

# Classifiez automatiquement des biens de consommation



Zaccaria Amillou

# Sommaire

- Contexte
- Données
- EDA
- Traitement de texte
- Traitement image
- Modélisation
- Script

# Contexte



## Objectifs

Automatiser la classification des images des produits

- Images
- Description

Labellisation automatique des objets via une image et une description.



Key Features of Elegance  
Polyester Multicolor  
Abstract Eyelet Door  
Curtain Floral...

Home Furnishing



Specifications of Sathi-  
yas Cotton Bath Towel  
(3 Bath Towel, Red,  
Yellow, Blue)...

Baby Care



# Données

```
In 13 1 # affichage sample df
2 display(df.sample(5))
Executed at 2024.03.12 09:22:03 in 8ms
```

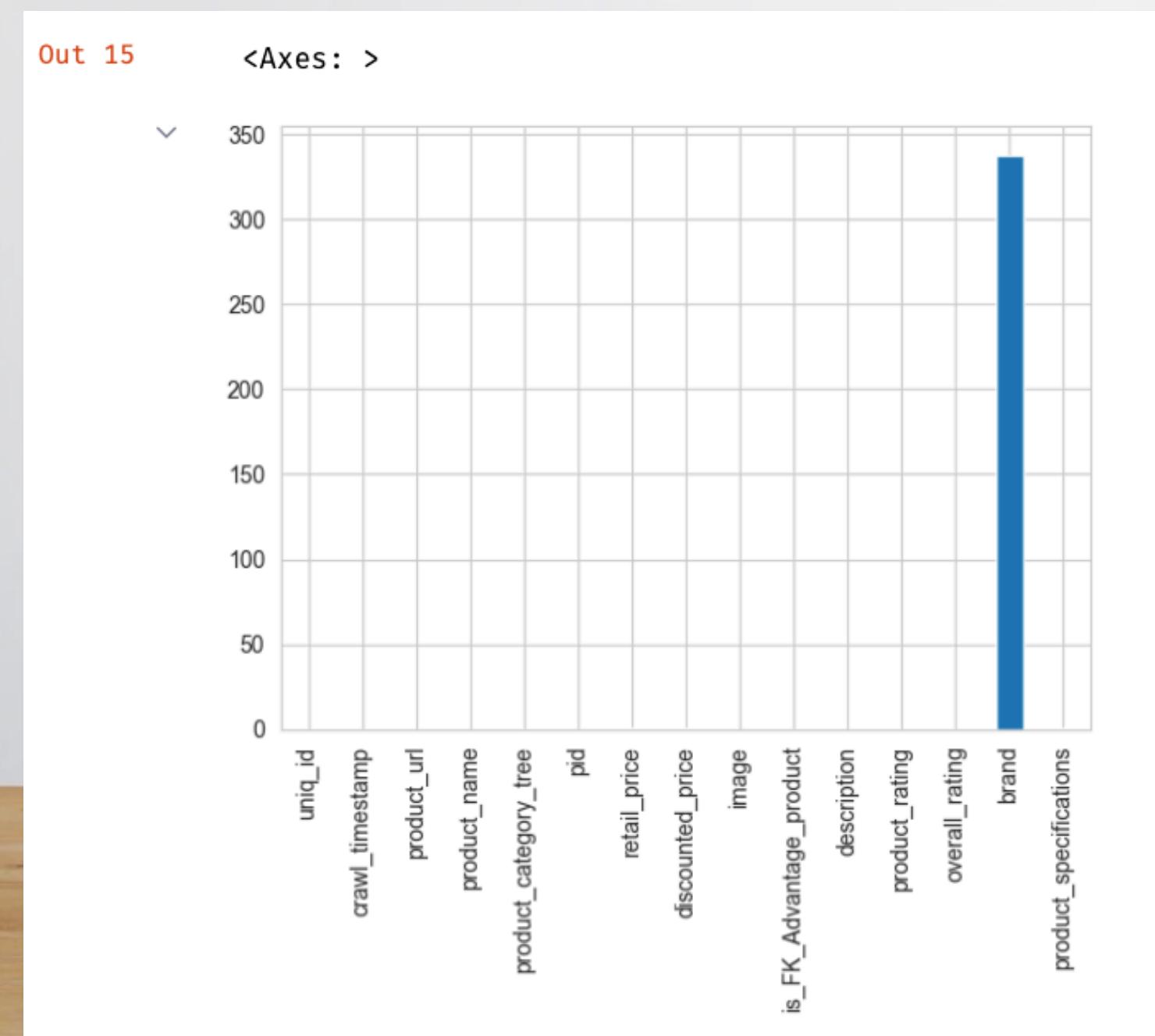
	uniq_id	crawl_timestamp	product_url	product_name	prod
71	f2658bad8c2b7d5b08984c6ac34267f	2016-04-19 05:00:32 +0000	http://www.flipkart.com/prime-printed-8-seater..	Prime Printed 8 Seater Table Cover	["Ho
430	a54b94096938252901d3f7f8de97bfff	2016-03-11 06:55:07 +0000	http://www.flipkart.com/ocean-vglass0051/p/itm..	ocean VGLASS0051	["Ki
306	7e11aaffb1d08f8091ea94598f7865a	2015-12-01 06:13:00 +0000	http://www.flipkart.com/lal-haveli-rajasthani-..	Lal Haveli Rajasthani Handmade Decorative Wood ..	["Ho
809	1d1be744e491ed61e705c20e4a72320	2015-12-12 11:46:53 +0000	http://www.flipkart.com/oxyglow-lacto-bleach-f..	Oxyglow Lacto Bleach & Fruit Massage Cream Wit ..	["Be
225	67d6b4a8aa19d57740b0440365b7762	2015-12-01 12:40:44 +0000	http://www.flipkart.com/lal-haveli-decorative-..	Lal Haveli Decorative Dholak Musician Showpiec ..	["Ho

le df est composée par 1050 lignes et 15 colonnes

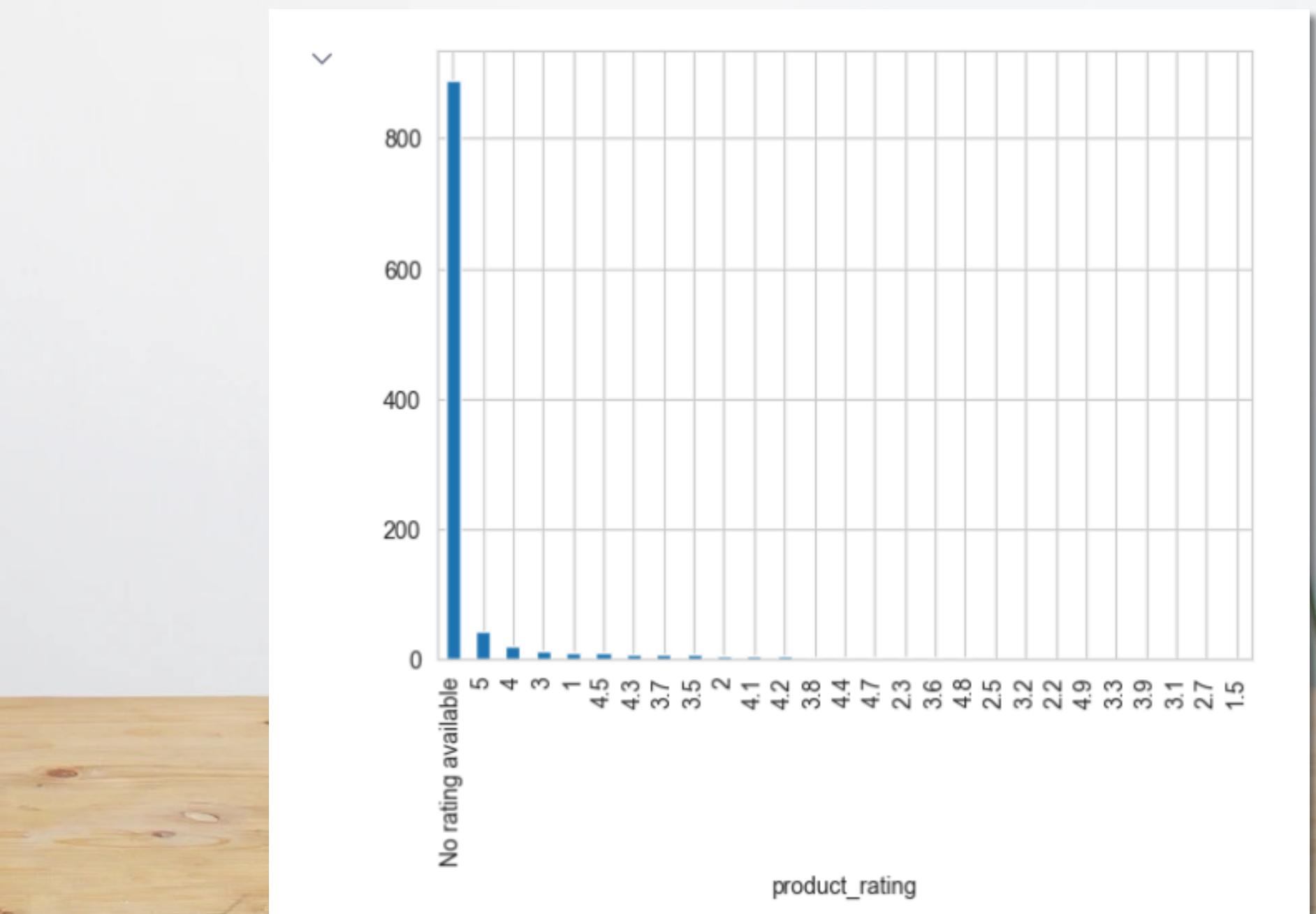
# EDA

## Exploration des données

Valeurs nulles

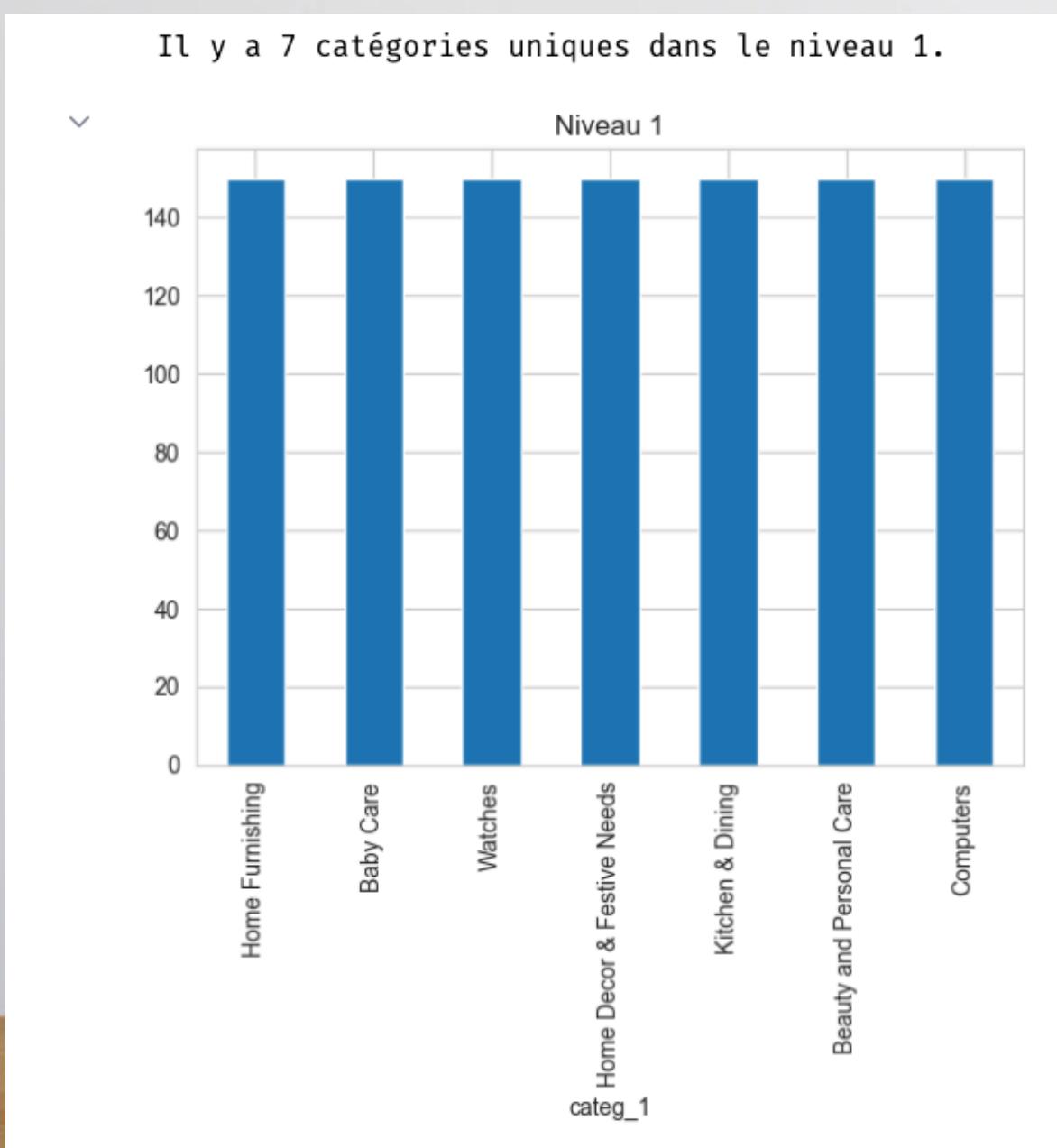


Contrôle notes clients

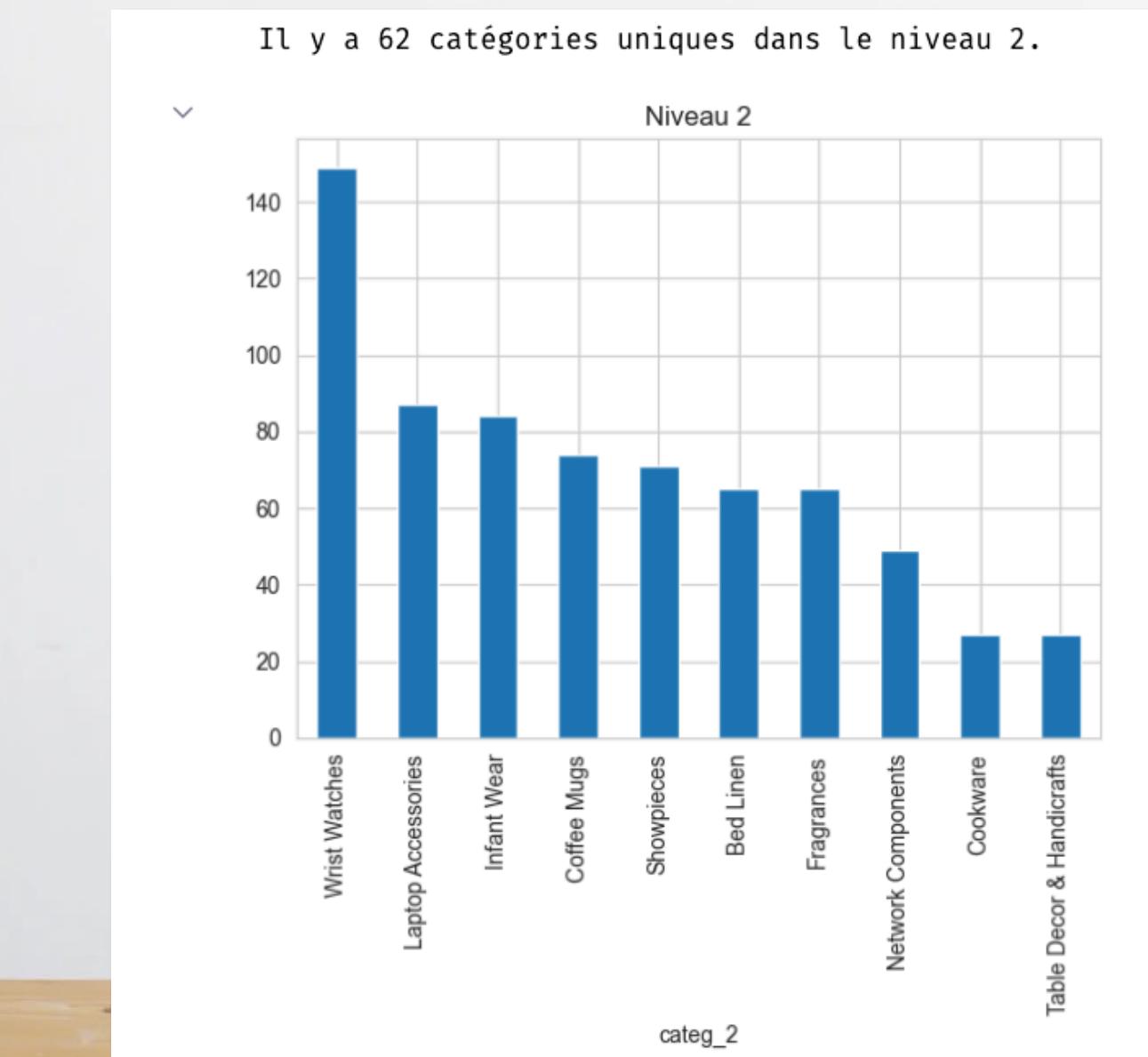


# EDA

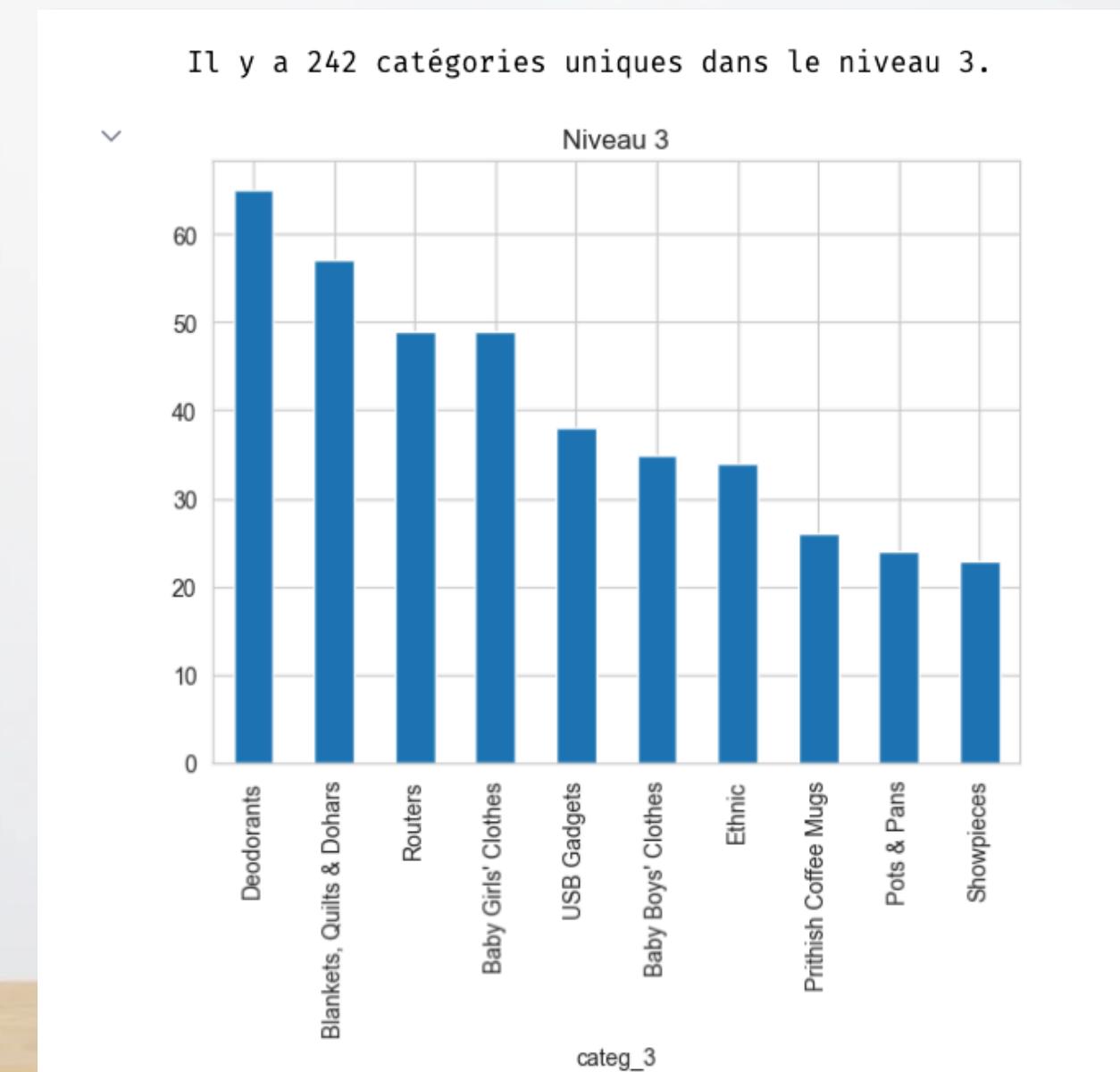
## Exploration des données



Niveaux 1

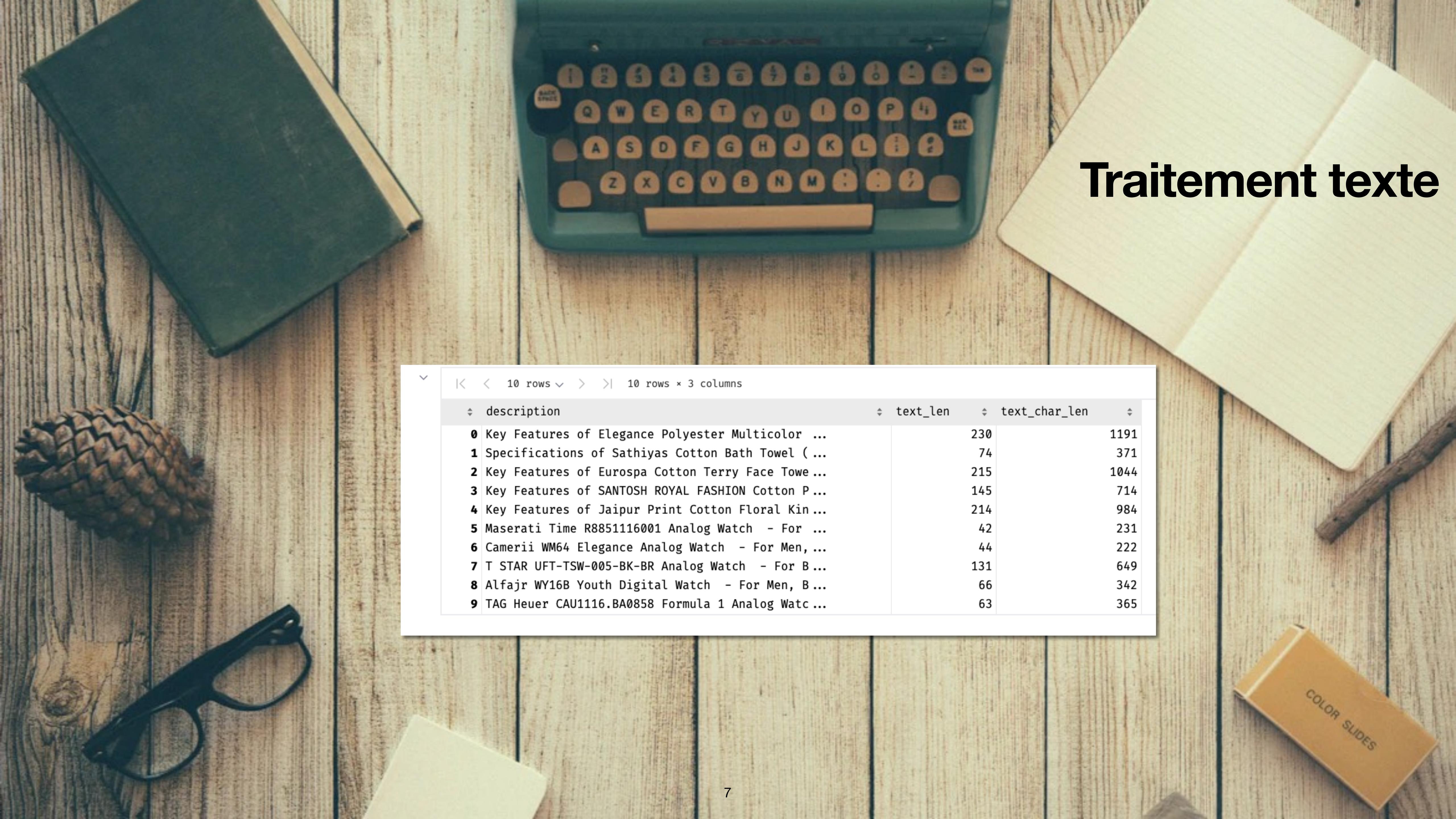


Niveaux 2



Niveaux 3

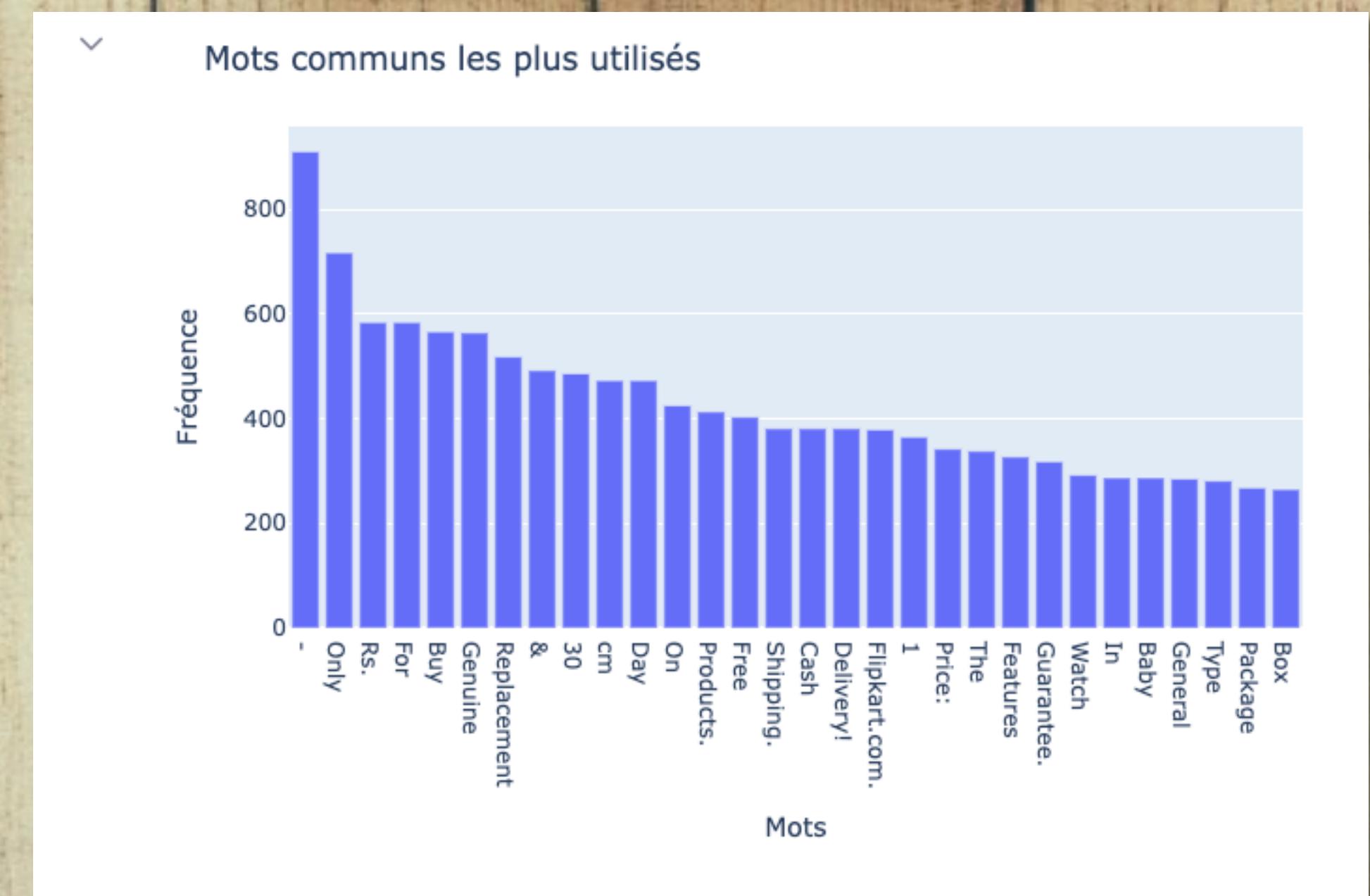
# Traitements de texte



A screenshot of a data visualization tool showing a table of 10 rows by 3 columns. The columns are labeled 'description', 'text\_len', and 'text\_char\_len'. The data consists of product descriptions and their lengths.

	description	text_len	text_char_len
0	Key Features of Elegance Polyester Multicolor ...	230	1191
1	Specifications of Sathiyas Cotton Bath Towel ( ...	74	371
2	Key Features of Eurospa Cotton Terry Face Towe ...	215	1044
3	Key Features of SANTOSH ROYAL FASHION Cotton P ...	145	714
4	Key Features of Jaipur Print Cotton Floral Kin ...	214	984
5	Maserati Time R8851116001 Analog Watch - For ...	42	231
6	Camerii WM64 Elegance Analog Watch - For Men, ...	44	222
7	T STAR UFT-TSW-005-BK-BR Analog Watch - For B ...	131	649
8	Alfajr WY16B Youth Digital Watch - For Men, B ...	66	342
9	TAG Heuer CAU1116.BA0858 Formula 1 Analog Watc ...	63	365

# Traitement texte



# Pré-traitement texte



```
In 4 1 def process_text(text, min_len_word=3, force_is_alpha=True):
2
3     tokenizer = RegexpTokenizer(r'\w+')
4     lemmatizer = WordNetLemmatizer()
5
6     # Convertit le texte en minuscules
7     text = str(text).lower()
8     raw_tokens_list = tokenizer.tokenize(text)
9     cleaned_tokens_list = [w for w in raw_tokens_list if w not in stopwords.words('english')]
10    non_rare_tokens = [w for w in cleaned_tokens_list if w not in list_rare_words]
11    more_than_N = [w for w in non_rare_tokens if len(w) ≥ min_len_word]
12
13    if force_is_alpha:
14        alpha_tokens = [w for w in more_than_N if w.isalpha()]
15    else:
16        alpha_tokens = more_than_N
17
18    lemmatized_tokens = [lemmatizer.lemmatize(w) for w in alpha_tokens]
19
20    if eng_words:
21        english_tokens = [w for w in lemmatized_tokens if w in eng_words]
22    else:
23        english_tokens = lemmatized_tokens
24
25    trans_text = ' '.join(english_tokens)
26
27    return trans_text
28
```

# Traitemetn texte

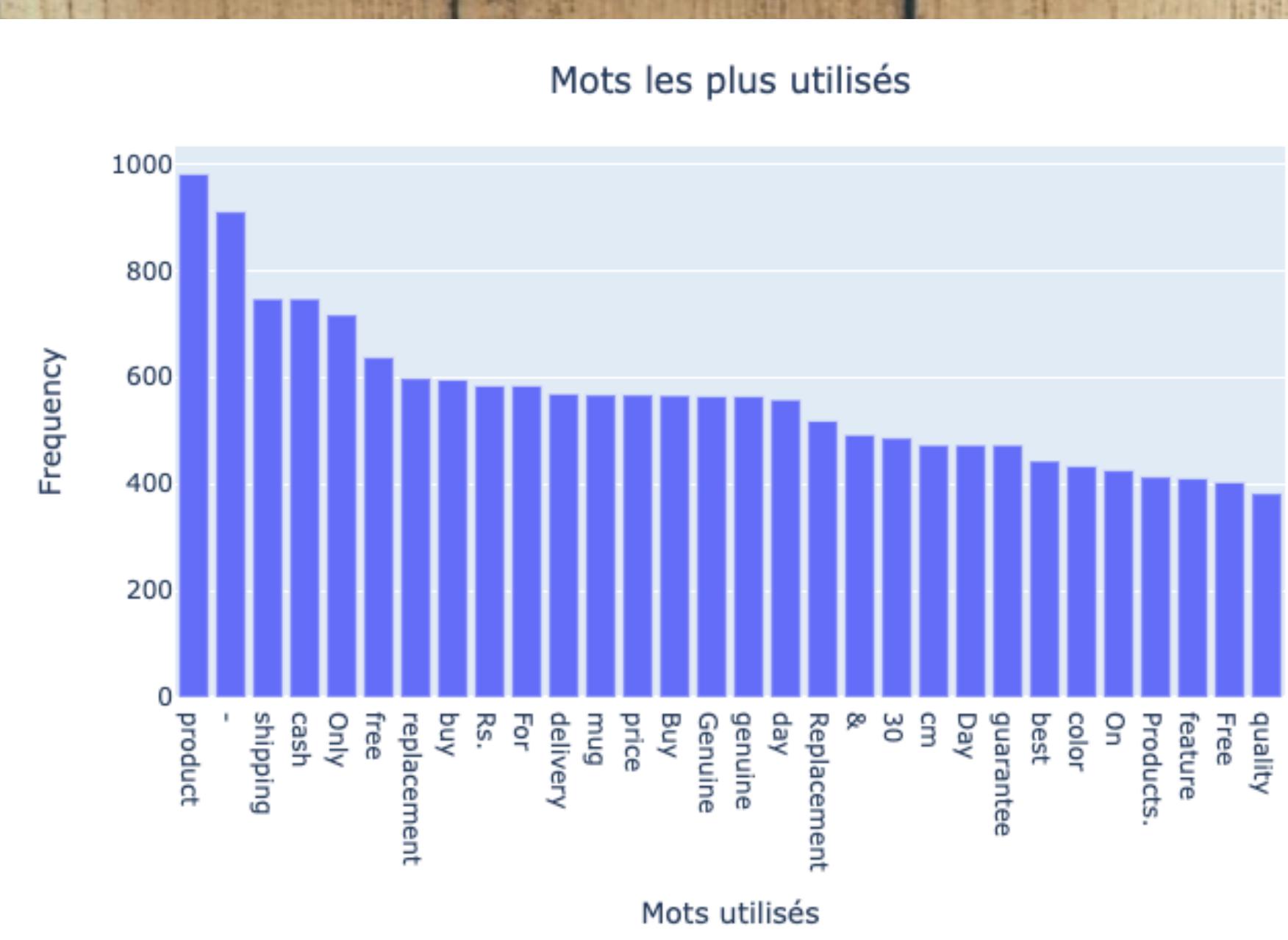
- Tokenization
- Stopwords
- Filtrage mots rares
- Filtrage longueur
- Lemmatisation
- Filtrage mots anglais
- Consturction texte final

# Traitements de texte



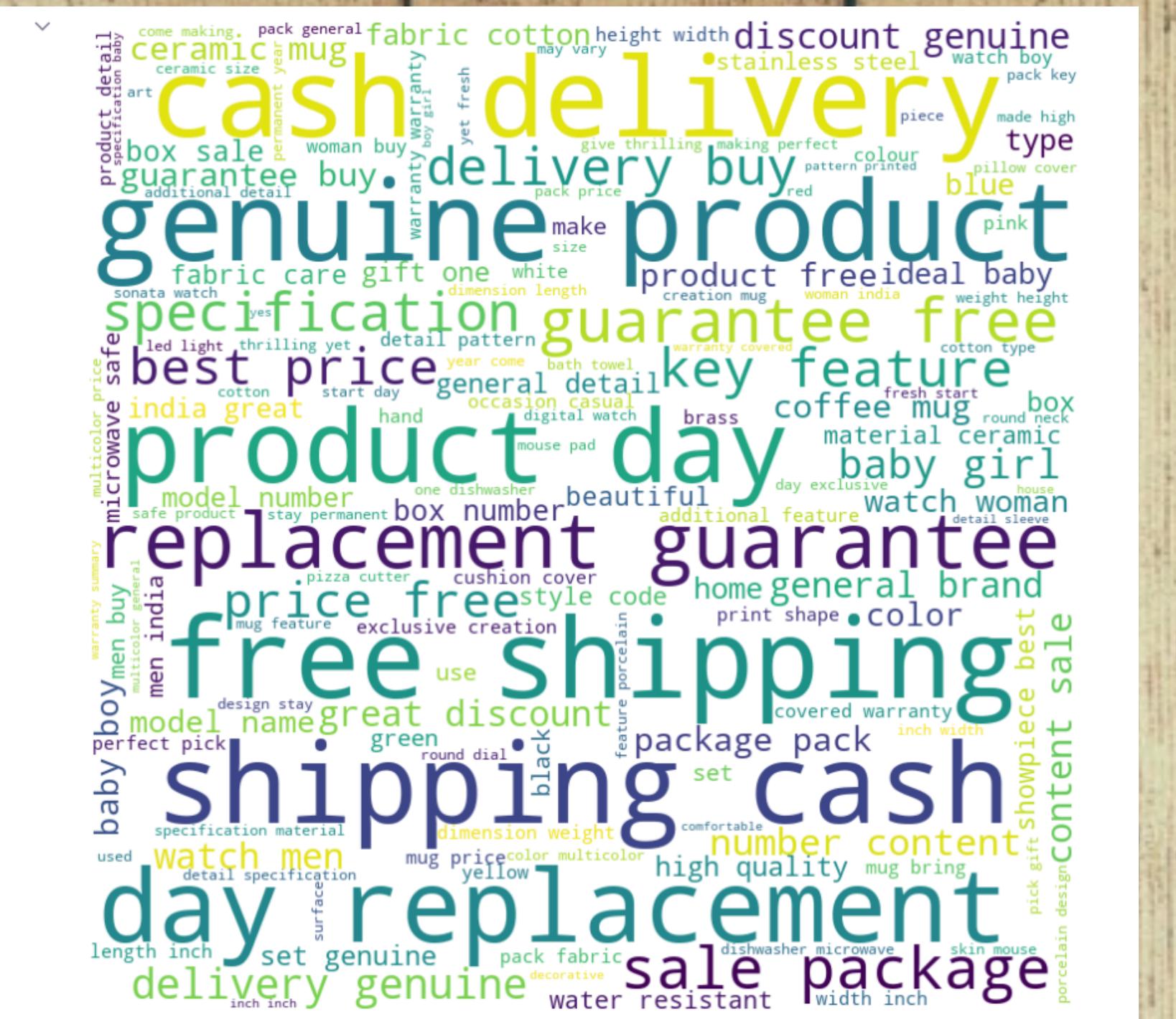
	description	clean_text	text_len	text_len_clean	text_char_len	text_char_len_clean
305	Buy Pg handicrafts Beaded two drawers sofa set ...	buy handicraft beaded two drawer sofa set show ...	45	28	231	173
245	Buy Kalash Kadhai 11.5 L for Rs.1584 online. K ...	buy best price free shipping cash delivery gen ...	29	12	144	73
563	Fluid FS201-BL01 Analog-Digital Watch - For W ...	fluid digital watch woman buy fluid digital wa ...	46	27	264	158
126	Flipkart.com: Buy Svayam Premium Manicure Duo ...	buy premium manicure duo stainless steel price ...	29	16	169	101
514	Maxima 04615CMGY Gold Analog Watch - For Men ...	maximum gold watch men buy maximum gold watch ...	50	29	263	163
59	Specifications of Brillare Science Dandruff Co ...	specification science dandruff control shampoo ...	41	24	228	152
682	Flipkart.com: Buy VLCC Natural Sciences Oil Fr ...	buy natural science oil free gel pack price ge ...	37	17	188	97
1039	Key Features of Reiki Crystal Products Showpie ...	key feature crystal product showpiece singing ...	270	140	1336	877
497	Times 284TMS284 Party-Wedding Analog Watch - ...	time party wedding watch woman buy time party ...	42	23	237	134
69	Specifications of GAYATRI CREATIONS ORANGE KID ...	specification gayatri creation orange set cont ...	48	33	225	184

# Traitement texte



# Traitement texte

## Wordcloud



# Traitement texte

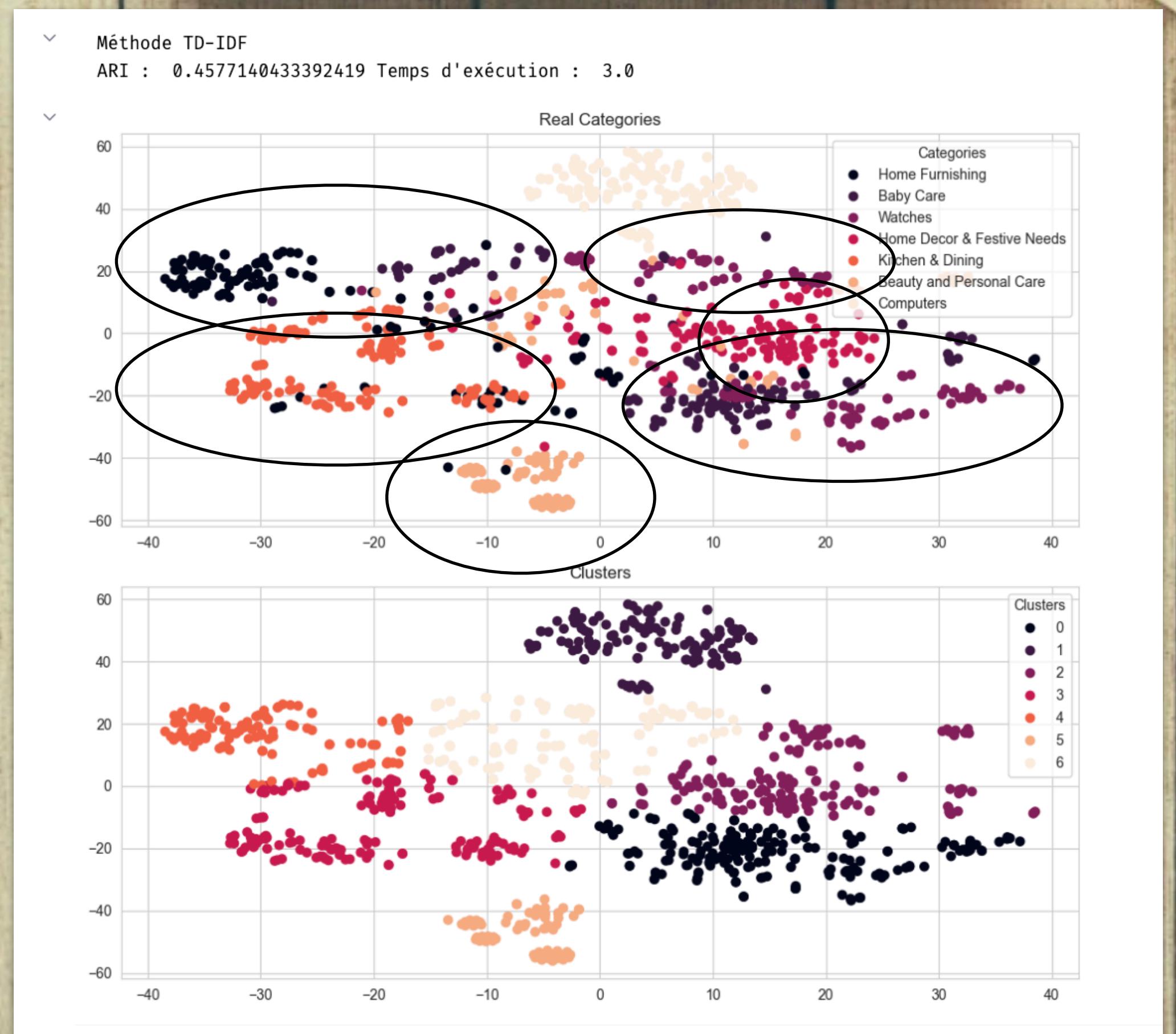
## Wordcloud



# Traitement texte

## Test Faisabilité

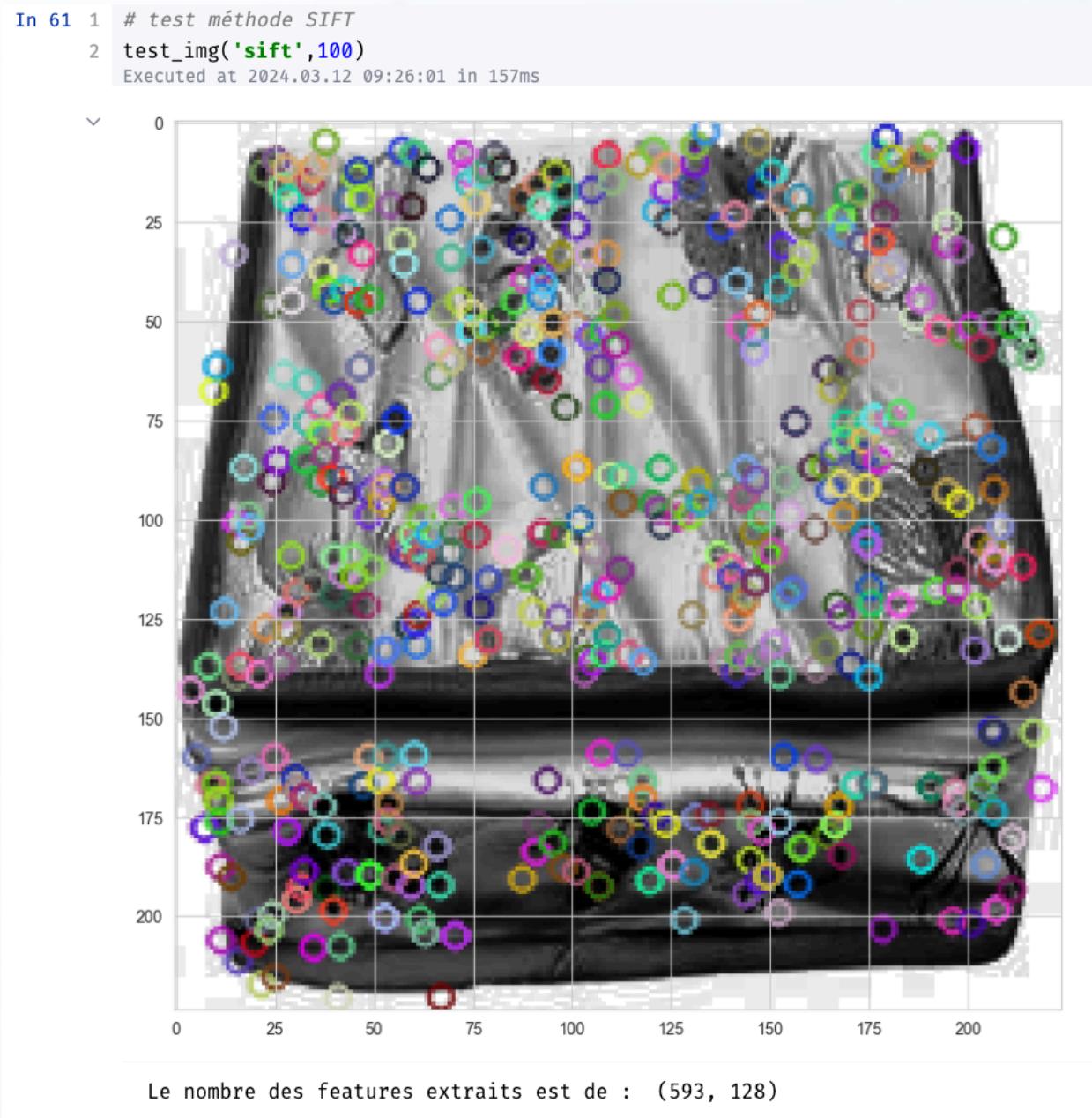
- TD-IDF
- Word2Vec
- BERT
- USE



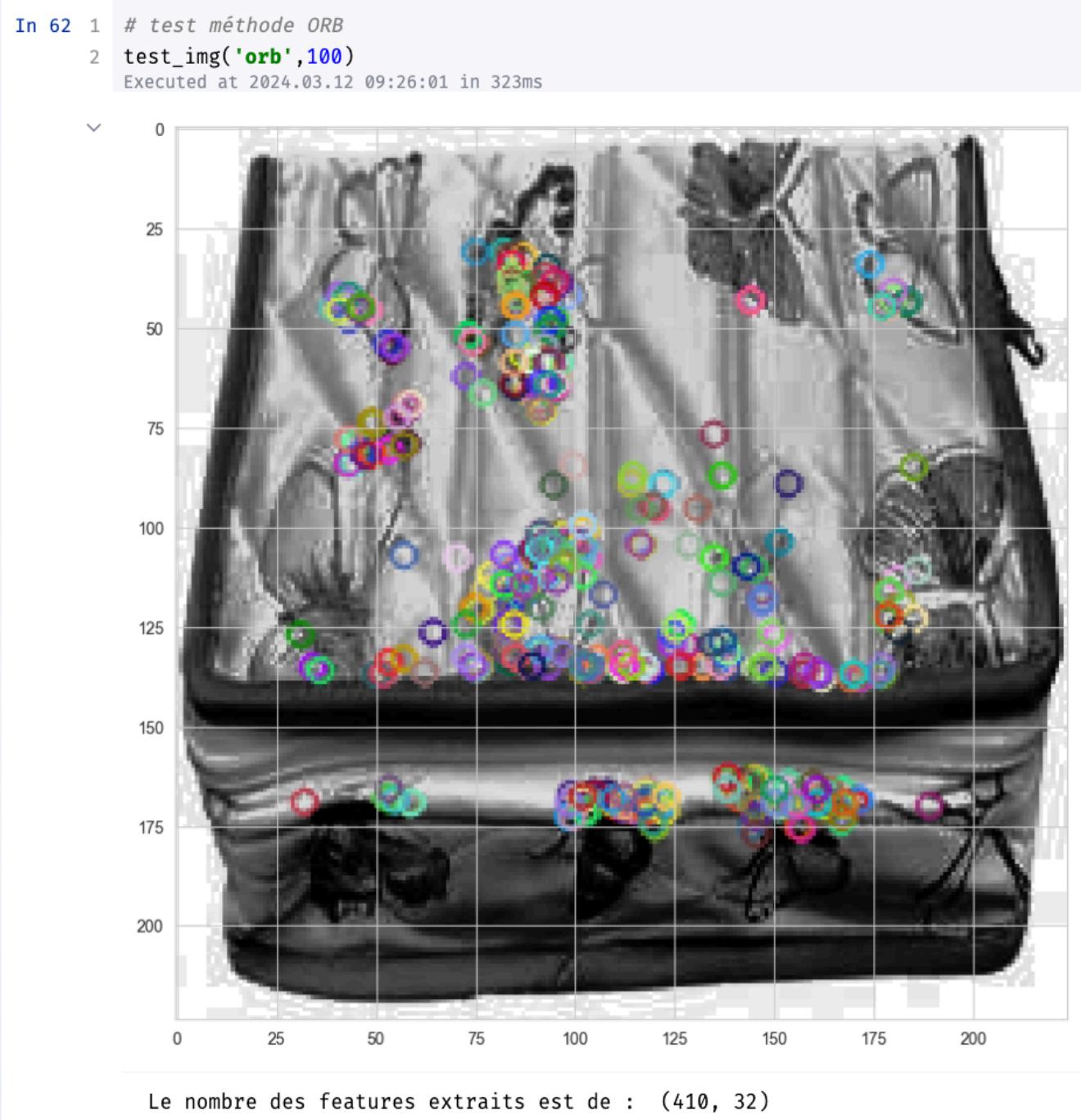
# Traitements images

## Extractions des features

Sift



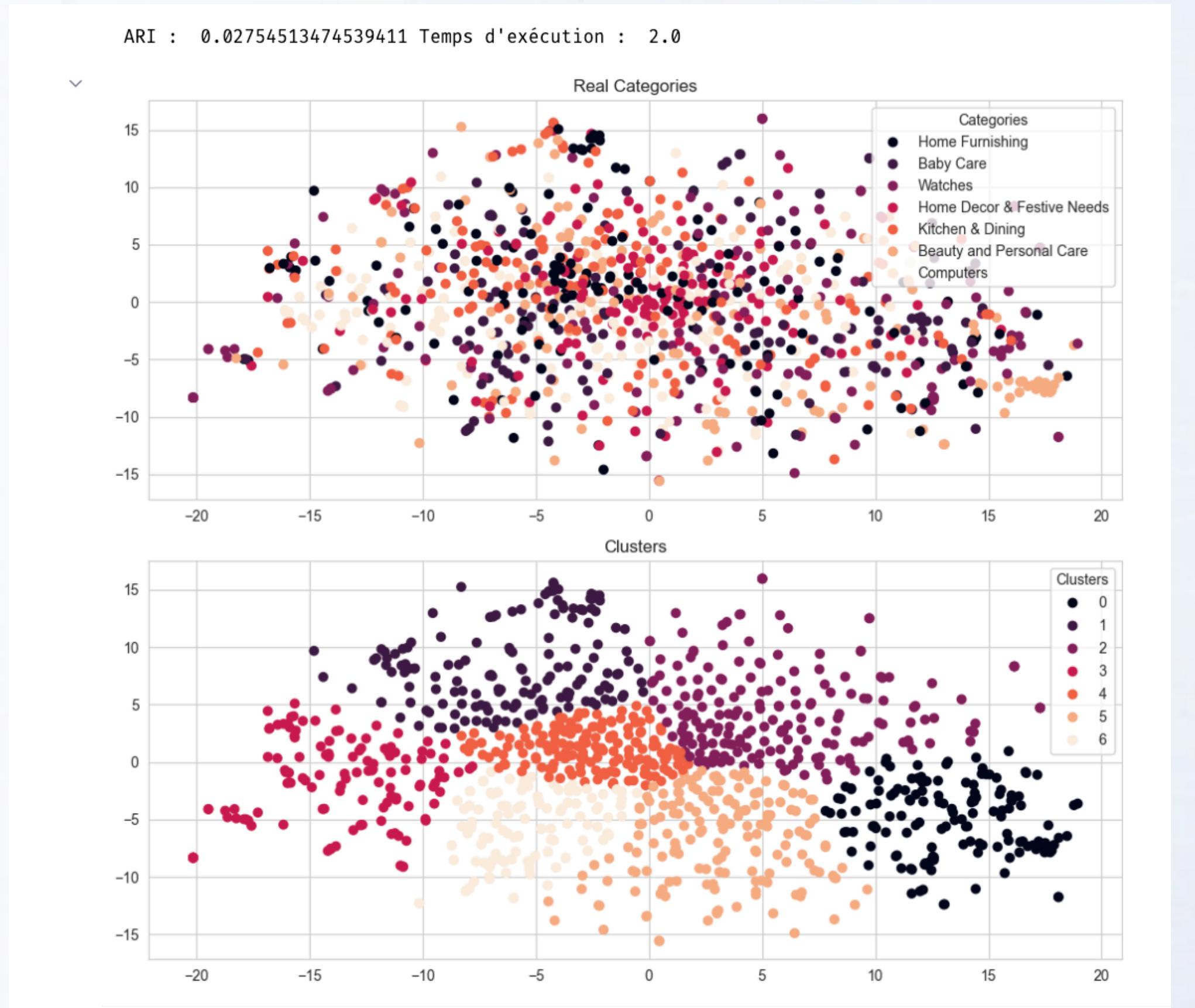
ORB



# Traitements images

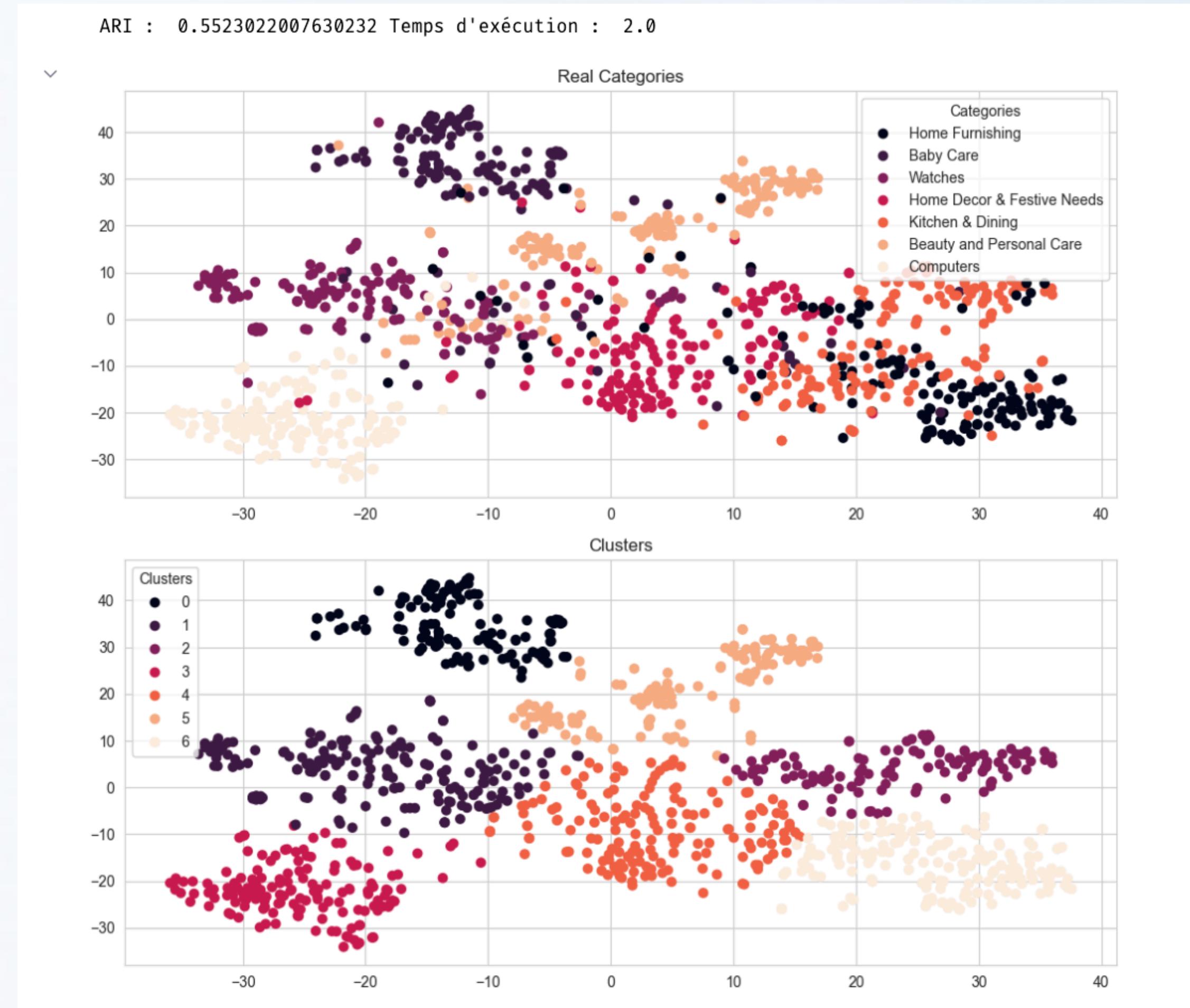
## Extractions des features

Sift



# Traitements images

VGG19



# Modélisation

## Data Generator

```
In 16 1 # Créer un générateur d'images
2 datagen = ImageDataGenerator(rotation_range=20,
3                               width_shift_range=0.2,
4                               height_shift_range=0.2,
5                               shear_range=0.2,
6                               zoom_range=0.2,
7                               horizontal_flip=True,
8                               rescale=1./255)
Executed at 2024.03.12 09:32:11 in 1ms

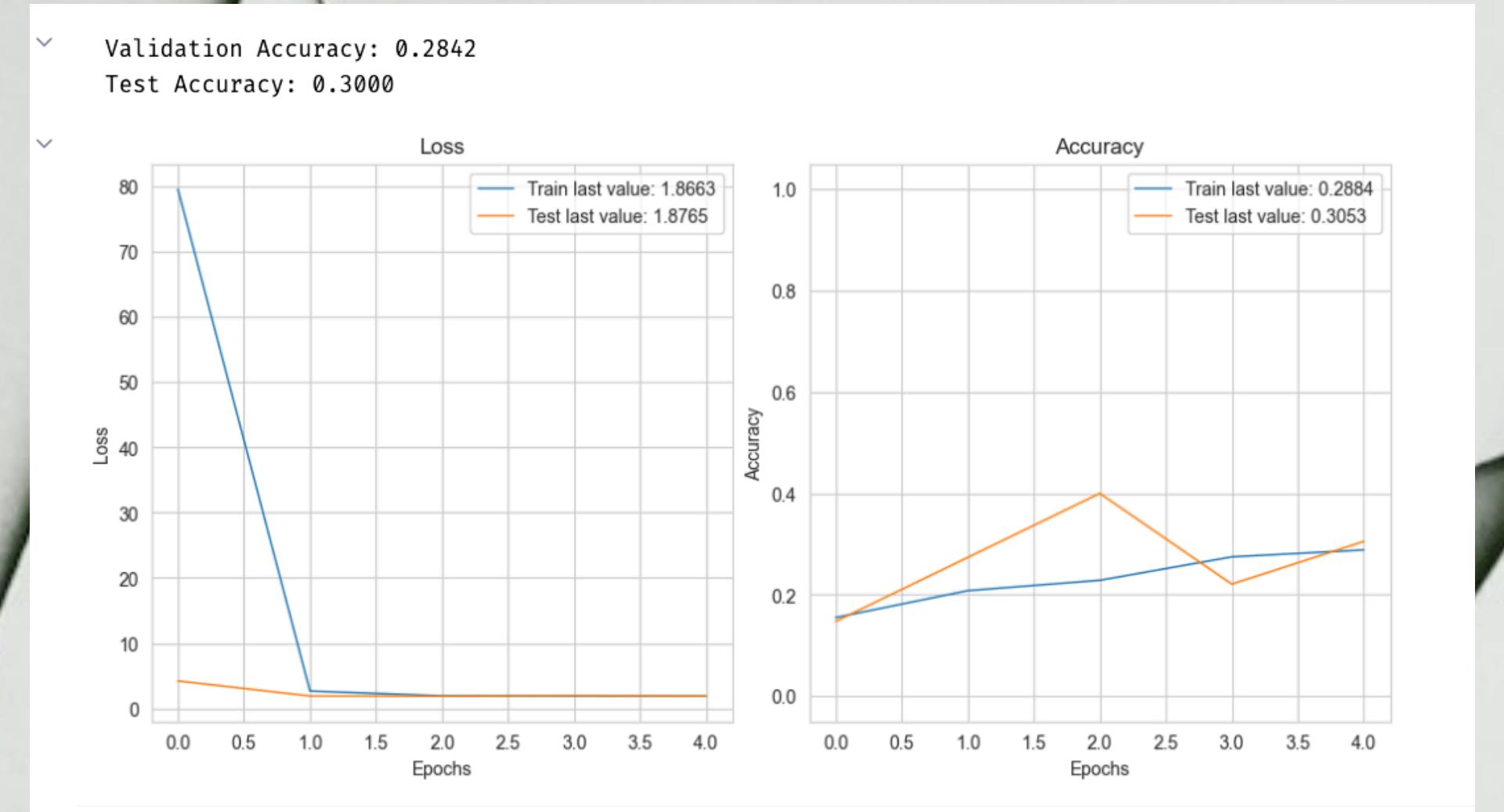
In 17 1 # création df pour train, test et validation
2 df_train = pd.DataFrame({"image": X_train, "label_categ": y_train})
3 df_test = pd.DataFrame({"image": X_test, "label_categ": y_test})
4 df_val = pd.DataFrame({"image": X_val, "label_categ": y_val})
Executed at 2024.03.12 09:32:11 in 2ms

In 18 1 # génération images pour set d'entraînement
2 train_generator = datagen.flow_from_dataframe(dataframe=df_train,
3                                               x_col="image",
4                                               y_col="label_categ",
5                                               target_size=(224, 224),
6                                               batch_size=32,
7                                               class_mode="categorical")
8
9 # Générer des images pour le set de test
10 test_generator = datagen.flow_from_dataframe(dataframe=df_test,
11                                              x_col="image",
12                                              y_col="label_categ",
13                                              target_size=(224, 224),
14                                              batch_size=32,
15                                              class_mode="categorical")
16
17 # générer images pour set de validation
18 val_generator = datagen.flow_from_dataframe(dataframe=df_val,
19                                              x_col="image",
20                                              y_col="label_categ",
21                                              target_size=(224, 224),
22                                              batch_size=32,
23                                              class_mode="categorical")
Executed at 2024.03.12 09:32:11 in 11ms
```

# Modélisation

## CNN

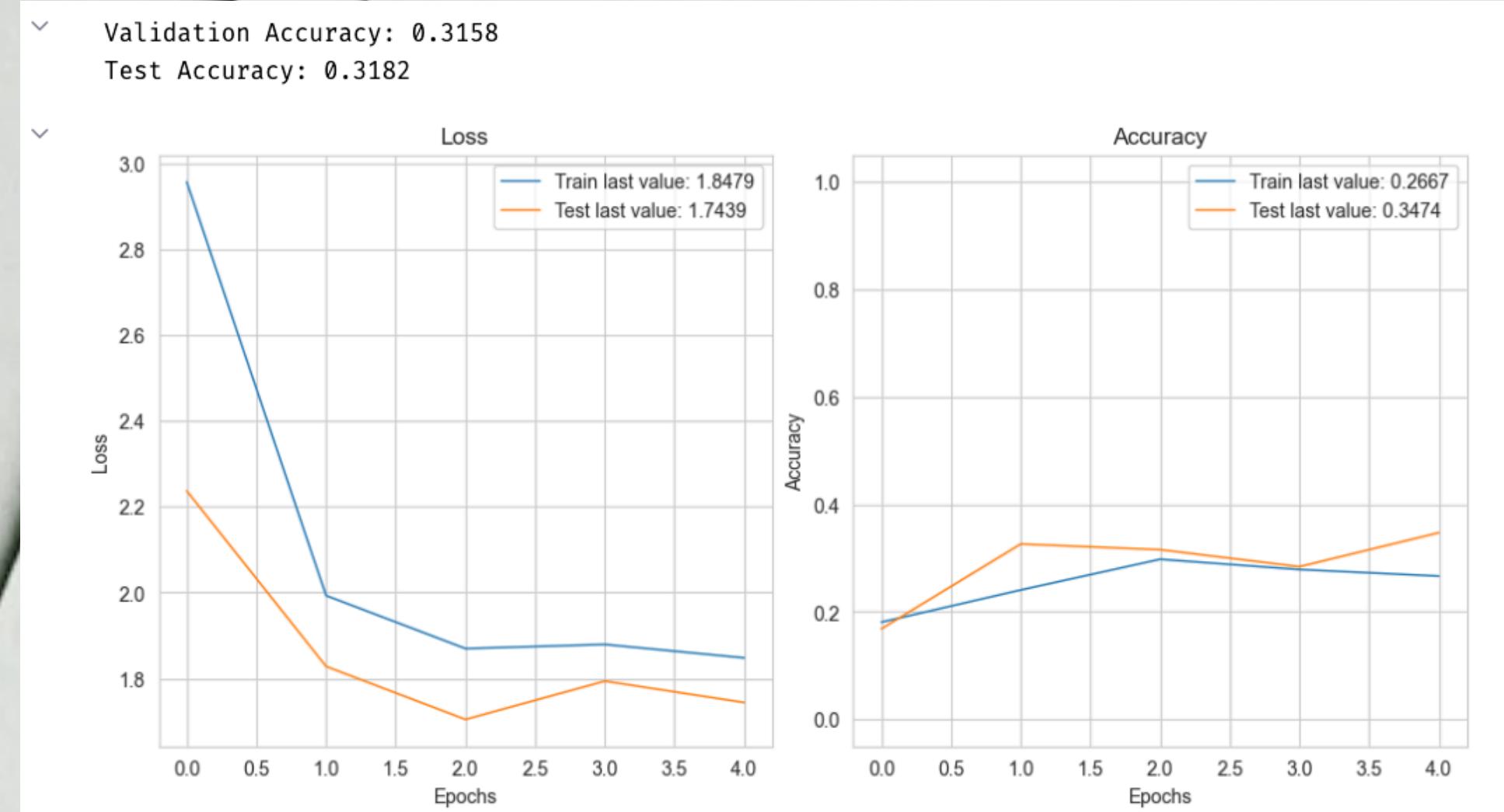
```
In 20 1 # Créer le modèle
2 model_cnn = Sequential()
3 model_cnn.add(Conv2D(64, (3, 3), activation='relu', input_shape=(224, 224, 3)))
4 model_cnn.add(MaxPooling2D((2, 2)))
5 model_cnn.add(Flatten())
6 model_cnn.add(Dense(128, activation='relu'))
7 model_cnn.add(Dense(7, activation='softmax'))
8
9 # Compiler le modèle
10 model_cnn.compile(loss='categorical_crossentropy', optimizer='rmsprop', metrics=['accuracy'])
11
12 # Afficher un résumé du modèle
13 model_cnn.summary()
14
15 checkpoint = ModelCheckpoint("../models/best_model_cnn.h5", monitor='val_accuracy', save_best_only=True, mode='max', verbose=1)
16
17 # Entrainer le modèle sur les données d'entraînement
18 history_cnn = model_cnn.fit(train_generator,
19                             batch_size=32,
20                             epochs=epochs,
21                             validation_data=val_generator,
22                             validation_batch_size=32,
23                             callbacks=[early_stop, checkpoint])
24
25 # Charger le meilleur modèle
26 best_model_cnn = load_model("../models/best_model_cnn.h5")
27
28 # Temps d'exécution
29 end_time = time.time()
30 execution_time = end_time - start_time
31 print("Temps d'exécution de la cellule : ", execution_time, " secondes")
32
```



# Modélisation

## RESNET 50

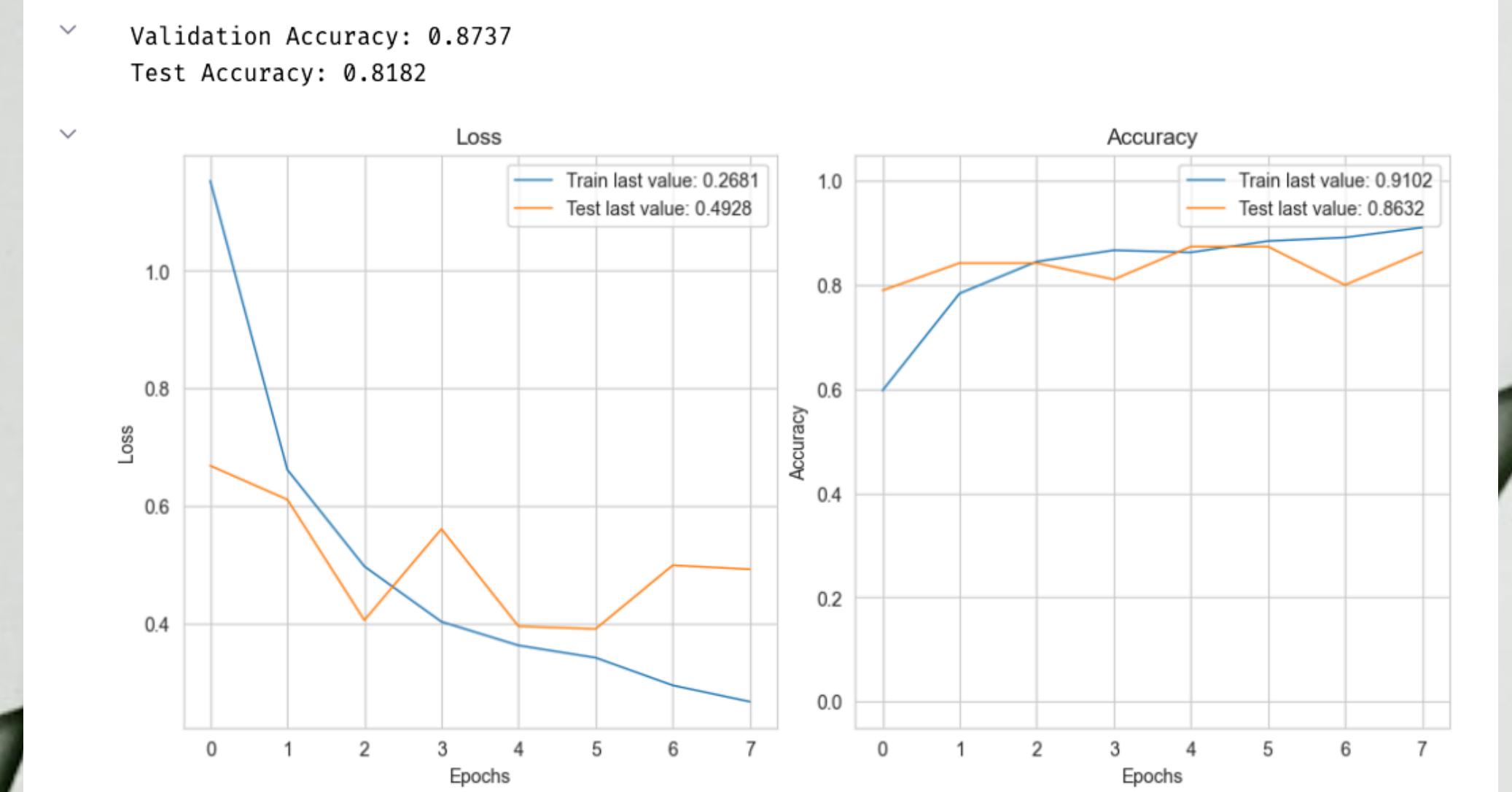
```
In 24 1 # Définir le tenseur d'entrée
2 inputs = Input(shape=(224, 224, 3))
3
4 # Charger le modèle ResNet50 pré-entraîné avec des poids pré-entraînés
5 resnet_base = ResNet50(weights='imagenet', include_top=False, input_tensor=inputs)
6
7 # Geler les couches pré-entraînées
8 for layer in resnet_base.layers:
9     layer.trainable = False
10
11 # Ajouter des couches supérieures personnalisées pour ma tâche spécifique
12 x = MaxPooling2D(pool_size=(2, 2))(resnet_base.output) # Couche de mise en commun
13 x = Flatten()(x) # Couche de dépliement
14 x = Dense(128, activation='relu')(x)
15 x = Dense(64, activation='relu')(x)
16 predictions = Dense(7, activation='softmax')(x) # Couche de sortie avec 7 neurones et activation softmax
17
18 # Créer le modèle final
19 model_resnet = Model(inputs=inputs, outputs=predictions)
20
21 # Compiler le modèle
22 model_resnet.compile(loss='categorical_crossentropy', optimizer=Adam(), metrics=['accuracy'])
23
24 # Afficher un résumé du modèle
25 model_resnet.summary()
26
27 # Point de contrôle pour sauvegarder le meilleur modèle
28 checkpoint = ModelCheckpoint("../models/best_model_resnet50.h5", monitor='val_accuracy', save_best_only=True, mode='max', verbose=1)
29 # Entrainer le modèle sur les données d'entraînement
30 history_resnet = model_resnet.fit(train_generator,
31                                     batch_size=32,
32                                     epochs=epochs,
33                                     validation_data=val_generator,
34                                     validation_batch_size=32,
35                                     callbacks=[early_stop, checkpoint])
36
37 # Charger le meilleur modèle basé sur la précision de validation
38 best_model_resnet = load_model("../models/best_model_resnet50.h5")
39 # Calculer et afficher le temps d'exécution
40 end_time = time.time()
41 execution_time = end_time - start_time
42 print("Temps d'exécution de la cellule : ", execution_time, " secondes")
```



# Modélisation

## Inception

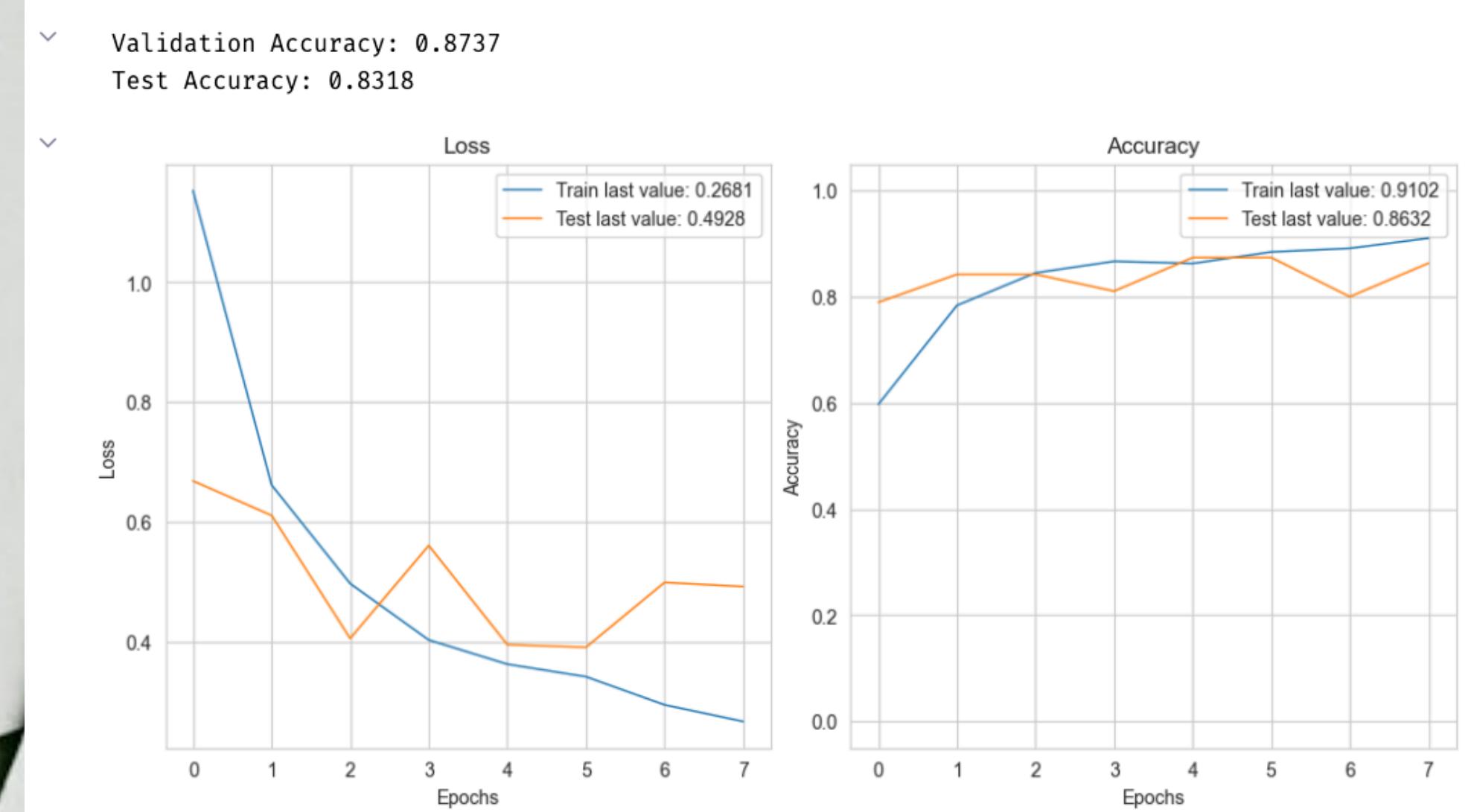
```
In 22 1 # Définir le tenseur d'entrée
2 inputs = Input(shape=(224, 224, 3))
3
4 # Charger le modèle InceptionV3 pré-entraîné avec les poids pré-entraînés
5 inception_base = InceptionV3(weights='imagenet', include_top=False, input_tensor=inputs)
6
7 # Geler les couches pré-entraînées
8 for layer in inception_base.layers:
9     layer.trainable = False
10
11 # Ajouter des couches supérieures personnalisées pour votre tâche spécifique
12 x = GlobalAveragePooling2D()(inception_base.output)
13 x = Dense(128, activation='relu')(x)
14 x = Dense(64, activation='relu')(x)
15 predictions = Dense(7, activation='softmax')(x) # Couche de sortie avec 7 neurones et activation softmax
16
17 # Créer le modèle final
18 model_inception = Model(inputs=inputs, outputs=predictions)
19
20 # Compiler le modèle
21 model_inception.compile(loss='categorical_crossentropy', optimizer=Adam(), metrics=['accuracy'])
22
23 # Afficher un résumé du modèle
24 model_inception.summary()
25
26 checkpoint = ModelCheckpoint("../models/best_model_inception.h5", monitor='val_accuracy', save_best_only=True, mode='max', verbose=1)
27
28 # Entrainer le modèle sur les données d'entraînement
29 history_inception = model_inception.fit(train_generator,
30                                         batch_size=32,
31                                         epochs=epochs,
32                                         validation_data=val_generator,
33                                         validation_batch_size=32,
34                                         callbacks=[early_stop, checkpoint])
35
36 # Charger le meilleur modèle basé sur la précision de validation
37 best_model_inception = load_model("../models/best_model_inception.h5")
38
39 # Calculer et afficher le temps d'exécution
40 end_time = time.time()
41 execution_time = end_time - start_time
42 print("Temps d'exécution de la cellule : ", execution_time, " secondes")
43
```



# Modélisation

## Xception

```
In 26 1 # Définir le tenseur d'entrée
2 inputs = Input(shape=(224, 224, 3))
3
4 # Charger le modèle Xception pré-entraîné avec les poids pré-entraînés
5 xception_base = Xception(weights='imagenet', include_top=False, input_tensor=inputs)
6
7 # Geler les couches pré-entraînées
8 for layer in xception_base.layers:
9     layer.trainable = False
10
11 # Ajouter des couches supérieures personnalisées pour votre tâche spécifique
12 x = GlobalAveragePooling2D()(xception_base.output)
13 x = Dense(128, activation='relu')(x)
14 x = Dense(64, activation='relu')(x)
15
16 predictions = Dense(7, activation='softmax')(x) # Couche de sortie avec 7 neurones et activation softmax
17
18 # Créer le modèle final
19 model_xception = Model(inputs=inputs, outputs=predictions)
20
21 # Compiler le modèle
22 model_xception.compile(loss='categorical_crossentropy', optimizer=Adam(), metrics=['accuracy'])
23
24 # Afficher un résumé du modèle
25 model_xception.summary()
26
27 checkpoint = ModelCheckpoint("../models/best_model_xception.h5", monitor='val_accuracy', save_best_only=True, mode='max', verbose=1)
28
29 # Entrainer le modèle sur les données d'entraînement
30 history_xception = model_xception.fit(train_generator,
31                                         batch_size=32,
32                                         epochs=epochs,
33                                         validation_data=val_generator,
34                                         validation_batch_size=32,
35                                         callbacks=[early_stop, checkpoint])
36
37 # Charger le meilleur modèle basé sur la précision de validation
38 best_model_xception = load_model("../models/best_model_xception.h5")
39
40 # Calculer et afficher le temps d'exécution
41 end_time = time.time()
42 execution_time = end_time - start_time
43 print("Temps d'exécution de la cellule : ", execution_time, " secondes")
44
```



# Modélisation

## Xception

## Inception

```
In 28 1 plot_bad_images(best_model_inception, val_generator, original_labels)
Executed at 2024.03.12 09:42:47 in 3s 211ms
```

3/3 [=====] - 3s 632ms/step

✓ Vraie classe : Computers  
Prédiction : Watches



Vraie classe : Watches  
Prédiction : Baby Care



Vraie classe : Watches Vraie classe : Home Decor & Festive Needs  
Prédiction : Home Decor & Festive Needs



```
In 29 1 plot_bad_images(best_model_xception, val_generator, original_labels)
Executed at 2024.03.12 09:42:52 in 4s 346ms
```

3/3 [=====] - 4s 1s/step

✓ Vraie classe : Computers  
Prédiction : Watches



Vraie classe : Computer  
Prédiction : Kitchen & Dining



Vraie classe : Home Decor & Festive Needs  
Prédiction : Baby Care

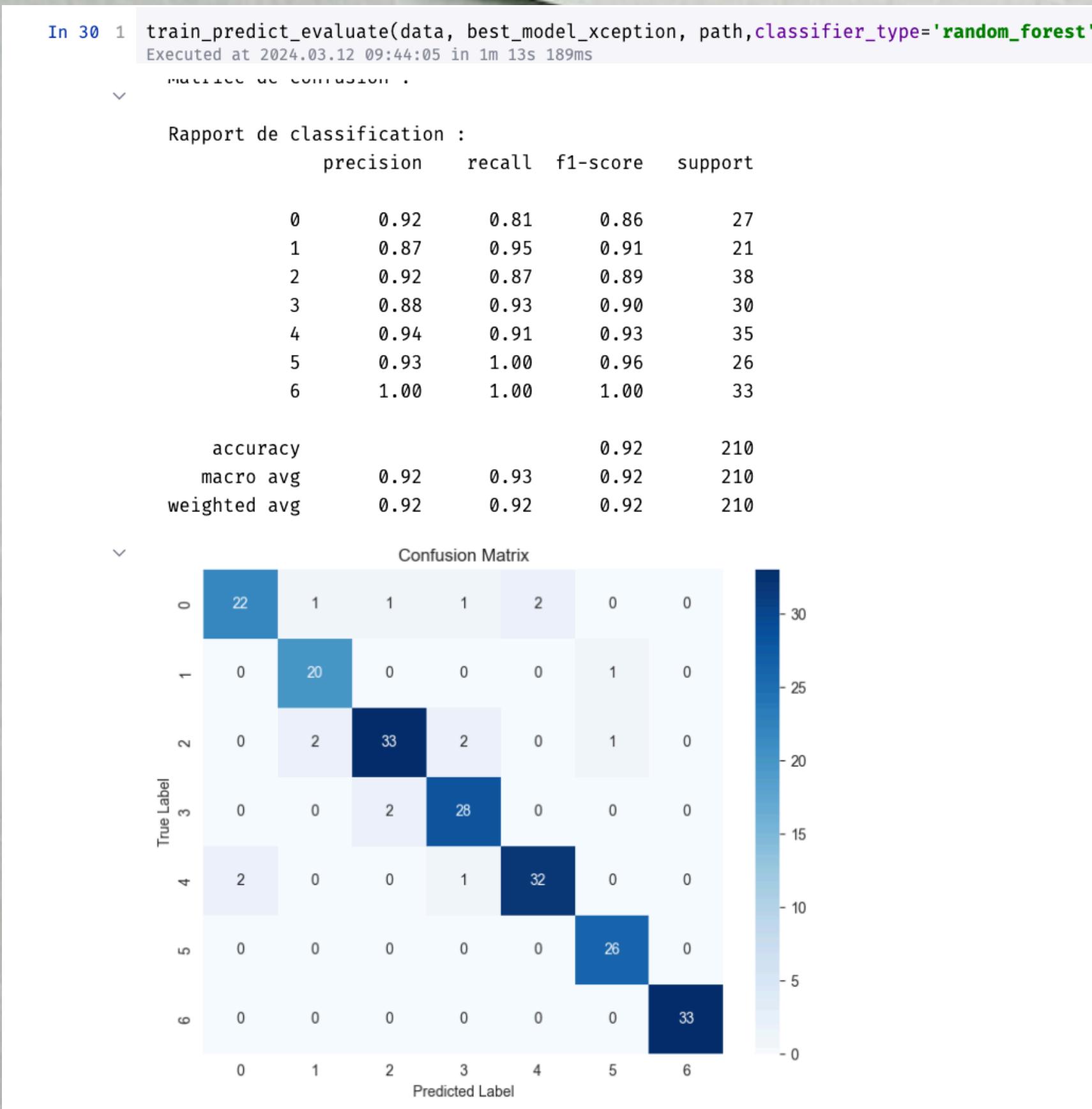


Vraie classe : Beauty and Personal Care  
Prédiction : Beauty and Personal Care



# Modélisation

## Modèle combiné



# Script

```
# URL de l'API
url = "https://edamam-food-and-grocery-database.p.rapidapi.com/api/food-database/v2/parser"

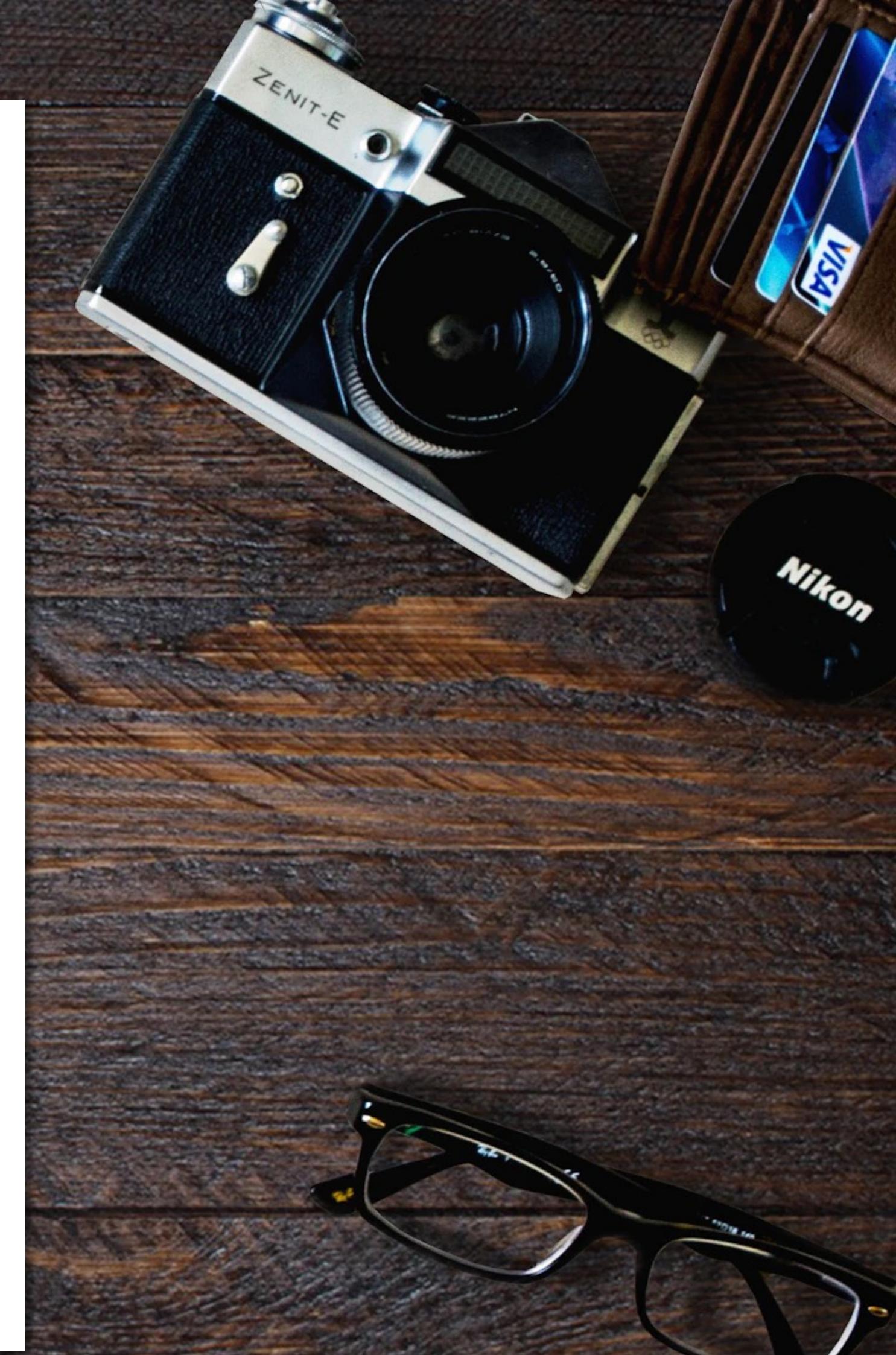
# Paramètres pour l'API
querystring = {"ingr": "champagne"}

# En-têtes de la requête pour l'API
headers = {
    "X-RapidAPI-Key": "279d3e1367msh23732bcf30a5598p1ec127jsn4ba6087d5aed",
    "X-RapidAPI-Host": "edamam-food-and-grocery-database.p.rapidapi.com"
}

# requête pour l'API et réponse
response = requests.request( method: "GET", url, headers=headers, params=querystring)

# Charger la réponse en format JSON
json_data = json.loads(response.text)

# Extraire les données pertinentes et les structurer en une liste de dict
csv_data = []
for item in json_data['hints']:
    row = {
        'foodId': item['food']['foodId'],
        'label': item['food']['label'],
        'category': item['food']['category'],
        'foodContentsLabel': item['food'].get('foodContentsLabel'),
        'image': item['food'].get('image')
    }
    csv_data.append(row)
```



# Script

```
# sauvegarde des données CSV dans un fichier
with open('output.csv', 'w', newline='') as csv_file:
    fieldnames = ['foodId', 'label', 'category', 'foodContentsLabel', 'image']
    writer = csv.DictWriter(csv_file, fieldnames=fieldnames)
    writer.writeheader()
    for row in csv_data:
        writer.writerow(row)

# Charger les données dans un dataframe et sélectionner colonnes souhaitées
df = pd.read_csv(filepath_or_buffer='output.csv', nrows=10, usecols=['foodId', 'label', 'category', 'foodContentsLabel', 'image'])

# Exporter le dataframe en fichier CSV
df.to_csv(path_or_buf='output_2.csv', index=False)

# Constat du fichier
df
```

10 rows × 5 columns						
Missing Count	0 Count	0 Count	0 Count	10 Count	9 Count	7 Count
10	7	Packaged ...	50%	10	nan	70%
Unique values	Unique values	Generic m...	40%	Unique values	https://w...	10%
Generic f...	10	Generic foods	10%	Other	Other	20%
0 food_a656mk2a5dmq...	Champagne	Generic foods	NaN	https://www.edama...		
1 food_b753ithamdb8...	Champagne Vinaigr...	Packaged foods	OLIVE OIL; BALSAM...	NaN		
2 food_b3dyababjo54...	Champagne Vinaigr...	Packaged foods	INGREDIENTS: WATE...	https://www.edama...		
3 food_a9e0ghsamvoc...	Champagne Vinaigr...	Packaged foods	CANOLA AND SOYBEA...	NaN		
4 food_an4jjueaucpu...	Champagne Vinaigr...	Packaged foods	WATER; CANOLA AND...	NaN		
5 food_bmu5dmkazwuv...	Champagne Dressin...	Packaged foods	SOYBEAN OIL; WHIT...	https://www.edama...		
6 food_alpl44taoyv1...	Champagne Butterc...	Generic meals	sugar; butter; sh...	NaN		

# Merci pour votre attention

