

Nom et Prénom : AMAR Papa
Matricule ULB / Filière : 000540605 / MA-STICS
Identifiant Github : Amar-ulb

STIC-B545 – Rapport TP3

Introduction

Ce document se veut un rapport du TP3 qui porte sur le clustering et l'entraînement d'un modèle Word2Vec sur un corpus.

Nous effectuerons d'abord un clustering des documents CAMille pour une décennie choisie et interpréterons les résultats. Ensuite il s'agira d'entraîner un modèle word2vec sur l'ensemble du corpus et d'explorer les relations entre vecteurs.

1. Clustering et interprétations des résultats obtenus

1.1. Clustering :

Decade = 1960 et Nombre de clusters désiré : K= 3

Une des principaux paramètres à régler dans le clustering K-means est le « n_clusters ».

J'aurais aimé utiliser la méthode Elbow pour déterminer le nombre optimal de clusters. Ne pouvant pas y parvenir (erreur sur mon code) , j'ai essayé arbitrairement plusieurs nombres comprise entre 2 et 10, avant de choisir le nombre 3. Nous verrons plus loin que le résultat fait sens.

Après avoir appliqué le clustering à l'aide de la fonction `'fit_predict'` j'obtiens 1000 documents qui correspondent aux articles de la décennie choisie. (Voir `output_s2_clustering.ipynb` dans TP3).

```
[67]: pprint(dict(clustering))
... Output exceeds the size limit. Open the full output data in a text editor
{0: ['KB_JB838_1960-01-21_01-00014.txt',
     'KB_JB838_1960-01-22_01-00015.txt',
     'KB_JB838_1960-02-24_01-00020.txt',
     'KB_JB838_1960-03-11_01-00022.txt',
```

Figure 1 - clustering

1.2. Visualisation des clusters

Pour visualiser les documents dans un espace 2D , les vecteurs sont réduits en 2 dimensions à l'aide de l'algorithme PCA.

```
pca = PCA(n_components=2)
reduced_vectors = pca.fit_transform(tfidf_vectors.toarray())
```

Figure 2-Algorithmme PCA

Des points (marker = "x") encore appelés centroïdes, sont définis pour former chaque cluster.

```
# Ajouter les centroïdes
centroids = pca.transform(km_model.cluster_centers_)
plt.scatter(centroids[:, 0], centroids[:, 1], marker = "x", s=100, linewidths = 2, color='black')
```

Figure 3- Ajout des centroïdes

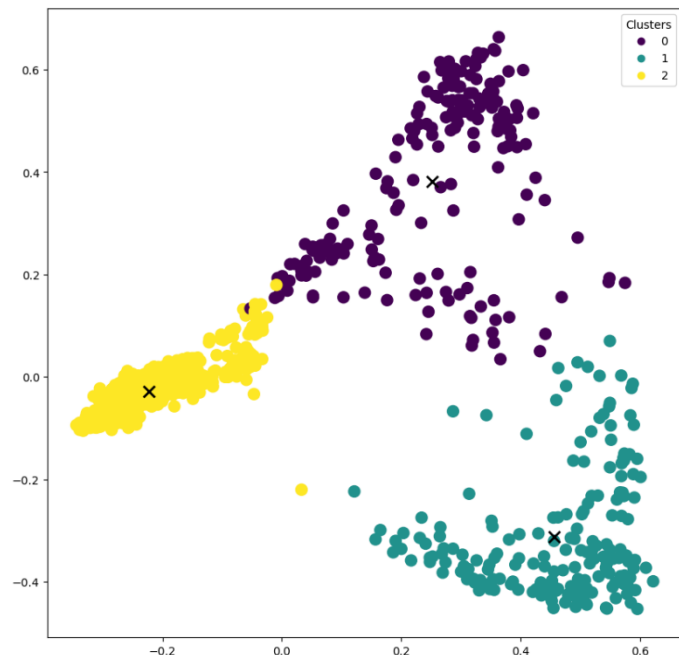
1.3. Interprétation des résultats obtenus

En utilisant K-Means avec k=3 nous avons obtenu le résultat de clustering représenté sur l'image ci-dessous. L'image de sortie montre clairement un nuage de points de notre ensemble de données

avec les 3 différents clusters identifiés par des couleurs. Chaque couleur représente en effet un cluster d'où la palette de 3 couleurs en référence au nombre de clusters prédéterminée.

Ce regroupement s'explique par le fait que le clustering K-means est une méthode de clustering itérative qui divise les données en entrée en k clusters dans lesquels chaque observation appartient au cluster avec la moyenne la plus proche (centroïde). Chaque centroïde (centre de gravité du cluster) est associé à des points de données. On constate 3 centroïdes éloignés les uns des autres. Le regroupement ainsi créé semble raisonnable selon moi, bien que la variance égale inégale dans chaque dimension rende la méthode moins adaptée à cet ensemble de données.

Figure 4 - Nuage de points



2. Entraînement d'un modèle word2vec et exploration des relations entre vecteurs

2.1. Entraînez un modèle word2vec (word embeddings)

Au vu de la quantité des données (pas très importante), nous avons changé quelques paramètres ci-après pour entraîner notre modèle.

Window = 2 : Pour éventuellement avoir des termes plus liés.

min_count = 10: Le modèle ignorera les mots qui ne satisfont pas le **min_count**. Les mots très peu fréquents sont généralement sans importance, il est donc préférable de s'en débarrasser.

workers=1 : Au vu du CPU cores de mon ordinateur, j'ai jugé sage de ne pas faire de parallélisation.

epochs=1 : J'ai constaté que l'augmentation du nombre de **epochs** augmente considérablement le temps de l'entraînement. Cela pourrait être envisager si nous avions de données plus importantes.

Finalement le modèle *word2vec* sauvegardé dans un fichier *"../data/newspapers.model"* est le suivant

```
Entraînement d'un modèle Word2Vec sur ce corpus

%%time
model = Word2Vec(
    corpus, # On passe le corpus de ngrams que nous venons de créer
    vector_size=32, # Le nombre de dimensions dans lesquelles le contexte des mots devra être réduit, aka. vector_size
    window=2, # La taille du "contexte", ici 2 mots avant et après le mot observé
    min_count=10, # On ignore les mots qui n'apparaissent pas au moins 10 fois dans le corpus
    workers=1, # Pas de parallélisation
    epochs=1 # Nombre d'itérations du réseau de neurones sur le jeu de données pour ajuster les paramètres avec la descente de gradient, aka. epochs
)

[145] ✓ 27m 19.5s

... CPU times: total: 16min 35s
Wall time: 27min 18s
```

Figure 5-Modèle word2vec

2.2. Explorer les relations entre vecteurs.

2.2.1. Calcul de la similarité entre deux termes avec la fonction *similarity* (3 exemples)

Les trois exemples ci-dessus calculent la similarité cosinusoidale entre les deux mots spécifiés en utilisant les vecteurs de mots de chacun. D'après les scores ci-dessus, il est logique que "directeur" soit similaire à "doyen" mais différent de "élève". À titre illustratif nous avons opter pour le 3ème exemple un terme identique pour montrer que lorsqu'on fait une similarité entre deux mots identiques, le score sera de 1,0 car la plage de similarité cosinus peut aller de [-1 à 1] et parfois délimitée entre [0,1] selon la façon dont elle est calculée et si les vecteurs ont à la fois des valeurs positives et négatives.

```
Calculer la similarité entre deux termes

model.wv.similarity("directeur", "doyen")
[73] ✓ 0.2s Python
... 0.8996973

model.wv.similarity("directeur", "eleve")
[86] ✓ 0.3s Python
... 0.012841117

model.wv.similarity("directeur", "directeur")
[94] ✓ 0.4s Python
... 1.0
```

Figure 6-Calcul de la similarité avec la fonction *similarity*

2.2.2. Recherche de mots les plus proches avec la fonction *most_similar* (3 exemples)

Dans l'ensemble, les résultats ont du sens. Tous les mots associés ont tendance à être utilisés dans le même contexte pour le mot de requête donné.

```
Chercher les mots les plus proches d'un terme donné

model.wv.most_similar("directeur", topn=5)
[152] ✓ 0.4s Python
... [('secretaire', 0.9204807877540588),
     ('directeur_general', 0.9079679250717163),
     ('membre', 0.9004807472229004),
     ('doyen', 0.8996972441673279),
     ('professeur', 0.8956307172775269)]

model.wv.most_similar("institut", topn=5)
[153] ✓ 0.2s Python
... [('academie', 0.8904618620872498),
     ('ecole', 0.8680539727210999),
     ('institut_national', 0.8586804270744324),
     ('association', 0.843223333587646),
     ('office', 0.829876184463501)]

model.wv.most_similar("femme", topn=5)
[154] ✓ 0.3s Python
... [('jeune_femme', 0.8930578231811523),
     ('jeune_fille', 0.8865285515785217),
     ('fille', 0.8428325653076172),
     ('femmo', 0.8267726898193359),
     ('bile', 0.8217665553092957)]
```

Figure 7-Recherche de mots les plus proches avec la fonction *most_similar*

2.2.3. Recherches complexes à travers l'espace vectoriel

- ✓ Exemple 1 : Analogie simple. L'homme est à la femme, ce que le roi est à...

L'exemple consiste à trouver des mots similaires à « femme » et « roi », et non similaires à « homme ». Le premier mot renvoyé par la requête est "queen" étant donné que le mot « roi » a déjà « homme » comme composant de sa représentation vectorielle.

- ✓ Exemple 2 : Analogies géographiques. Trouver dans quel pays se trouve une ville
- ✓ Exemple 3 : Analogies géographiques. La requête capture les relations entre un pays déterminé et les villes associées.

```
Faire des recherches complexes à travers l'espace vectoriel

[110] ✓ 0.2s Python
... [('reine', 0.7927921414375305)]

[116] ✓ 0.3s Python
... [('france', 0.8610690832138062)]

[155] ✓ 0.5s Python
... [('ath', 0.8340957164764404), ('etterbeek', 0.8246972560882568), ('evere', 0.813856840133667), ('enghien', 0.804882287979126), ('anderlecht', 0.8044556975364685), ('vilvorde', 0.7966938614845276), ('ixelles', 0.7956671118736267), ('ixelles', 0.7878463864326477), ('ypres', 0.7858937382698059), ('forest', 0.7854872941970825)]
```

Conclusion

À l'issu de ce TP, j'appréhende mieux trois enseignements vus en cours de traitement automatique de corpus :

- ✓ Prétraiter les données à utiliser avec un modèle Word2Vec
- ✓ Former un modèle Word2Vec
- ✓ Utilisez des métriques qualitatives, telles que *similarity* et *most_similar* pour évaluer la qualité des clusters formés

Par ailleurs je retiens aussi que le clustering avec K-Means est un type d'apprentissage automatique non-supervisé approprié pour regrouper selon un lien de similarité, une grande quantité de données non étiquetées en plusieurs sous-ensembles. Toutefois, je constate d'une part que c'est une méthode sensible aux valeurs aberrantes dans le paramétrage du clustering et l'entraînement du modèle et d'autre part le choix arbitraire de la valeur optimale de k n'est pas toujours facile. Enfin j'ajouterai en dernier lieu une petite remarque personnelle et non moins importante : il faut une machine relativement puissante pour la création d'un corpus et l'entraînement un modèle Word2Vec.