# Classifying consumer goods automatically

MAY 2022

NAJWA MOULINE

# Presentation Outline

1. Objectives

2. Dataset Preparation

3. Pre-processing and Clustering

4. Classification Feasibility

# Objectives

## Context

- Company interested to launch an e-commerce marketplace
- Sellers offer items to buyers by posting a photo and description
- Assignment of item categories manually by sellers – tedious and unreliable

## Business Problem

- Improving user experience for sellers and buyers through reliable item classification
- Automating the assignment of the category of products for sale

## Mission

- Perform a feasibility study of an image or description-based item classification engine to automate:
  - Pre-processing of the dataset (textual and visual)
  - Dimension reduction
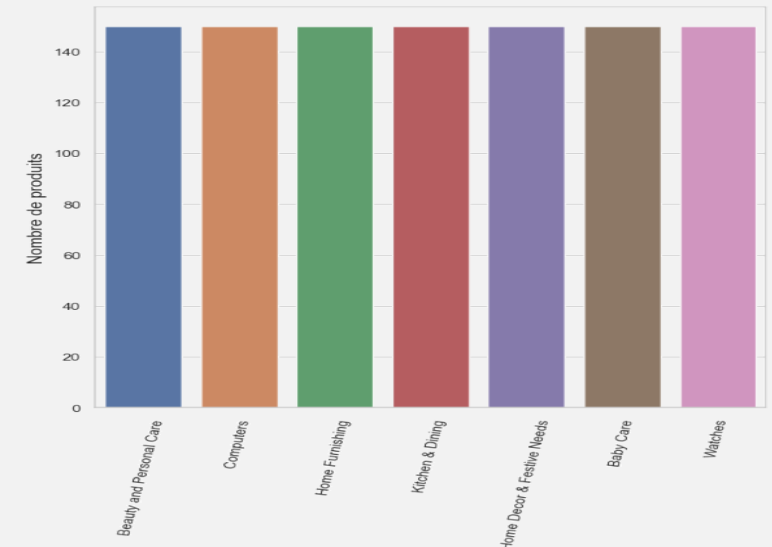  - Clustering and representation in 2D

# Dataset Preparation

## Dataset

| Dataset | A csv format file and jpg images |
|---------|----------------------------------|
| Observations | 1050 articles |
| Variables | 15 relating to the identifier of the item, description, price, rating, image, … |
| Duplicates | None |
| Filling | Completed for the data considered |
| Features | Description (text), image (visual) |
| **Target** | |

- Category class 1 – 7 categories

- Category class 2 – 62 categories

- The dataset is good for our mission

- It is well distributed by category 1, which should make it possible to identify all categories with modeling



4

# Pre-processing and clustering

Methodology

(textual data)

## Methodology – classification of textual data

### Pre-processing

| Standardisation and Simplification | → | Features Extraction |

2 options:
- NLTK
- Spacy

4 options:
- BoW (CV / TF-IDF)
- Word2Vec
- BERT / USE

### Modeling

| Reduction | → | Modeling and Evaluation |

PCA (80%)
TSNE (2D)

2 options:
- Clustering
- Supervised classification

# Pre-processing and clustering

Pre-processing

(textual data)

## Standardisation and Simplification (NLTK) – from 9591 unique words to 2796

| Steps | Resulting Worldcloud |
|---|---|
| 1. Switching to Lowercase<br>2. Removal of punctuation<br>3. Tokenisation<br>4. Removal of StopWords<br>5. Lemmatisation<br>6. Removal of words with less than 3 letters<br>7. Removal of words occurring once in the corpus |  |

### Example

Phrase de base : 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 interiors.This curtain is made from 100% high quality polyester fabric.It features an eyelet style stitch with Metal Ring.It makes the room envi

Phrase prétraitée : feature elegance polyester multicolor abstract eyelet door curtain floral curtainelegance polyester multicolor abstract eyelet door curtain height pack price curtain enhance look interiorsthis curtain make high quality polyester fabricit feature eyelet style stitch metal ri

# Pre-processing and clustering

## Pre-processing

## (textual data)

| Steps | Resulting Worldcloud |
|---|---|

1. Switching to Lowercase
2. Tokenisation
3. Lemmatisation
4. Keeping words that are pronouns, adj, noun, and verb
5. Removal of StopWords
6. Removal of words with less than 3 letters
7. Removal of words occurring once in the corpus



## Example

```
Phrase de base : 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) P
rice: Rs. 899 This curtain enhances the look of the interiors.This curtain is made from 100% high
quality polyester fabric.It features an eyelet style stitch with Metal Ring.It makes the room envi
```

```
Phrase prétraitée : [ feature', 'elegance', 'polyester', 'multicolor', 'abstract', 'eyelet', 'doo
r', 'curtain', 'floral', 'curtain', 'elegance', 'polyester', 'multicolor', 'abstract', 'eyelet',
'door', 'curtain', 'height', 'pack', 'price', 'curtain', 'enhance', 'look', 'curtain', 'high', 'qu
ality', 'polyester', 'feature', 'eyelet', 'style', 'stitch', 'metal', 'ring.it', 'room', 'environm
```

# Pre-processing and clustering

## Pre-processing

(textual data)

## Features Extraction – 5 options considered

### Bag of Words - CV

Vector formed from the occurrence of the word in the description.
Vector (1050, number of unique tokens)

### Bag of Words – TF-IDF

Vector formed from the occurrence of the word in the description and in the corpus.
Vector (1050, number of unique tokens)

### Word2Vec

Vector formed from the "meaning" of the word.
Vector (1050, 300)

### BERT

Vector formed from the position and context of the word in the description (transfer learning). Vector (1050, 768)

### USE

Vector formed from the position and context of the word in the description (transfer learning). Vector (1050, 512)
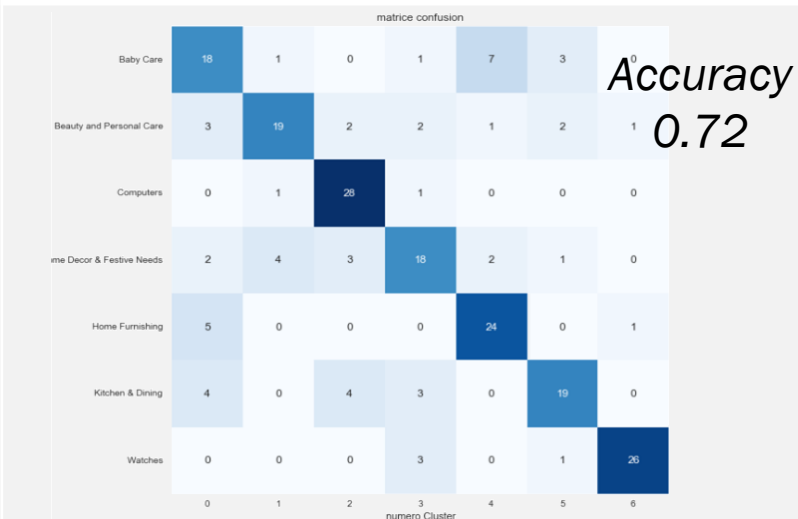
# Pre-processing and clustering

Modeling

(textual data)



Unsupervised classification with KMeans

**Visualisation NLTK / TFIDF** — ARI 0.39

**Visualisation Spacy / Word2Vec** — ARI 0.29

**Visualisation Spacy / BERT** — ARI 0.28

**Visualisation NLTK / USE** — ARI 0.37

# Pre-processing and clustering

Modeling

(textual data)

## Supervised classification with KNN

### Visualisation NLTK / TFIDF



*Accuracy 0.72*

### Visualisation Spacy / Word2Vec



*Accuracy 0.82*

### Visualisation Spacy / BERT



*Accuracy 0.77*

### Visualisation NLTK / USE



*Accuracy 0.90*

# Pre-processing and clustering

Methodology

(visual data)

**Pre-processing**

Standardisation and Simplification

OpenCV

Features Extraction

2 options:
- ORB
- CNN (Keras)

**Modeling**

Reduction

PCA (80%)
TSNE (2D)

Modeling and Evaluation

2 options:
- Clustering
- Supervised classification

# Pre-processing and clustering

Pre-processing

(visual data)

## Standardisation and Simplification (OpenCV)

### Steps

1. Image redimensioning
2. Transformation to black and white
3. Noise cancellation
4. Frequencies equalisation

# Pre-processing and clustering

## Pre-processing

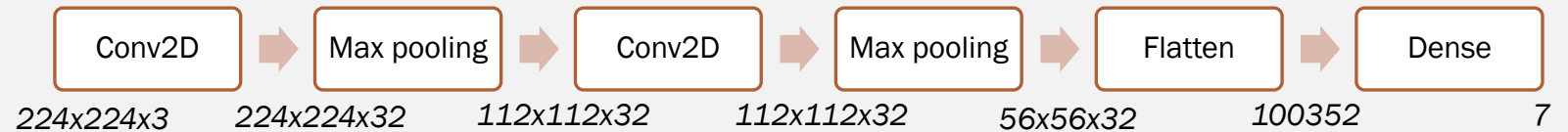## (visual data)

## Features Extraction

### ORB

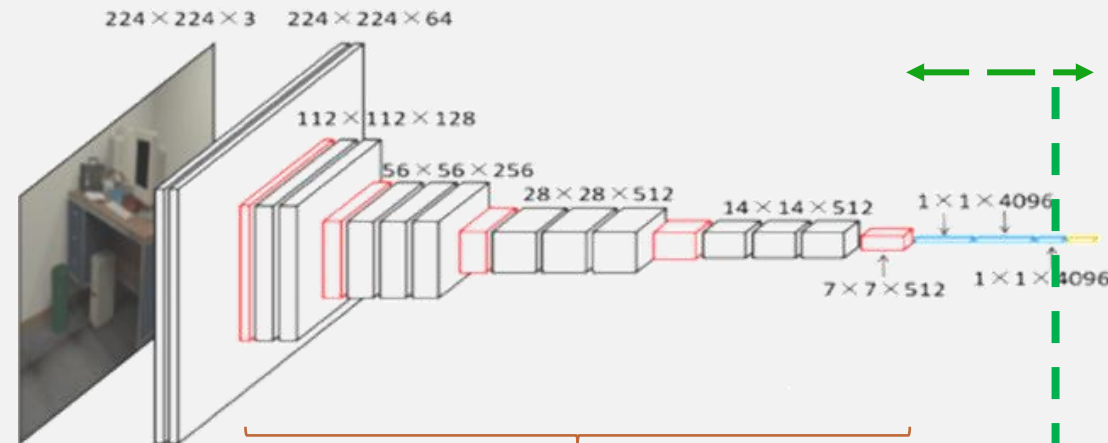| Detection of descriptors | → | Segmentation on 1000 image clusters | → | Feature = Histogram of the image |
|---|---|---|---|---|

### CCN Simple

| Conv2D | → | Max pooling | → | Conv2D | → | Max pooling | → | Flatten | → | Dense |
|---|---|---|---|---|---|---|---|---|---|---|

*224x224x3*  *224x224x32*  *112x112x32*  *112x112x32*  *56x56x32*  *100352*  *7*

### CCN transfer learning - VGG16 model



Tested scenarios:
- Deletion of one or more "high" layers
- Adding a classification of 7 layers
- Layers training

224 × 224 × 3   224 × 224 × 64

112 × 112 × 128

56 × 56 × 256

28 × 28 × 512

14 × 14 × 512

1 × 1 × 4096

7 × 7 × 512   1 × 1 × 4096

Convolution + ReLU
max pooling
fully connected + ReLU
softmax

Convolutional base: pre-trained on ImageNet

# Pre-processing and clustering
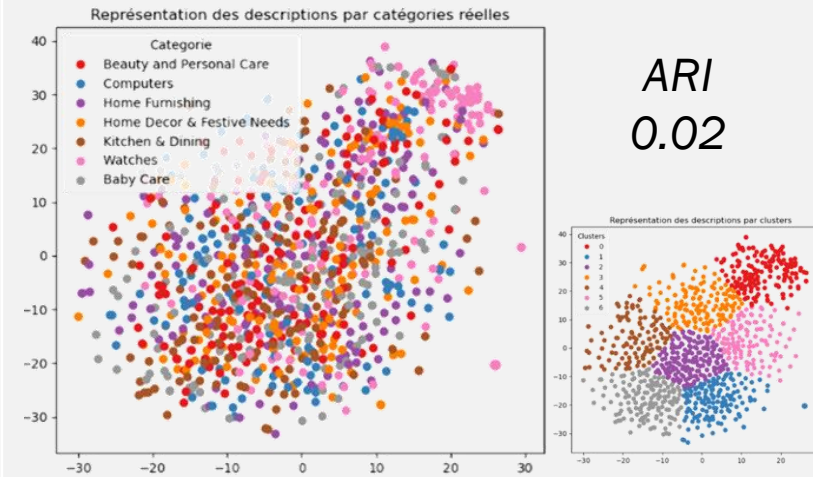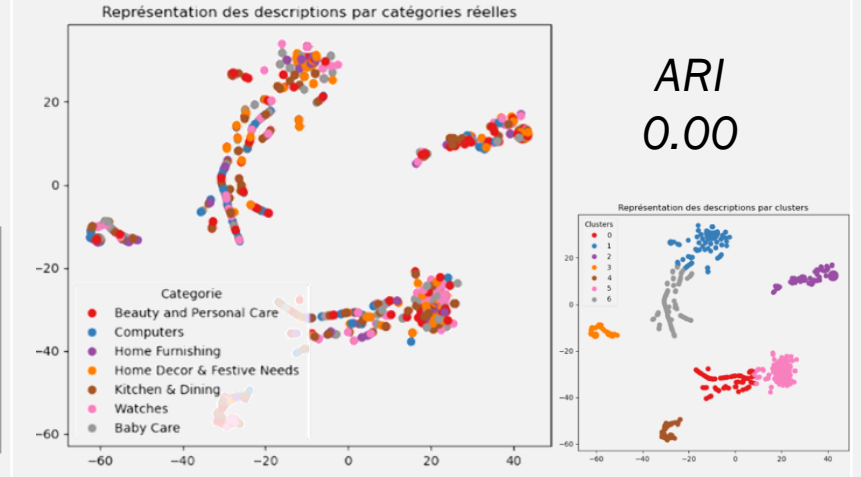
Modeling

(visual data)

# Pre-processing and clustering

Modeling

(visual data)

## Supervised classification with KNN

**Visualisation *ORB***

*Accuracy 0.23*

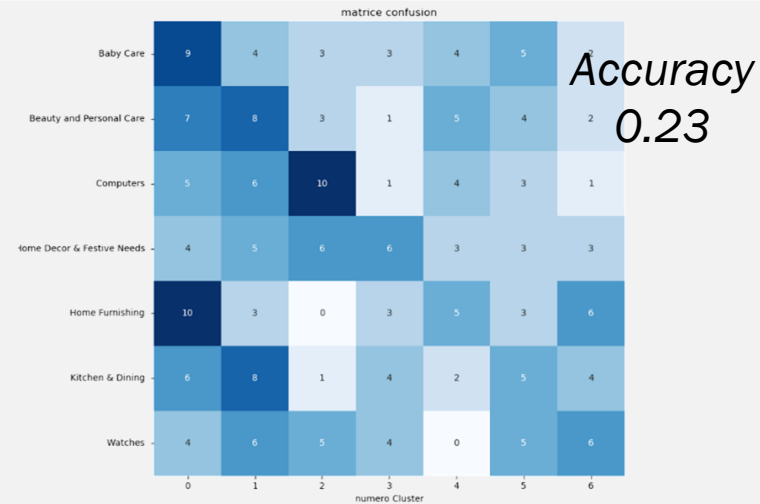**Visualisation CNN Simple**

*Accuracy 0.22*

**Visualisation CNN VGG16**

*Accuracy 0.7*

**Visualisation CNN VGG16 (without Dense)**

*Accuracy 0.8*

# Classification Feasibility

Feasibility and Conclusions

## Feasibility study – unsupervised classification

- Clustering performed by the Kmeans model with 7 clusters
- In some cases, clusters spaced and evenly distributed
- But low ARIs

|  | ARI |  | ARI |
|---|---|---|---|
| Description - NLTK / CV | 0.33 | Image - ORB | 0.02 |
| Description - NLTK / TFIDF | 0.39 | Image – CNN Simple | 0.01 |
| Description - NLTK / Word2Vec | 0.28 | Image – CNN VGG16 | 0.33 |
| Description - NLTK / BERT | 0.32 | Image – CNN VGG16 (Dense 7) | 0.03 |
| Description - NLTK / USE | 0.37 | Image – CNN VGG16 (without Dense) | 0.46 |

=> Inconclusive feasibility study

# Classification Feasibility

Feasibility and Conclusions

## Feasibility study – supervised classification

- Prediction of the classification issued from the KNN model trained on the 7 categories
- For the majority of cases, clusters spaced and distributed equally
- With right accuracy for all categories except BabyCare

|  | Acc |  | Acc |
|---|---|---|---|
| Description - NLTK / CV | 0.74 | Image - ORB | 0.23 |
| Description - Spacy / TFIDF | 0.75 | Image – CNN Simple | 0.20 |
| Description - Spacy / Word2Vec | 0.82 | Image – CNN VGG16 | 0.70 |
| Description - NLTK / BERT | 0.84 | Image – CNN VGG16 (Dense with 7) | 0.27 |
| Description - NLTK / USE | 0.90 | Image – CNN VGG16 (without Dense) | 0.79 |

=> Conclusive feasibility study

# Classification Feasibility

## Recommendations

---

**Classification engine recommendations**

## Improved performance

### Data

- Consideration of a larger database (150 articles per category is too little)
- Review of the assignment of categories to certain products initially miscategorised
- Subdivision of a category – example of Baby Care into 3 categories like Baby Furniture, Baby Care and Baby Clothes.

### Pre-processing

- n-grams for descriptions
- Other methods specific to e-commerce

### Features Extraction

- Testing other features extraction methods
- Aggregation of image/text features

### Modeling

- Modeling optimisation – supervised and unsupervised algo, and associated hyperparameters

*Thank you for your attention!*