

## Etape 2 : Modèle de base

- Construction d'une baseline Zero-Shot a partir d'un modèle LLM
- Finetuning de ce modèle à partir du dataset (contrastive loss)

```
In [1]: import json
import matplotlib.pyplot as plt
from tqdm import tqdm
import numpy as np
import pickle

from sentence_transformers import SentenceTransformer
from sentence_transformers.util import cos_sim
from sentence_transformers.quantization import quantize_embeddings
from sentence_transformers import losses
from sentence_transformers.readers import InputExample
from torch.utils.data import DataLoader
```

```
/Users/mbp004/dev/test_yxir/yxir/lib/python3.11/site-packages/threadpoolctl.
py:1214: RuntimeWarning:
Found Intel OpenMP ('libiomp') and LLVM OpenMP ('libomp') loaded at
the same time. Both libraries are known to be incompatible and this
can cause random crashes or deadlocks on Linux when loaded in the
same Python program.
Using threadpoolctl may cause crashes or deadlocks. For more
information and possible workarounds, please see
https://github.com/joblib/threadpoolctl/blob/master/multiple\_openmp.md
warnings.warn(msg, RuntimeWarning)
```

### Methode Zero-Shot

On prend un modele du MTEB leaderboard pour realiser les embeddings de documents  
On teste le modele suivant : mxbai-embed-large-v1, proposé par mixedbread ai. Le  
modèle a été entraîné sur la tache de STS (Semantic Textual Similarity)

```
In [2]: # 1. Specify preferred dimensions
dimensions = 512

# 2. load model
model_name = 'intfloat/e5-small-v2'
# model_name = "mixedbread-ai/mxbai-embed-large-v1"
model = SentenceTransformer(model_name, truncate_dim=dimensions)
```

```
/Users/mbp004/dev/test_yxir/yxir/lib/python3.11/site-packages/huggingface_hu
b/file_download.py:1132: FutureWarning: `resume_download` is deprecated and
will be removed in version 1.0.0. Downloads always resume when possible. If
you want to force a new download, use `force_download=True`.
warnings.warn(
```



```
pickle.dump(dic_all_embeddings, fh)
fh.close()
```

```
In [7]: with open('../data/list_all_embeddings_e5-small-v2.pickle', 'rb') as fh:
        dic_all_embeddings = pickle.load(fh)
        fh.close()
        list_all_embeddings = dic_all_embeddings['embeddings']
```

```
In [8]: nb_good_embeddings = 0
for i in range(len(list_all_embeddings)):
    embeddings = list_all_embeddings[i]
    similarities = cos_sim(embeddings[0], embeddings[1:])
    sim_pos, sim_neg = similarities.flatten()
    if sim_pos > sim_neg :
        nb_good_embeddings+=1
perc_good_embeddings = round(100*nb_good_embeddings/len(list_all_embeddings))
print('Embeddings de documents compatibles avec la query: {}, {} %'.format(n
```

Embeddings de documents compatibles avec la query: 370, 74.15 %

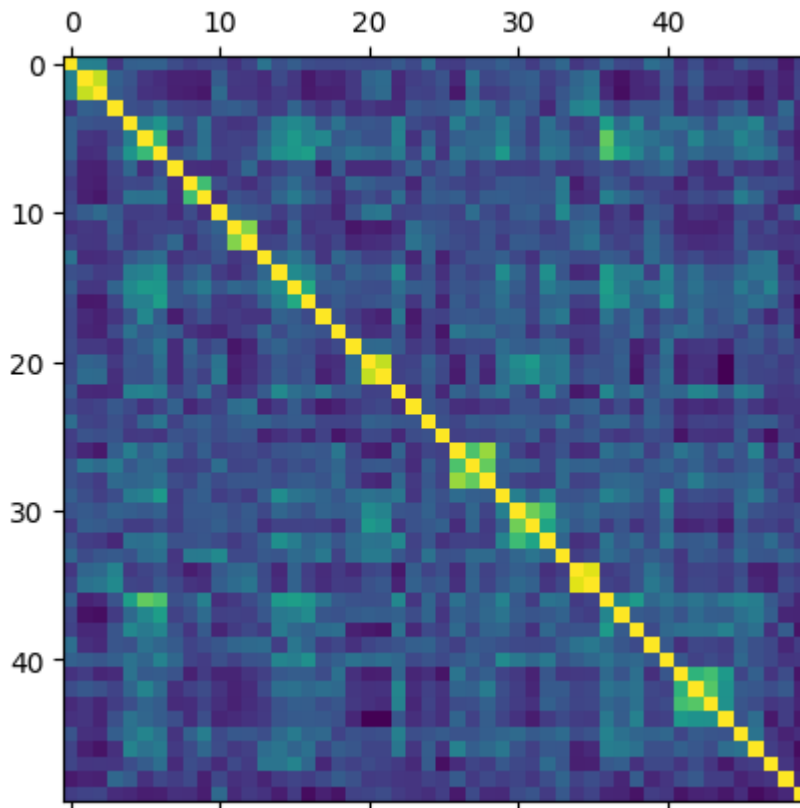
```
In [9]: # Heatmap sur les similarites entre contenus de brevet
list_embeddings_heatmap = []
for i in range(len(list_all_embeddings)):
    list_embeddings_heatmap.append(list_all_embeddings[i][1])
    list_embeddings_heatmap.append(list_all_embeddings[i][2])

list_similarities_heatmap = []
for i in tqdm(range(len(list_embeddings_heatmap)), desc="Calcul des similarites"):
    list_similarities_heatmap_i = []
    for j in range(len(list_embeddings_heatmap)):
        sim = float(cos_sim(list_embeddings_heatmap[i], list_embeddings_heatmap[j]))
        list_similarities_heatmap_i.append(sim)
    list_similarities_heatmap.append(list_similarities_heatmap_i)
similarity_matrix = np.array(list_similarities_heatmap)

plt.matshow(similarity_matrix[0:50,0:50])
plt.show()
```

Calcul des similarites: 100% |██████████  
██████████  
██████████ | 998/998 [01:07<00:00, 14.87it/s]





On voit un peu cet effet sur les embeddings des queries meme si c'est moins flagrant. A priori ceci est du à :

- la taille des queries qui est tres faible donc plus facile a discriminer

Les pistes que ça nous amène à regarder :

- restreindre les documents a certaines sections, pour eviter que le vecteur d'embedding soit noyé dans pleins d'informations
- finetuner le modele avec les données de brevet

## Performances en retrieval

Calcul des performances en retrieval sur une liste d'embeddings de type : `[[emb_query, emb_pos, emb_neg]...]` On cherche a calculer :

- le `top_K_accuracy` sur les labels positifs (a maximiser)
- le `top_K_accuracy` sur les labels negatifs (a minimiser)

```
In [35]: def compute_top_K_accuracy_score(list_embeddings, K=5):
    """
    Fonction pour calculer le top_K_accuracy score a partir d'une liste d'em
    [[emb_query, emb_pos, emb_neg]...]
    list_embeddings -- list, list des embeddings des query, positive, negati
    K -- int, le top K accuracy
    """
    list_embeddings_query = [list_embeddings[i][0] for i in range(len(list_a
```



```
model.fit(
    [(train_dataloader, train_loss)],
    epochs=10,
)
```

[illegible]

```
In [38]: # Recalculer les embeddings
list_all_embeddings_fineting = []
for i in tqdm(range(len(data)), desc="Calcul des embeddings apres finetuning")
    query = data[i]['query']
    pos_text = data[i]['pos']
    negative_text = data[i]['negative']

    docs = [
        query,
        pos_text,
        negative_text
    ]

    embeddings = model.encode(docs)
    list_all_embeddings_fineting.append(embeddings)

dic_all_embeddings = {'embeddings': list_all_embeddings_fineting}
with open('../data/list_all_embeddings_e5-small-v2_fineting.pickle', 'wb')
    pickle.dump(dic_all_embeddings, fh)
    fh.close()
```

```
Calcul des embeddings apres finetuning: 100%|██████████  
██████████  
██████████ | 499/499 [06:20<00:00, 1.31it/s]
```

```
In [39]: # Bonne classification
nb_good_embeddings = 0
for i in range(len(list_all_embeddings_finetuning)):
    embeddings = list_all_embeddings_finetuning[i]
    similarities = cos_sim(embeddings[0], embeddings[1:])
    sim_pos, sim_neg = similarities.flatten()
    if sim_pos > sim_neg :
        nb_good_embeddings+=1
perc_good_embeddings = round(100*nb_good_embeddings/len(list_all_embeddings_finetuning))
print('Embeddings de documents compatibles avec la query: {}, {} %'.format(nb_good_embeddings, perc_good_embeddings))
```

Embeddings de documents compatibles avec la query: 499, 100.0 %

```
In [40]: # top_K_accuracy
acc_K_pos, acc_K_neg = compute_top_K_accuracy_score(list_all_embeddings_fine
print('acc_K_pos : {} (a maximiser)'.format(acc_K_pos))
print('acc_K_neg : {} (a minimiser)'.format(acc_K_neg))
```

```
Calcul du top_K_accuracy score: 100%|███████████  
███████████  
███████████ | 499/499 [00:24<00:00, 20.68it/s]  
acc_K_pos : 0.19438877755511022 (a maximiser)  
acc_K_neg : 0.0 (a minimiser)
```

In [ ]:

In [ ]: