 Conclusions

INDEX

 place de marché 

1 Mission objective

2 Data preparation

- Data cleaning

3 NLP (Text)

- Pre-processing

- Modeling

4 Machine Learning (CV)

- Pre-processing
- Feature extraction

5 Deep Learning (CV)

- VGG16

6 Conclusions

DATA CLEANING

place de marché

1 Basic cleaning

- Import 1 csv (1050 rows × 14 columns)
- Lowercases of features name
- Convert features to best possible dtypes
- Split : product_category_tree (>>) → Level 0 → **TARGET**
- Delete of duplicate values
- Delete of missing values

INDEX

place de marché

- 1 Mission objective
- 2 Data preparation
 - Data cleaning
- 3 NLP (Text)
 - Pre-processing
 - Modeling
- 4 Machine Learning (CV)
 - Pre-processing
 - Feature extraction
- 5 Deep Learning (CV)
 - VGG16
- 6 Conclusions

PRE-PROCESSING

place de marché

1 Preprocessing pipeline

- Tokenization
- StopWords
- Lemmatization
- Stemming

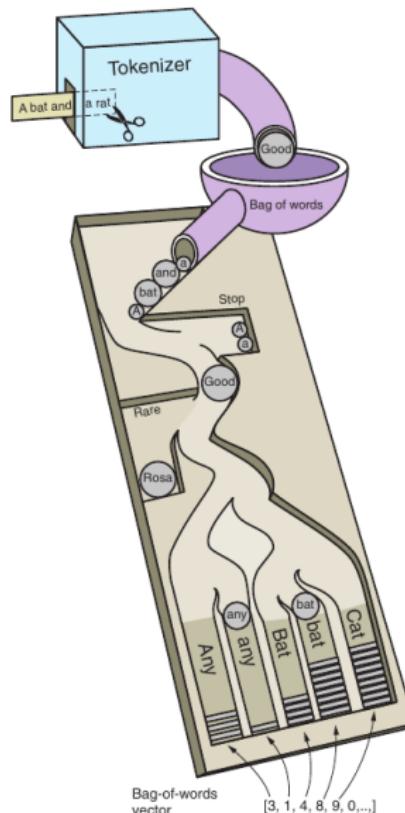
1 Bag of words

- Singel-gram
- Range-Gram

1 TF-IDF

- Singel-gram
- Range-Gram

1 T-SNE



PRE-PROCESSING


place de marché


1 Preprocessing pipeline

- Tokenization
- StopWords
- Lemmatization
- Stemming

1 Bag of words

- Singel-gram
- Range-Gram

1 TF-IDF

- Singel-gram
- Range-Gram

1 T-SNE

ORIGINAL

specifications of sathiya's cotton bath towel _3 bath towel, red, yellow, blue_ bath towel features machine washable yes material cotton design self design general brand sathiya's type bath towel gsm 500 model name sathiya's cotton bath towel ideal for men, women, boys, girls model id asvtwl322 color red, yellow, blue size medium dimensions length 30 inch width 60 inch in the box number of contents in sales package 3 sales package 3 bath towel

TOKENIZATION

specifications of sathiya's cotton bath towel bath towel red yellow blue bath towel features machine washable yes material cotton design self design general brand sathiya's type bath towel gsm model name sathiya's cotton bath towel ideal for men women boys girls model id asvtwl color red yellow blue size medium dimensions length inch width inch in the box number of contents in sales package sales package bath towel

Removal of special characters/punctuations & case folding

PRE-PROCESSING

place de marché

1 Preprocessing pipeline

- Tokenization
- StopWords
- Lemmatization
- Stemming

1 Bag of words

- Singel-gram
- Range-Gram

1 TF-IDF

- Singel-gram
- Range-Gram

1 T-SNE

TOKENIZATION

specifications of sathiya's cotton bath towel bath towel red yellow blue bath towel features machine washable yes material cotton design self design general brand sathiya's type bath towel gsm model name sathiya's cotton bath towel ideal for men women boys girls model id asvtwl color red yellow blue size medium dimensions length inch width inch in the box number of contents in sales package sales package bath towel
REMOVE_STOPWORDS

```
[['specifications', 'sathiya', 'cotton', 'bath', 'towel', 'bath', 'towel', 'red', 'yellow', 'blue', 'bath', 'towel', 'features', 'machine', 'washable', 'yes', 'material', 'cotton', 'design', 'self', 'design', 'general', 'brand', 'sathiya', 'type', 'bath', 'towel', 'gsm', 'model', 'name', 'sathiya', 'cotton', 'bath', 'towel', 'id', 'ideal', 'men', 'women', 'boys', 'girls', 'model', 'id', 'asvtwl', 'color', 'red', 'yellow', 'blue', 'size', 'medium', 'dimensions', 'length', 'inch', 'width', 'inch', 'box', 'number', 'contents', 'sales', 'package', 'sales', 'package', 'bath', 'towel']]
```

Stop words : 'who', 'what', 'when', 'why', 'how', 'which', 'where', 'whom'

PRE-PROCESSING

place de marché

1 Preprocessing pipeline

- Tokenization
- StopWords
- Lemmatization
- Stemming

1 Bag of words

- Singel-gram
- Range-Gram

1 TF-IDF

- Singel-gram
- Range-Gram

1 T-SNE

REMOVE_STOPWORDS

```
[['specifications', 'sathiyas', 'cotton', 'bath', 'towel', 'bath',
'towel', 'red', 'yellow', 'blue', 'bath', 'towel', 'features', 'ma-
chine', 'washable', 'yes', 'material', 'cotton', 'design', 'self',
'design', 'general', 'brand', 'sathiyas', 'type', 'bath', 'towel',
'gsm', 'model', 'name', 'sathiyas', 'cotton', 'bath', 'towel', 'id-
eal', 'men', 'women', 'boys', 'girls', 'model', 'id', 'asvtwl', 'c-
olor', 'red', 'yellow', 'blue', 'size', 'medium', 'dimensions', 'l-
ength', 'inch', 'width', 'inch', 'box', 'number', 'contents', 'sal-
es', 'package', 'sales', 'package', 'bath', 'towel']]
```

LEMMATIZATION

```
[['specifications', 'sathiyas', 'cotton', 'bath', 'towel', 'bath',
'towel', 'red', 'yellow', 'blue', 'bath', 'towel', 'feature', 'mac-
hine', 'washable', 'yes', 'material', 'cotton', 'design', 'self',
'design', 'general', 'brand', 'sathiyas', 'type', 'bath', 'towel',
'gsm', 'model', 'name', 'sathiyas', 'cotton', 'bath', 'towel', 'id-
eal', 'men', 'women', 'boys', 'girls', 'model', 'id', 'asvtwl', 'c-
olor', 'red', 'yellow', 'blue', 'size', 'medium', 'dimension', 'le-
ngth', 'inch', 'width', 'inch', 'box', 'number', 'content', 'sales',
'package', 'sales', 'package', 'bath', 'towel']]
```

Ex: computer, computerization or computerize → compute

PRE-PROCESSING

place de marché

1 Preprocessing pipeline

- Tokenization
- StopWords
- Lemmatization
- Stemming

1 Bag of words

- Singel-gram
- Range-Gram

1 TF-IDF

- Singel-gram
- Range-Gram

1 T-SNE

LEMMATIZATION

```
[['specifications', 'sathiya', 'cotton', 'bath', 'towel', 'bath', 'towel', 'red', 'yellow', 'blue', 'bath', 'towel', 'feature', 'machine', 'washable', 'yes', 'material', 'cotton', 'design', 'self', 'design', 'general', 'brand', 'sathiya', 'type', 'bath', 'towel', 'gsm', 'model', 'name', 'sathiya', 'cotton', 'bath', 'towel', 'ideal', 'men', 'women', 'boys', 'girls', 'model', 'id', 'asvtwl', 'color', 'red', 'yellow', 'blue', 'size', 'medium', 'dimension', 'length', 'inch', 'width', 'inch', 'box', 'number', 'content', 'sales', 'package', 'sales', 'package', 'bath', 'towel']]
```

STEMMING

```
[['specif', 'sathiya', 'cotton', 'bath', 'towel', 'bath', 'towel', 'red', 'yellow', 'blue', 'bath', 'towel', 'featur', 'machin', 'was habl', 'ye', 'materi', 'cotton', 'design', 'self', 'design', 'gene r', 'brand', 'sathiya', 'type', 'bath', 'towel', 'gsm', 'model', 'name', 'sathiya', 'cotton', 'bath', 'towel', 'ideal', 'men', 'wome n', 'boy', 'girl', 'model', 'id', 'asvtwl', 'color', 'red', 'yellow', 'blue', 'size', 'medium', 'dimens', 'length', 'inch', 'width', 'inch', 'box', 'number', 'content', 'sale', 'packag', 'sale', 'pac kag', 'bath', 'towel']]
```

To convert a word to its meaningful base form by removing few characters

PRE-PROCESSING

place de marché

1 Preprocessing pipeline

- Tokenization
- StopWords
- Lemmatization
- Stemming

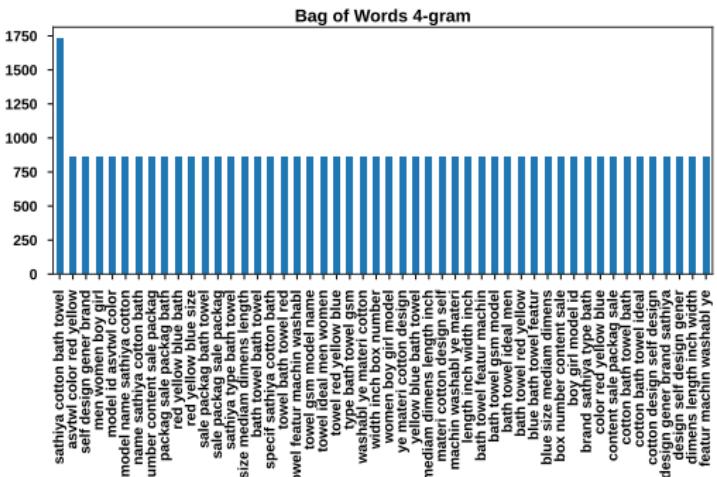
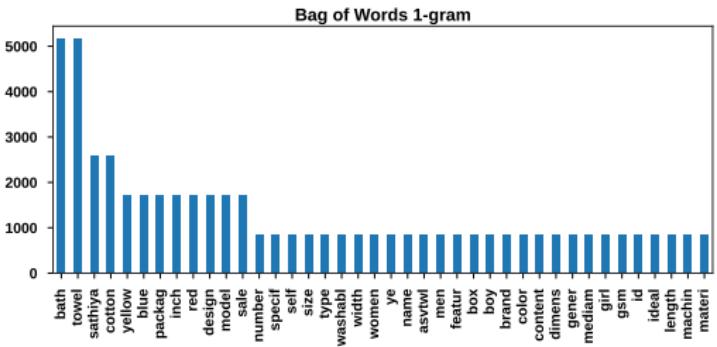
1 Bag of words

- Singel-gram
- Range-Gram

1 TF-IDF

- Singel-gram
- Range-Gram

1 T-SNE



PRE-PROCESSING

place de marché

1 Preprocessing pipeline

- Tokenization
- StopWords
- Lemmatization
- Stemming

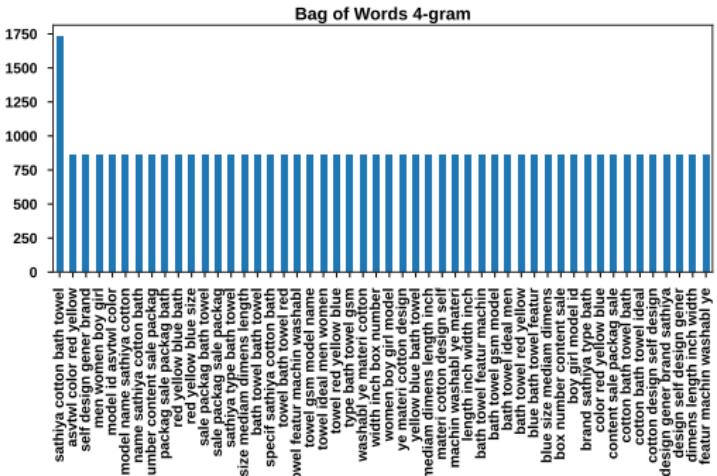
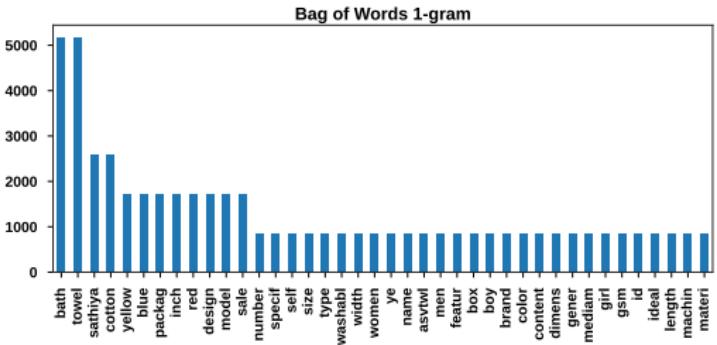
1 Bag of words

- Singel-gram
- Range-Gram

1 TF-IDF

- Singel-gram
- Range-Gram

1 T-SNE



PRE-PROCESSING

1 Preprocessing pipeline

- Tokenization
- StopWords
- Lemmatization
- Stemming

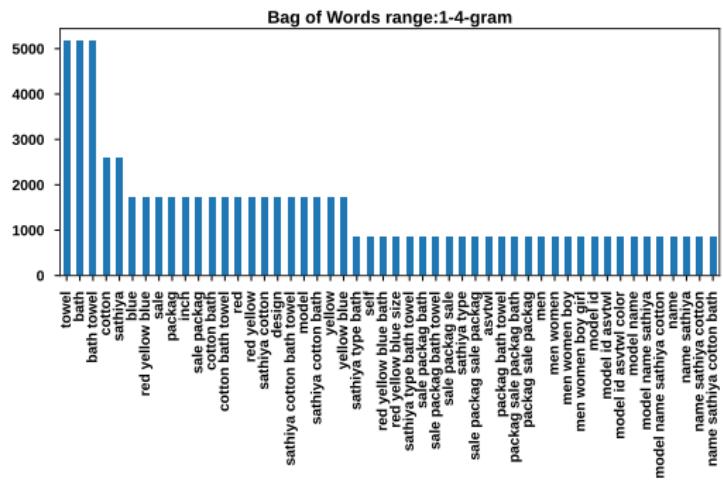
1 Bag of words

- Singel-gram
- Range-Gram

1 TF-IDF

- Singel-gram
- Range-Gram

1 T-SNE



PRE-PROCESSING

place de marché

1 Preprocessing pipeline

- Tokenization
- StopWords
- Lemmatization
- Stemming

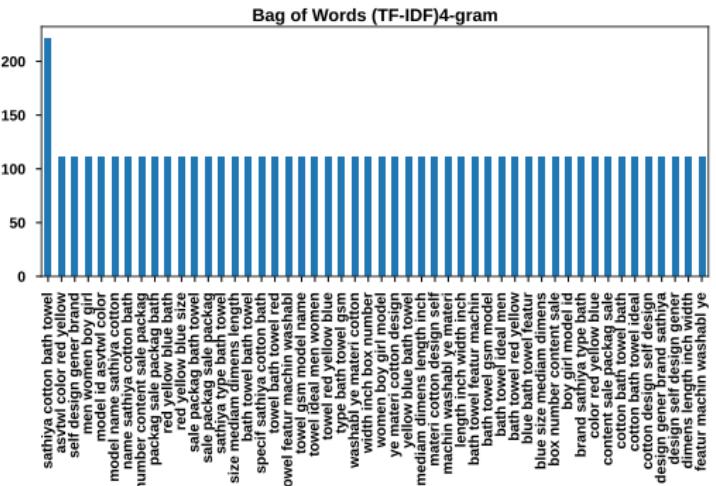
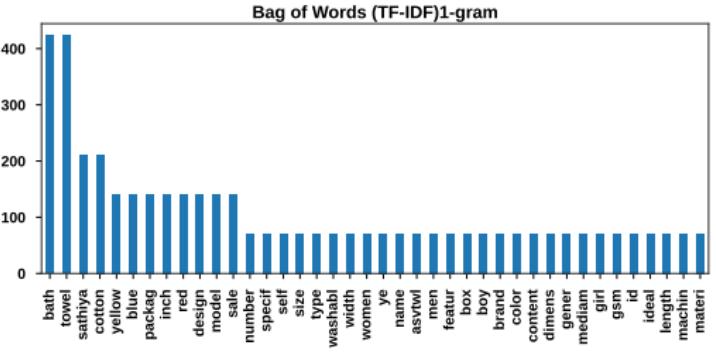
1 Bag of words

- Singel-gram
- Range-Gram

1 TF-IDF

- Singel-gram
- Range-Gram

1 T-SNE



PRE-PROCESSING

place de marché

1 Preprocessing pipeline

- Tokenization
- StopWords
- Lemmatization
- Stemming

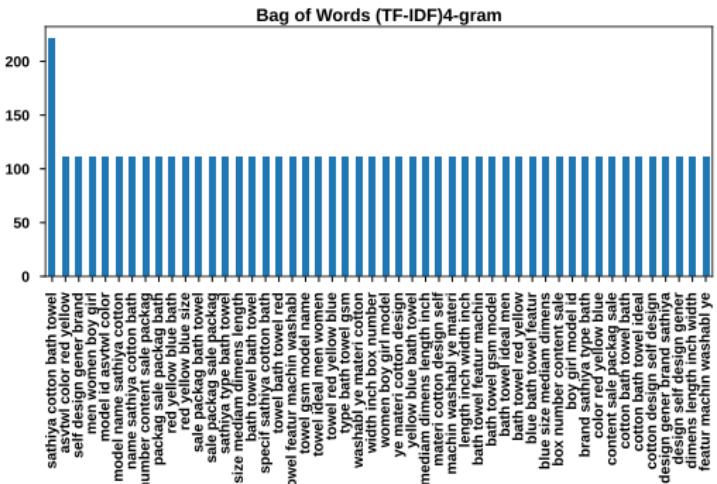
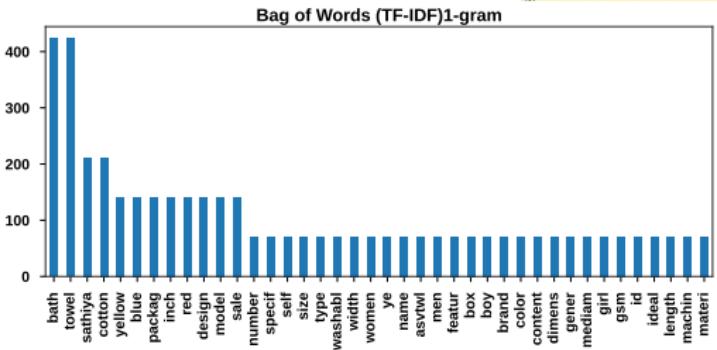
1 Bag of words

- Singel-gram
- Range-Gram

1 TF-IDF

- Singel-gram
- Range-Gram

1 T-SNE



PRE-PROCESSING

1 Preprocessing pipeline

- Tokenization
- StopWords
- Lemmatization
- Stemming

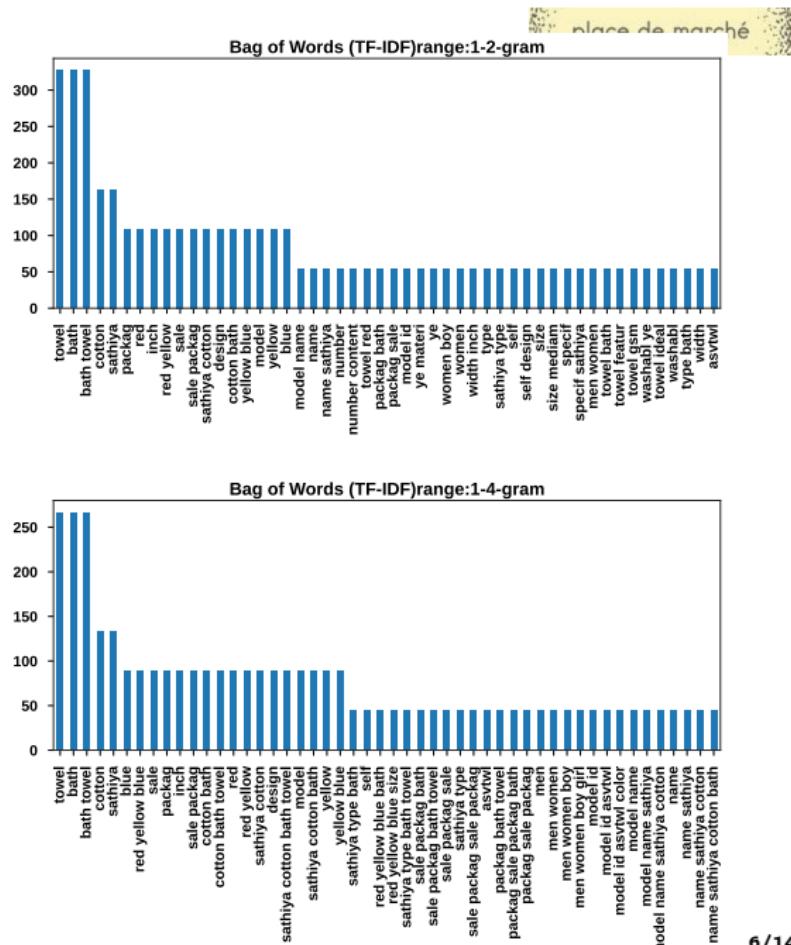
1 Bag of words

- Singel-gram
- Range-Gram

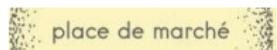
1 TF-IDF

- Singel-gram
- Range-Gram

1 T-SNE



PRE-PROCESSING



1 Preprocessing pipeline

- Tokenization
- StopWords
- Lemmatization
- Stemming

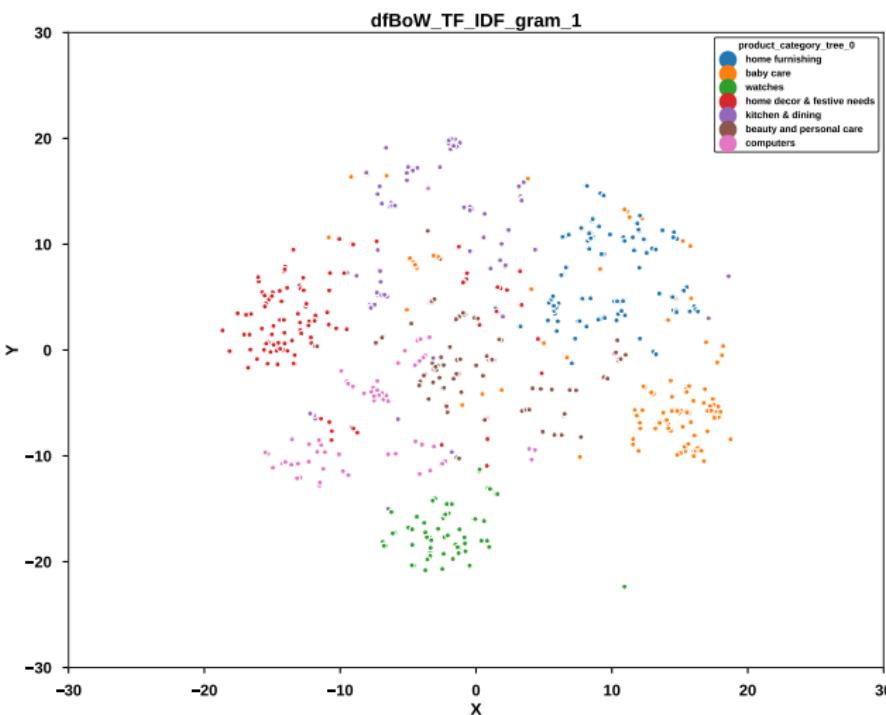
1 Bag of words

- Singel-gram
- Range-Gram

1 TF-IDF

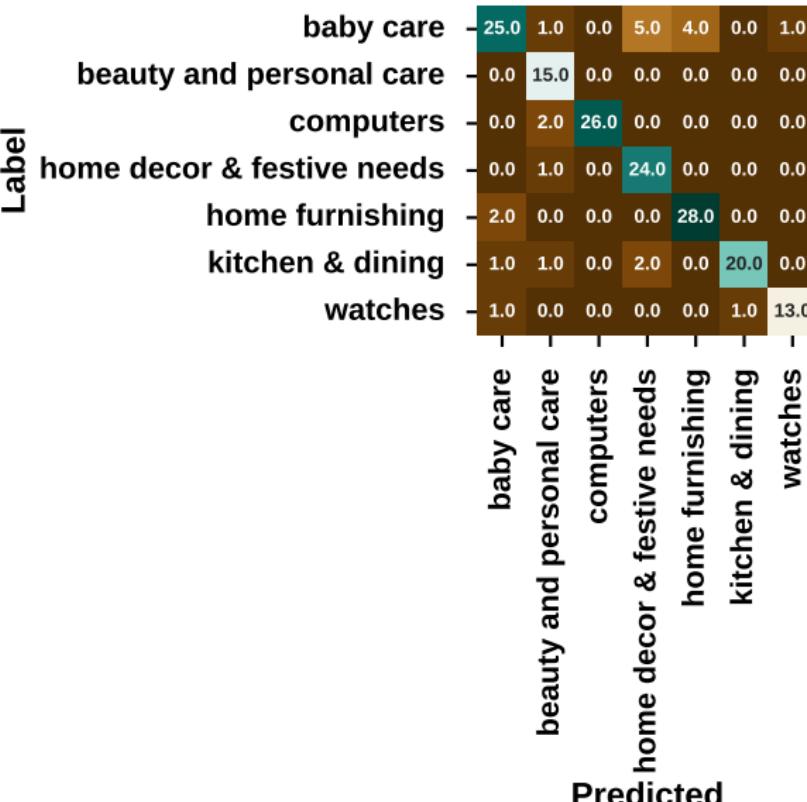
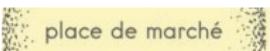
- Singel-gram
- Range-Gram

1 T-SNE



MODELING

Accuracy : 87%



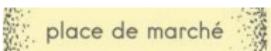
1 Naive Bayes

- Confusion matrix
- Weight
- Cloud of words

2 LDA & NMF*

MODELING

Accuracy : 87%



Label	Predicted						
	baby care	beauty and personal care	computers	home decor & festive needs	home furnishing	kitchen & dining	watches
baby care	25.0	1.0	0.0	5.0	4.0	0.0	1.0
beauty and personal care	0.0	15.0	0.0	0.0	0.0	0.0	0.0
computers	0.0	2.0	26.0	0.0	0.0	0.0	0.0
home decor & festive needs	0.0	1.0	0.0	24.0	0.0	0.0	0.0
home furnishing	2.0	0.0	0.0	0.0	28.0	0.0	0.0
kitchen & dining	1.0	1.0	0.0	2.0	0.0	20.0	0.0
watches	1.0	0.0	0.0	0.0	0.0	1.0	13.0

1 Naive Bayes

- Confusion matrix
- Weight
- Cloud of words

2 LDA & NMF*

MODELING

place de marché

1 Naive Bayes

- Confusion matrix
- Weight
- Cloud of words

2 LDA & NMF*

y=baby care top features	y=beauty and personal care top features	y=computers top features	y=home decor & festive needs top features	y=home furnishing top features	y=kitchen & dining top features	y=watches top features	
Weight?	Feature	Weight?	Feature	Weight?	Feature	Weight?	
+14.717	baby	+12.234	combo	+11.503	laptop	+17.023	
-7.798	cotton	+8.700	jewellery	-5.798	battery	-9.197	
-7.519	photo	+8.689	vanity	+5.440	light	+6.453	
-7.175	frame	+7.568	hair	+5.303	front	+6.468	
-7.080	small	+6.724	lipkart	+5.236	only	+5.753	
-6.934	offering	+6.530	com	+4.426	keyboard	+6.127	
-6.703	girls	+6.361	lowest	+4.226	admax	+5.956	
-6.142	width	+6.350	set	+4.130	replacement	+5.772	
-6.074	980	+5.925	lip	+4.017	warranty	+5.296	
-6.074	carter	+5.626	surgical	+4.017	painting	+5.260	
-5.969	girl	+5.570	eye	+4.017	bluetooth	+4.985	
-5.916	479	+5.553	traits	-3.958	pad	+4.603	
-5.680	size	+5.372	kit	-3.936	flexible	+4.573	
-5.635	wallmantra	+5.345	oil	-3.849	xyzel	+4.390	
-5.561	ireeya	+5.149	and	-3.814	led	+4.178	
-5.425	ideal	+5.141	blush	-1000 more positive ...	brass	... 753 more positive ...	
-5.316	brush	+4.963	massage	-4026 more negative ...	paper	... 4273 more negative ...	
-5.268	coral	+4.939	body	-3.834	combo	-4.125	
-5.244	water	+4.939	flow	-3.834	with	-4.125	
-5.186	529	... 1171 more positive ...	-3.927	... 1341 more positive ...	-5.209	... 479	
... 950 more positive 3855 more negative ...	-3.927	... 3000 more negative ...	-5.312	baby	... 2493	
... 4076 more negative -5.911	at	-5.390	... 3000 more negative ...	-4.249	juicer	... 2241
			-4.407	... 3000 more negative ...	-4.253	lock	... 2334
			-5.468	vanity	-4.28	hands	... 2100 more positive ...
					-4.611	showpiece	... 4615 more negative ...
					-2.489	cm	... 2489 cm
					-2.503	win	... 2503 win

Weights

MODELING

place de marché

1 Naive Bayes

- Confusion matrix
- Weight
- Cloud of words

2 LDA & NMF*



Cloud of words

MODELING

place de marché

1 Naive Bayes

- Confusion matrix
- Weight
- Cloud of words

2 LDA & NMF*

Topic 0:

com flipkart product ship day

Topic 1:

com flipkart ship product replac

Topic 2:

flipkart com ship product rs

Topic 3:

com flipkart ship product genuin

Topic 4:

flipkart com ship product deliveri

Topic 0:

com flipkart product

Topic 1:

box materi set watch

Topic 2:

box materi set watch

Topic 3:

box materi set watch

Topic 4:

box materi set watch

LDA

NMF

INDEX

place de marché

- 1 Mission objective
- 2 Data preparation
 - Data cleaning
- 3 NLP (Text)
 - Pre-processing
 - Modeling

- 4 Machine Learning (CV)
 - Pre-processing
 - Feature extraction
- 5 Deep Learning (CV)
 - VGG16
- 6 Conclusions

PRE-PROCESSING

place de marché

- 1 GRAY
- 2 Scaling
- 3 Gaussian filtering
- 4 CLAHE



PRE-PROCESSING

place de marché

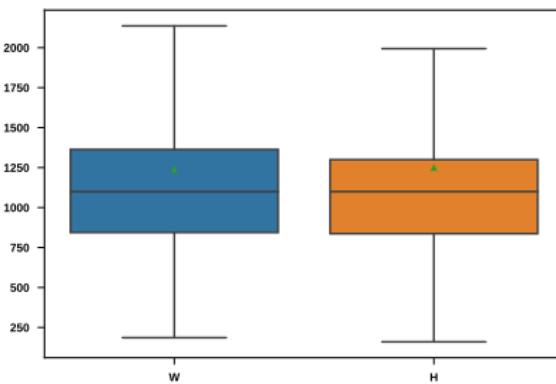


1 GRAY

2 Scaling

3 Gaussian filtering

4 CLAHE



PRE-PROCESSING

place de marché

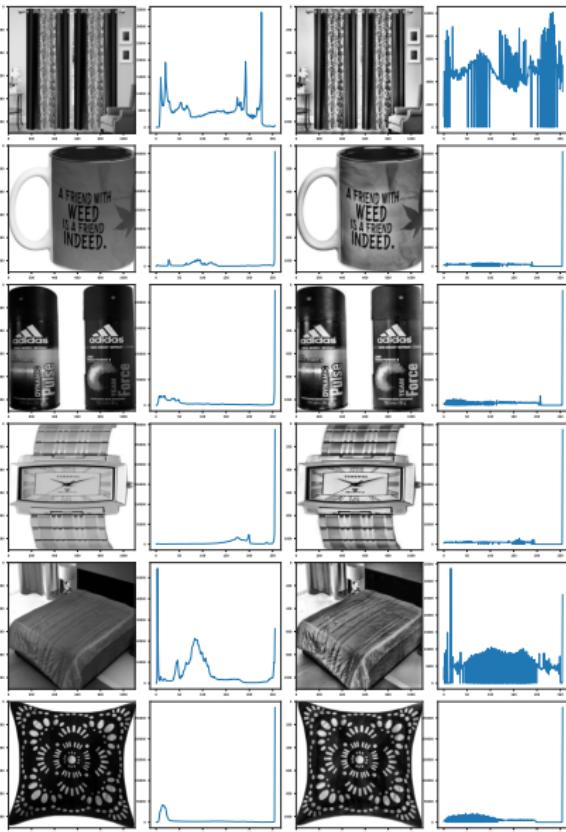
- 1 GRAY
- 2 Scaling
- 3 Gaussian filtering
- 4 CLAHE



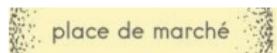
PRE-PROCESSING

place de marché

- 1 GRAY
- 2 Scaling
- 3 Gaussian filtering
- 4 CLAHE



FEATURE EXTRACTION



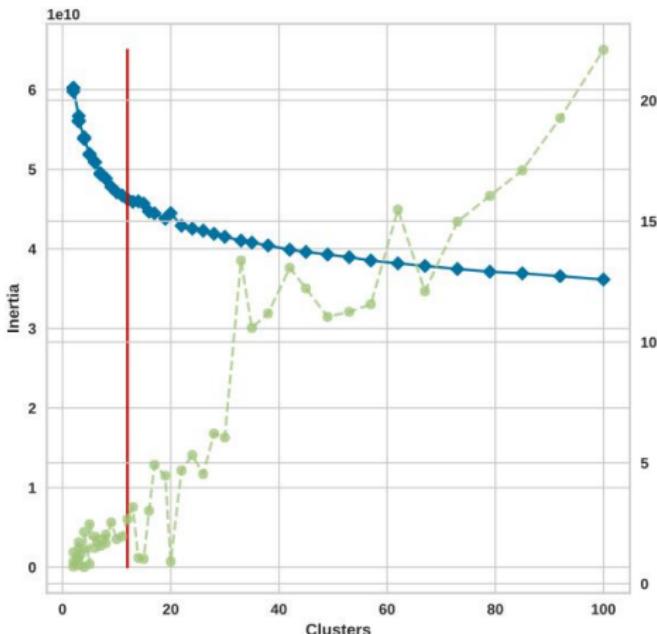
ORB (386913 Features)

1 KMeans

2 BoVW

■ T-SNE

3 Naive Bayes



MiniBatchKMeans (k=12)

FEATURE EXTRACTION

place de marché

1 KMeans

2 BoVW

■ T-SNE

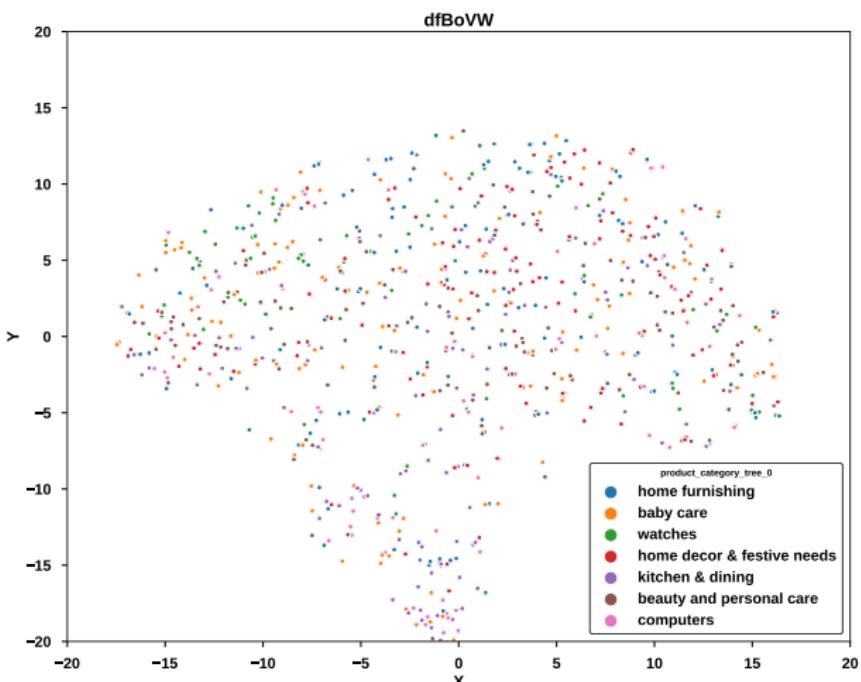
3 Naive Bayes

	0.0	1.0	2.0	3.0	4.0	5.0	6.0	7.0	8.0	9.0	10.0	11.0
0.0	66	33	37	19	27	24	41	50	59	14	68	62
1.0	56	24	20	25	23	72	38	74	65	25	38	40
2.0	18	40	33	19	51	32	60	52	48	35	48	64
3.0	37	50	15	105	15	79	17	45	78	19	21	19
4.0	33	26	34	31	31	55	45	73	48	21	66	37
...
858.0	11	38	16	106	45	99	34	35	22	26	46	22
859.0	51	27	16	89	29	73	44	57	45	27	28	14
860.0	36	35	28	91	63	61	28	31	23	32	28	22
861.0	27	48	35	105	29	45	51	38	15	26	46	35
862.0	32	57	18	67	40	53	38	41	31	37	28	17

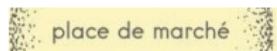
FEATURE EXTRACTION

place de marché

- 1 KMeans
- 2 BoVW
- T-SNE
- 3 Naive Bayes



FEATURE EXTRACTION



Accuracy : 27%

Label	Predicted							
	baby care	beauty and personal care	computers	home decor & festive needs	home furnishing	kitchen & dining	watches	
baby care	19.0	2.0	22.0	29.0	7.0	10.0	15.0	
beauty and personal care	20.0	9.0	15.0	34.0	10.0	13.0	9.0	
computers	6.0	3.0	30.0	21.0	6.0	8.0	5.0	
home decor & festive needs	4.0	3.0	10.0	52.0	25.0	15.0	3.0	
home furnishing	12.0	3.0	14.0	33.0	27.0	8.0	8.0	
kitchen & dining	14.0	1.0	27.0	33.0	7.0	21.0	6.0	
watches	12.0	1.0	5.0	23.0	9.0	4.0	17.0	
	baby care	beauty and personal care	computers	home decor & festive needs	home furnishing	kitchen & dining	watches	

- 1 KMeans
- 2 BoVW
 - T-SNE
- 3 Naive Bayes

INDEX

place de marché

- 1 Mission objective
- 2 Data preparation
 - Data cleaning
- 3 NLP (Text)
 - Pre-processing
 - Modeling
- 4 Machine Learning (CV)
 - Pre-processing
 - Feature extraction
- 5 Deep Learning (CV)
 - VGG16
- 6 Conclusions

VGG16

place de marché

- Freeze the feature extraction part

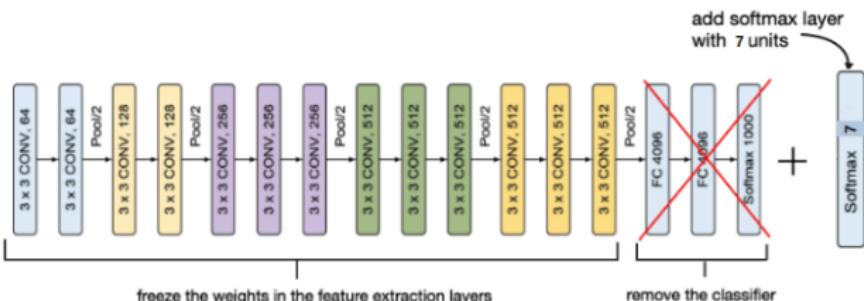
- Remove the classifier part

- Add our new classifier softmax layer 7 hidden units

1 Model

2 Results

- Validation
- Confusion matrix



Total params: 14,890,311

Trainable params: 175,623

Non-trainable params: 14,714,688

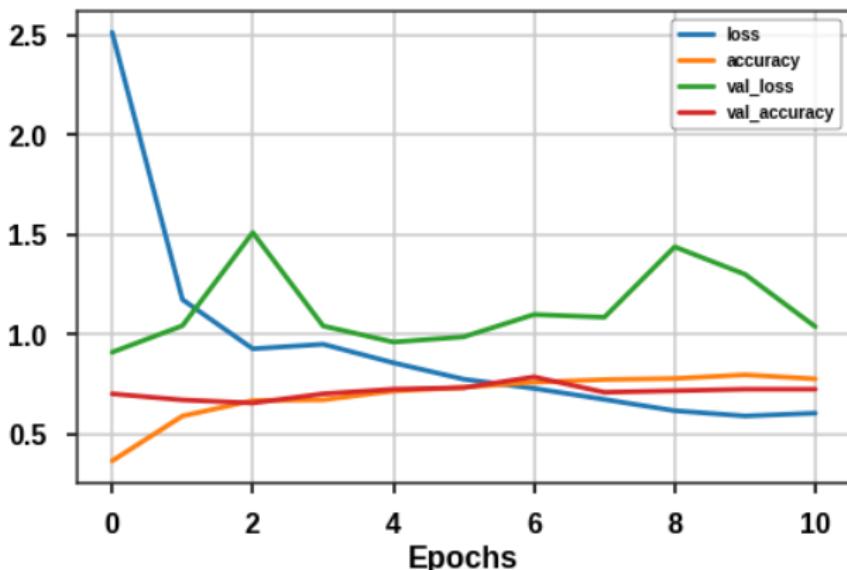
VG**G**16

place de marché

1 Model

2 Results

- Validation
- Confusion matrix



Validation accuracy : 78%

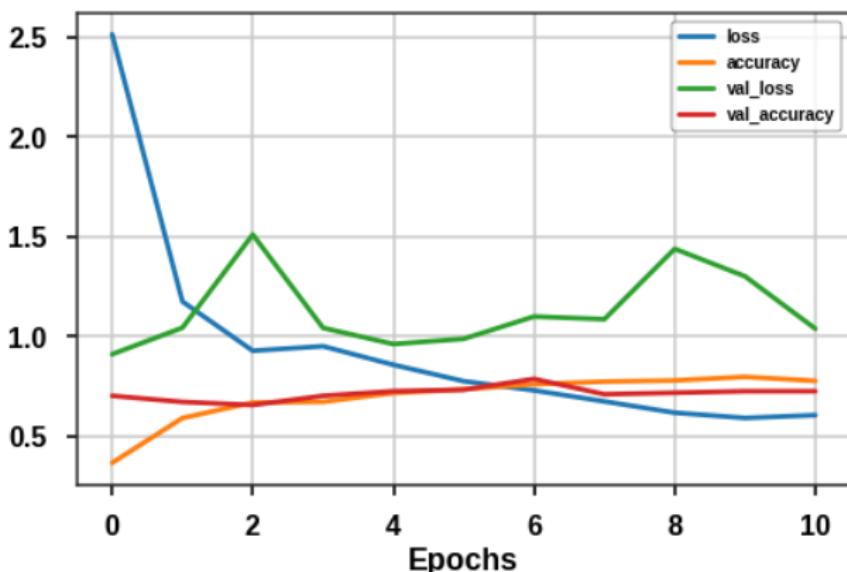
VGG16

place de marché

1 Model

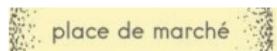
2 Results

- Validation
- Confusion matrix

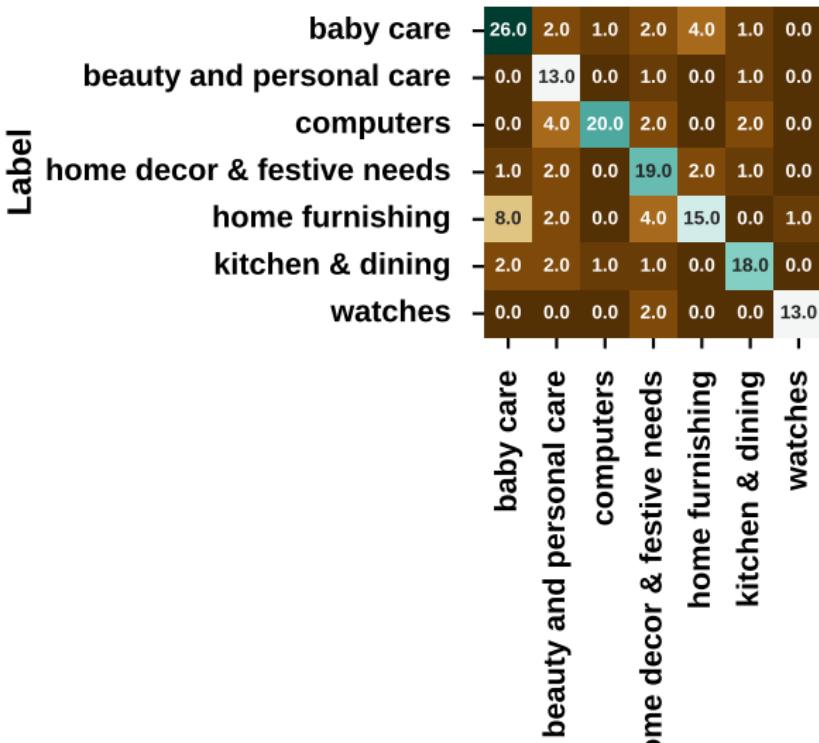


Validation accuracy : 78%

VGG16



Accuracy : 72%



- 1 Model
- 2 Results

- Validation
- Confusion matrix

INDEX

 place de marché 

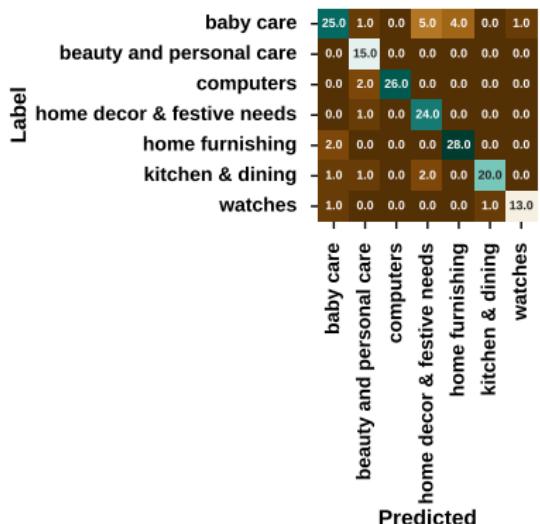
- 1 Mission objective
- 2 Data preparation
 - Data cleaning
- 3 NLP (Text)
 - Pre-processing
 - Modeling
- 4 Machine Learning (CV)
 - Pre-processing
 - Feature extraction
- 5 Deep Learning (CV)
 - VGG16
- 6 Conclusions

CONCLUSIONS

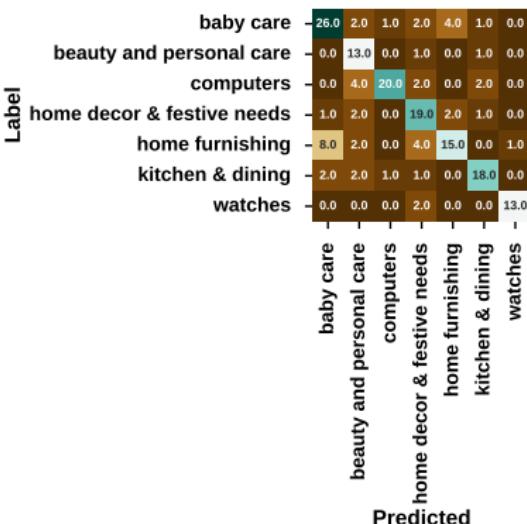
place de marché

- Dans l'ordre croissant, les meilleurs résultats sont les suivant :

Accuracy : 87%



Accuracy : 72%



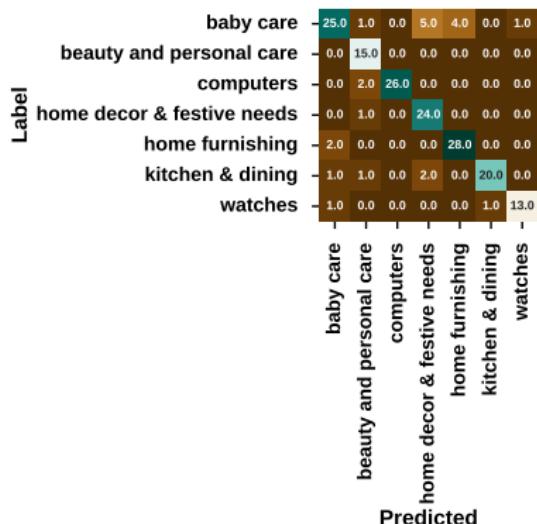
CONCLUSIONS

place de marché

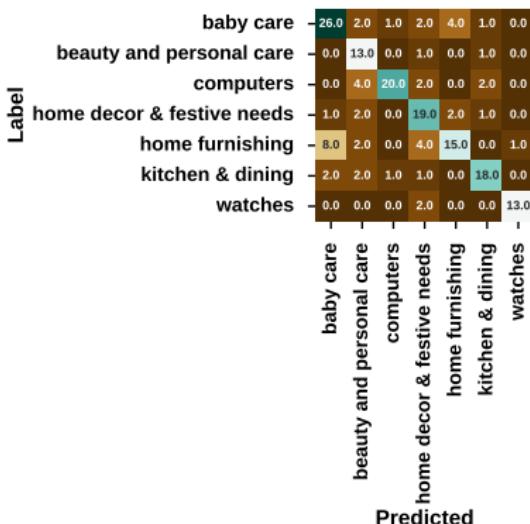
1 Dans l'ordre croissant, les meilleurs résultats sont les suivant :

- NLP (87%), Deep Learning (VGG16, 72%), Machine learning (OpenCV, 27%)

Accuracy : 87%



Accuracy : 72%



NLP

DL : VGG16