

# Repositionnement des médicaments

Cédric MARTINEZ, Roman RIEUNIER

Janvier 2020

## Introduction

Le repositionnement de médicaments est une problématique importante dans le corps médical. En effet une fois l'étude préclinique d'un nouveau médicament réalisée, il peut se passer beaucoup de temps avant qu'il soit mis sur le marché car il est nécessaire d'analyser tous les effets secondaires qu'il est susceptible de provoquer.

Le repositionnement de médicament permet alors de gagner du temps en permettant de réutiliser un médicament déjà mis sur le marché pour une nouvelle utilisation. Cette méthode s'appuie alors sur l'ensemble des bases de connaissances qu'il est possible d'avoir sur les médicaments.

Cependant cette mise en commun des bases de connaissance est limitée par la diversité des format dans lesquels sont enregistrées ces informations. Le cadre de notre projet était alors de réussir, grâce à des méthodes de [Machine Learning](#), d'automatiser la mise en commun d'informations provenant de bases de données différentes puis de créer un algorithme capable de prédire les repositionnements possibles des médicaments.

# 1 Mise en commun des bases de connaissances

## 1.1 Traitement des données

La première étape de notre projet est de recenser l'ensemble des bases de données que nous avons à disposition ainsi que de prévoir les formats qui pourront être attendues par de nouvelles bases de données.

Par conséquent il est important de prévoir un algorithme qui permet à la fois d'analyser des données *xml* comme dans la base de données de [DrugBank](#) ou des données *sdf* comme dans [DrugCentral](#).

### 1.1.1 Conversion *sdf* en *csv*

Nous avons donc décidé pour commencer de faire une recherche des différentes bibliothèques *Python* permettant de traiter les données qui nous intéressent. Le format *Structure Data File* ou *sdf* développé par Elsevier Molecular Design Limited est utilisé pour décrire des structures moléculaires en 3 dimensions. La bibliothèque [rdkit](#) permet de traiter facilement des données *sdf* en les stockant dans un *DataFrame* et en les convertissant éventuellement en *csv* comme dans le code ci-dessous.

Le format *csv* est intéressant car il peut être facilement adapté en *rdf* comme on le verra dans la suite de ce rapport.

```
1 import pandas as pd
2 from rdkit import Chem
3 from rdkit.Chem import PandasTools
4
5
6 my_sdf_file = './structures.molV2.sdf'
7
8 df = PandasTools.LoadSDF(my_sdf_file,
9                           smilesName='SMILES',
10                          molColName='Molecule',
11                          includeFingerprints=False)
12
13 df.to_csv('./structures.molV2.csv')
```

### 1.1.2 Conversion *xml* en *csv*

La bibliothèque [ElementTree](#) permet de charger et traiter des fichiers au format *xml*. On peut ainsi en extraire de l'information et l'exporter sous un autre format tel que le *csv*, importable dans des triplestores populaires comme Neo4J. Cette librairie est essentielle dans notre contexte puisque la grande majorité des données mises à disposition publiquement le sont en *xml*. Le jeu de données sur lequel nous travaillons principalement disponible sur [DrugBank](#) est justement au format *xml* et recense des informations sur plus de 6800 médicaments dont

1500 approuvés par la [FDA](#).

TODO

```
1 import xml.etree.ElementTree as ET
2 import pandas as pd
3
4 root = ET.parse('./full_database.xml').getroot()
5 ns = '{http://www.drugbank.ca}'
6
7 drugs = []
8 for drug in root.findall(f"{ns}drug"):
9     for product in drug.find(f'{ns}products').findall(f'{ns}product
10         '):
11         drugs.append({"Name": drug.find(f"{ns}name").text, "
12             Labeller": drug.find(f"{ns}description").text, "Countries":
13             product.find(f"{ns}country").text, "Drug": product.find(f"{ns}
14                 name").text})
15
16 pd.DataFrame(drugs).to_csv("drugbank.csv", sep=";")
```

## 1.2 Stockage des données

Une fois les données au format *csv* nous pouvons les stocker sous forme de graphe au format *rdf* grâce au logiciel [Neo4J](#).

En effet Neo4J, permet de charger des fichiers *csv* et de stocker les informations de ces fichiers dans un graphe de connaissances grâce à des commandes en langage [Cypher Query](#) développé exclusivement pour le logiciel. Pour charger un fichier *csv* il suffit alors d'utiliser la commande :

```
1 LOAD CSV FROM "file:///E:/pfe_data/drugs.csv"
2
3 // "file:/" va chercher le fichier dans le répertoire
4 // import/ du dossier de la base de données associée à notre
5 // projet Neo4J
6
7 // Sur mon ordinateur le fichier drugs.csv est placé dans le
8 // répertoire " E:\
9 //         neo4jDatabases\
10 //         database-fce6df9e-3f23-4c35-bc7b-5609a7078b86\
11 //         installation-3.5.14\import\
12 //         pfe_data"
13
14 // LOAD CSV FROM n'est pas utilisable seule
```

Dans notre version du logiciel nous récupérons les données de DrugBank sous la forme de 2 fichiers *csv*.

*drugs.csv* qui contient les champs :

<i>Name</i>	<i>Direct_parent</i>	<i>Kingdom</i>	<i>Superclass</i>
Lepirudin	Peptides	Organic Compounds	Organic Acids
...	...	...	...
<i>Class</i>	<i>Subclass</i>	<i>Synonyms</i>	<i>Description</i>
Carboxylic Acids and Derivatives	Amino Acids, Peptides, and Analogues	Hirudin variant-1, Lepirudin recombinant	Lepirudin is identical ...
...	...	...	...

*products.csv* qui contient les champs :

<i>Name</i>	<i>Labellers</i>	<i>Drug</i>
Refludan	Celgene Europe Limited EU, Bayer US, Bayer Canada	Lepirudin
...	...	...

Pour cela on utilise un premier script qui permet de récupérer les  $N$  premières molécules de la base de données de *DrugBank*.

```

1 TOTAL_DRUG = 30 # Nombre de médicaments qu'on veut
2                 # récupérer depuis la base de données complète
3
4 file = open('./drugbank_all_full_database.xml/full_database.xml',
5             'r', encoding="utf8")
6 count = 0
7 drug_count = 0
8 drug_nb = 0
9
10 with open('./drugbank_all_full_database.xml/partial_database.xml',
11           'w', encoding="utf8") as sub_file:
12     while drug_nb < TOTAL_DRUG:
13         count += 1
14
15         line = file.readline()
16
17         if not line:
18             break
19         sub_file.write(line)
20
21         if "<drug>" in line or "<drug " in line:
22             drug_count += 1
23         if "</drug>" in line:
24             drug_count -= 1
25             if drug_count == 0:
26                 drug_nb += 1
27         sub_file.write('</drugbank>')
28         sub_file.close()
29
30 file.close()

```

Puis on met ces données au format *csv* de manière à ce qu'elles puissent être facilement mise ensuite sous forme de graphe dans *Neo4J* et utilisées par un

algorithme de *Graph embedding* (dont on verra le fonctionnement dans la suite du rapport).

```
1 import xml.etree.ElementTree as ET
2 import pandas as pd
3
4 ns = '{http://www.drugbank.ca}'
5
6 root = ET.parse(
7     './drugbank_all_full_database.xml/\partial database.xml'
8 ).getroot()
9 drugs = []
10 products = []
11
12 # Transforme un tableau en string
13 # Utilisable facilement dans Neo4J
14 def a2str(array):
15     str = ""
16     for e in array:
17         str = str + e + ","
18     return str[:-1]
19
20 for drug in root.findall(f"{ns}drug"):
21     # Drug.csv
22
23     # Liste des attributs de classification
24     d_parent = "No parent"
25     kingdom = "No kingdom"
26     superclass = "No super c"
27     _class = "No class"
28     subclass = "No sub c"
29     if drug.find(f"{ns}classification") is not None:
30
31         if drug.find(f"{ns}classification").find(
32             f"{ns}direct-parent") is not None:
33             d_parent = drug.find(
34                 f"{ns}classification").find(
35                     f"{ns}direct-parent").text
36
37         if drug.find(f"{ns}classification").find(
38             f"{ns}kingdom") is not None:
39             kingdom = drug.find(
40                 f"{ns}classification").find(
41                     f"{ns}kingdom").text
42
43         if drug.find(f"{ns}classification").find(
44             f"{ns}superclass") is not None:
45             superclass = drug.find(
46                 f"{ns}classification").find(
47                     f"{ns}superclass").text
48
49         if drug.find(f"{ns}classification").find(
50             f"{ns}class") is not None:
51             _class = drug.find(
52                 f"{ns}classification").find(
53                     f"{ns}class").text
```

```

54         if drug.find(f"{ns}classification").find(
55             f"{ns}subclass") is not None:
56             subclass = drug.find(
57                 f"{ns}classification").find(
58                     f"{ns}subclass").text
59
60     synonyms = []
61     for syn in drug.find(
62         f'{ns}synonyms').findall(f'{ns}synonym'):
63         synonyms.append(syn.text)
64
65     drugs.append({"Name": drug.find(f"{ns}name").text,
66                 "Direct_parent": d_parent,
67                 "Kingdom": kingdom,
68                 "Superclass": superclass,
69                 "Class": _class,
70                 "Subclass": subclass,
71                 "Synonyms": a2str(synonyms),
72                 "Description": drug.find(
73                     f"{ns}description").text})
74
75 # Products.csv
76 labellers = {}
77 for product in drug.find(
78     f'{ns}products').findall(f'{ns}product'):
79     name = product.find(f"{ns}name").text
80
81     # On évite les doublons si les produits
82     # ont le même nom
83     if name not in labellers.keys():
84         labellers[name] = []
85
86     # On associe le labeller à son pays
87     # Si il y a plusieurs pays il y aura
88     # plusieurs noeuds ex : Bayer US / Bayer Canada
89     lab_str = product.find(f"{ns}labeller").text \
90         + " " \
91         + product.find(f"{ns}country").text
92     # On évite les doublons de labeller
93     if lab_str not in labellers[name]:
94         labellers[name].append(lab_str)
95
96     # On écris dans product.csv après avoir recensé tous
97     # les produits associés à une molécule pour éviter les
98     # doublons
99     for k in labellers: # k : nom du produit
100         products.append({"Name": k,
101                         "Labellers": a2str(labellers[k]),
102                         "Drug": drug.find(f"{ns}name").text})
103
104 pd.DataFrame(drugs).to_csv("E:/neo4jDatabases/database-fce6df9e-3
105     f23-4c35-bc7b-5609a7078b86/installation-3.5.14/import/pfe_data/
106     drugs.csv", sep=";")
107 pd.DataFrame(products).to_csv("E:/neo4jDatabases/database-fce6df9e
108     -3f23-4c35-bc7b-5609a7078b86/installation-3.5.14/import/
109     pfe_data/products.csv", sep=";")

```

Une fois qu'on a chargé le fichier *csv* on crée les entités de notre base de donnée. Etant donné que Neo4J ne stocke pas les valeurs *null* il faut également s'assurer que cela ne posera pas de problèmes :

```

1 LOAD CSV WITH HEADERS
2 FROM "file:///E:/pfe_data/drugs.csv" AS row
3 FIELDTERMINATOR ";"
4 WITH row WHERE row.Name IS NOT NULL
5 MERGE (m:Molecule {name: row.Name,
6       description: coalesce(row.Description, "No description")
7     })
8 MERGE (p:Molecule {name: coalesce(row.Direct_parent, "No parent")})
9 MERGE (m)-[:CHILD_OF]->(p)
10 MERGE (k:Kingdom {name: coalesce(row.Kingdom, "No kingdom")})
11 MERGE (m)-[:IN_KINGDOM]->(k)
12 MERGE (supc:Class {name: coalesce(row.Superclass, "No super c")})
13 MERGE (m)-[:IS_SUB_C]->(supc)
14 MERGE (c:Class {name: coalesce(row.Class, "No class")})
15 MERGE (m)-[:IS_CLASS]->(c)
16 MERGE (subc:Class {name: coalesce(row.Subclass, "No sub c")})
17 MERGE (m)-[:IS_SUPER_C]->(subc);
18
19 // Supprime les valeurs où le nom est null et
20 // met une description null à "No description"

```

Comme dans une base de données classique on peut également ajouter des contraintes pour éviter de future erreurs.

```

1 CREATE CONSTRAINT ON (d:Drug)
2 ASSERT d.name IS UNIQUE

```

```

1 CREATE CONSTRAINT ON (p:Product)
2 ASSERT p.name IS UNIQUE

```

```

1 CREATE CONSTRAINT ON (l:Labeller)
2 ASSERT l.name IS UNIQUE

```

Neo4J permet d'utiliser 2 commandes pour créer un nouveau noeud, **MERGE** ou **CREATE**, ici on préfère utiliser **MERGE** qui regarde d'abord si le noeud n'existe pas, au quel cas, il le crée.

Ensuite il faut préciser au logiciel les champs composés d'une liste d'items, par exemple ici les champs *Labellers* de *products.csv* et *Synonyms* de *drugs.csv*. On utilise alors le script ci-dessous.

```

1 LOAD CSV WITH HEADERS
2 FROM "file:///E:/pfe_data/products.csv" AS row
3 FIELDTERMINATOR ";"
4 WITH row WHERE row.Name IS NOT NULL
5 WITH row, split(row.Labellers, ',') AS labs

```

```

6 UNWIND labs AS lab
7 MERGE (p:Product {name: row.Name})
8 MERGE (l:Labeller {name: lab})
9 MERGE (p)-[:LABELLED_BY]->(l);

```

```

1 LOAD CSV WITH HEADERS
2 FROM "file:///E:/pfe_data/drugs.csv" AS row
3 FIELDTERMINATOR ";"
4 WITH row WHERE row.Name IS NOT NULL
5 WITH row, split(row.Synonyms, ',') AS syns
6 UNWIND syns AS syn
7 MERGE (m:Molecule {name: row.Name})
8 MERGE (s:Molecule {name: syn})
9 MERGE (m)-[:IS_SYNONYM]->(s)
10 MERGE (s)-[:IS_SYNONYM]->(m);

```

Enfin on relie les produits aux médicaments grâce au champ *Drug* de *products.csv* :

```

1 LOAD CSV WITH HEADERS
2 FROM "file:///E:/pfe_data/products.csv" AS row
3 FIELDTERMINATOR ";"
4 MATCH (p:Product {name: row.Name})
5 MATCH (m:Molecule {name: row.Drug})
6 MERGE (p)-[:CONTAINS]->(m);

```

Enfin, pour visualiser notre graphe, on peut utiliser la commande **MATCH** qui sélectionne des éléments du graphe.

```

1 MATCH (e) RETURN e
2 // Affiche la totalité du graphe

```

```

1 MATCH e=()-[r:CONTAINS]->(m:Molecule)
2 WHERE m.name="Lepirudin"
3 RETURN e
4 // Affiche uniquement les éléments dans une relation CONTAINS
5 // qui cible un noeud Lepirudin

```

Une fois ces scripts exécutés notre graphe de connaissance contient suffisamment de données pour pouvoir effectuer des prédictions de repositionnement. Le but de notre projet est donc maintenant de trouver un moyen, en utilisant des méthodes d'intelligence artificielle, de faire ces prédictions de manière automatique.



## 2 Prédiction de repositionnements

Maintenant que nous avons un graphe de connaissance contenant toutes les données dont nous avons besoin nous pouvons y appliquer un algorithme qui détecte automatiquement les repositionnements possibles en utilisant des méthodes de *Machine Learning*.

### 2.1 Plongement de mots

La méthode d'IA principale que nous avons utilisée dans notre projet est le Graph embedding.

Cette méthode repose sur une autre méthode, le [plongement de mots](#) ("Word embedding" en anglais), qui permet d'analyser des textes et de trouver les mots avec le plus de similitudes. Pour comprendre son fonctionnement nous avons regardé la bibliothèque [Word2vec](#) utilisée pour faire du plongement de mots.

L'idée principale de Word2Vec est de représenter les mots par des vecteurs qui sont similaires si les mots ont la même utilisation dans le texte utilisé pour l'entraînement du réseau de neurones, c'est à dire si leur entourage est similaire. Dans le code suivant :

```
1 data = [  
2     ["paracétamol", "molécule", "antalgique",  
3     "antipyrétique", "médicament"],  
4  
5     ["aspirine", "molécule", "antalgique",  
6     "antipyrétique", "médicament"],  
7  
8     ["estazolam", "molécule", "anxiolytique",  
9     "hypnotique", "sédatif"],  
10 ]  
11  
12 from gensim.test.utils import common_texts, get_tmpfile  
13 from gensim.models import Word2Vec  
14  
15 model = Word2Vec(data,  
16     min_count=1, # Nombre minimale d'occurrences  
17     # du mot pour qu'il soit pris en compte  
18     size= 10, # Nombre de dimension dans le  
19     # vecteur de similarités  
20     workers=3, # Nombre de partitions lors de  
21     # l'entraînement du modèle  
22     window =5, # Nombre de mots avant et après pour  
23     # établir l'entourage  
24     sg = 0) # Algorithme d'entraînement  
25     # CBOW(0) ou Skip-Gram(1)
```

On donne en entrée au modèle *Word2Vec* des descriptions des molécules de paracétamol, aspirine et estazolam par mots clés. Le modèle va attribuer un

vecteur de dimension 10 à chaque mot et les mots ayant un sens proche auront également un vecteur proche. La bibliothèque permet ensuite d'observer les similitudes entre les mots entre 100% et -100% et on obtient après entraînement les résultats suivants :

```
1 model.similarity('paracétamol', 'aspirine')
2 # 24%
3
4 model.similarity('paracétamol', 'estazolam')
5 # -19%
6
7 model.similarity('estazolam', 'aspirine')
8 # -42%
```

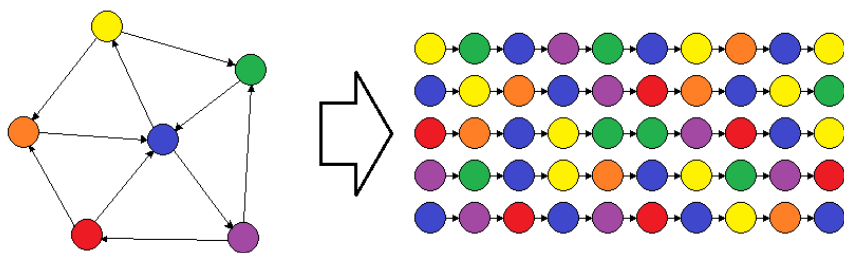
Par conséquent le modèle a bien compris que le paracétamol et l'aspirine sont plus proches entre eux que de l'estazolam qui est un somnifère.

## 2.2 Plongement sur des graphes (Graph embedding)

### 2.2.1 Graph Neural Network (GNN)

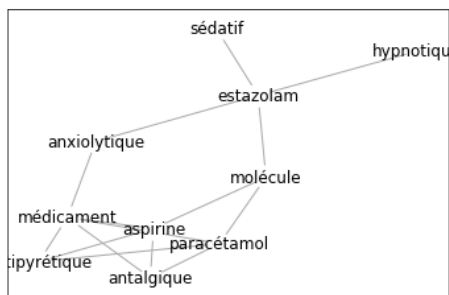
Ce même principe peut alors être appliqué sur des données sous forme de graphe comme le fait notamment la bibliothèque [Node2vec](#).

Pour cela il faut considérer les noeuds comme des mots et les connections comme l'entourage du noeud. Pour obtenir les vecteurs attendus par l'algorithme *Skip-Gram* il suffit alors de parcourir le graphe et d'enregistrer les positions successives comme dans la figure ci-dessous :



Sous forme de graphe, les données que nous avons utilisées précédemment prendraient environ la même forme que sur le graphe de droite.

A partir de ce format de données il est possible d'utiliser la bibliothèque *Node2vec* et ainsi retrouver les similitudes entre les mots.



Pour cela on exécute le script suivant :

```
1 from node2vec import Node2Vec
2
3 node2vec = Node2Vec(graph, dimensions=10, workers=3,
4                     num_walks=100, # Nombre de parcours du graphe
5                     walk_length=50) # Taille des parcours
6
7 model = node2vec.fit(window=5, min_count=1, batch_words=4)
```

Comme avec *Word2vec* on peut regarder la similitude entre paracétamol, aspirine et estazolam, et vérifier que le paracétamol est plus proche de l'aspirine que de l'estazolam pour être sûr que le modèle marche :

```
1 # Avec walk_length=10
2 model.wv.most_similar('paracétamol')
3 # Aspirine :    99%
4 # Estazolam :   96%
5
6 # Avec walk_length=20
7 model.wv.most_similar('paracétamol')
8 # Aspirine :   100%
9 # Estazolam :   43%
10
11 # Avec walk_length=50
12 model.wv.most_similar('paracétamol')
13 # Aspirine :   99%
14 # Estazolam :  -7%
```

Comme on peut le voir le modèle fonctionne très bien et on peut également remarquer qu'augmenter la longueur des parcours augmente significativement la précision du modèle.

Cette caractéristique est importante car elle veut dire que même avec un très petit jeu de données on peut avoir un modèle performant et simulant une quantité importante de parcours comme données d'entraînement ce qu'il n'était pas possible de faire avec les données sous forme de tableaux.

### 2.2.2 Visualisation t-SNE

Enfin, pour mieux visualiser les similitudes entre les mots nous pouvons utiliser l'algorithme [t-SNE](#).

Premièrement on peut afficher l'intégralité des mots sur un repère 2 dimensions tout en conservant les similitudes. Pour cela le code suivant récupère respectivement dans els tableaux *word\_cluster* et *embedding\_cluster* les différents mots du vocabulaire ainsi que leurs plongements.

```

1 word_cluster = []
2 embedding_cluster = []
3
4 for word in model.wv.vocab:
5     word_cluster.append(word)
6     embedding_cluster.append(model[word])

```

Ces tableaux sont ensuite passés à un modèle de *t-SNE* qui va fournir une représentation en 2 dimensions des similitudes entre les mots :

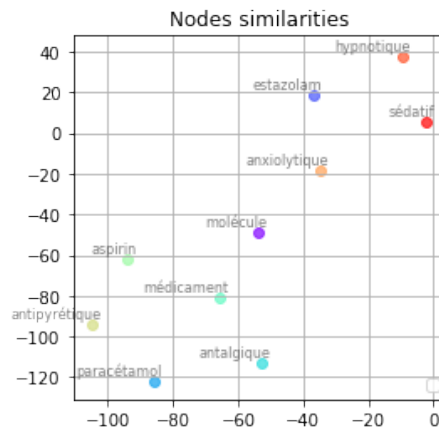
```

1 from sklearn.manifold import TSNE
2 import numpy as np
3
4 embedding_cluster = np.array(embedding_cluster)
5 n, m = embedding_cluster.shape
6 tsne_model_en_2d = TSNE(
7     perplexity=6, # Valeur relié à la taille des données
8                  # d'entrées, normalement entre 5 et 50,
9                  # ici 6 car le graphe utilisé comme
10                  # exemple est très petit
11     n_components=2, # Dimension attendue en résultat
12     init='pca', # PCA ou aléatoire, PCA est plus stable
13     n_iter=3500, # Nombre d'itérations pour
14                 # l'optimisation du résultat
15     random_state=32) # Graine pour le facteur aléatoire
16
17 embeddings_en_2d = np.array(tsne_model_en_2d.fit_transform(
18     embedding_cluster.reshape(n, m))).reshape(n, 2)

```

Enfin le code ci-dessous permet d'afficher le résultat obtenu en utilisant la bibliothèque *matplotlib* et on obtient alors la figure de droite.

On peut voir que les mots sont bien globalement placés en fonction de leur similitudes avec, d'un côté, le groupe lié à l'estazolam et de l'autre le groupe lié au paracétamol et à l'aspirine.



```

1 import matplotlib.pyplot as plt
2 import matplotlib.cm as cm
3
4 def tsne_plot_words(title, embedding_cluster, word_cluster,
5                     a, filename=None):
6     plt.figure(figsize=(4, 4))
7     colors = cm.rainbow(np.linspace(0, 1, len(word_cluster)))

```

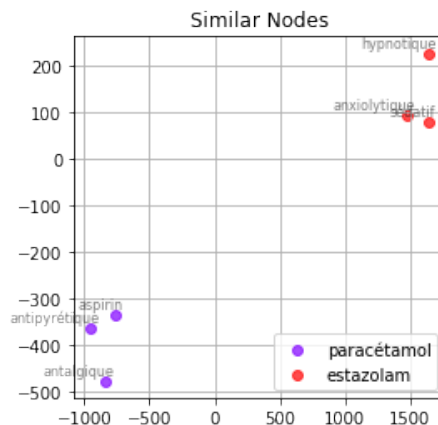
```

8
9     for i in range(len(word_cluster)):
10         x = embedding_cluster[i][0]
11         y = embedding_cluster[i][1]
12         plt.scatter(x, y, c=colors[i], alpha=a)
13         plt.annotate(word_cluster[i], alpha=0.5, xy=(x, y),
14                       xytext=(5, 2), textcoords='offset points',
15                       ha='right', va='bottom', size=8)
16
17     plt.legend(loc=4)
18     plt.title(title)
19     plt.grid(True)
20     if filename:
21         plt.savefig(filename, format='png', dpi=150,
22                     bbox_inches='tight')
23
24     plt.show()
25
26 tsne_plot_words('Nodes similarities', embeddings_en_2d,
27                 word_cluster, 0.7,
28                 'nodes_similarities.png')

```

En modifiant un peu le code précédent on peut également obtenir une représentation qui montre pour les mots qu'on a choisi, ici paracétamol et estazolam, un nuage de point correspondant aux mots avec le sens le plus proche. Cette méthode est pratique lorsqu'on a énormément de mots et qu'on veut en sélectionner une partie.

L'argument "topn=3" à la ligne 8 permet ici de dire combien de mot similaires doivent être affichés sur le graphique pour chaque mot clé.



```

1 keys = ['paracétamol', 'estazolam']
2
3 embedding_clusters = []
4 word_clusters = []
5 for word in keys:
6     embeddings = []
7     words = []
8     for similar_word, _ in model.most_similar(word, topn=3):
9         words.append(similar_word)
10        embeddings.append(model[similar_word])
11    embedding_clusters.append(embeddings)
12    word_clusters.append(words)
13
14 from sklearn.manifold import TSNE
15 import numpy as np
16

```

```

17 embedding_clusters = np.array(embedding_clusters)
18 n, m, k = embedding_clusters.shape
19 tsne_model_en_2d = TSNE(
20     perplexity=2, n_components=2, init='pca',
21     n_iter=3500, random_state=32)
22
23 embeddings_en_2d = np.array(tsne_model_en_2d.fit_transform(
24     embedding_clusters.reshape(n * m, k))
25     ).reshape(n, m, 2)
26
27 import matplotlib.pyplot as plt
28 import matplotlib.cm as cm
29
30 def tsne_plot_similar_words(title, labels, embedding_clusters,
31                             word_clusters, a, filename=None):
32     plt.figure(figsize=(4, 4))
33     colors = cm.rainbow(np.linspace(0, 1, len(labels)))
34     for label, embeddings, words, color
35         in zip(labels, embedding_clusters, word_clusters, colors):
36
37         x = embeddings[:, 0]
38         y = embeddings[:, 1]
39         plt.scatter(x, y, c=color, alpha=a, label=label)
40         for i, word in enumerate(words):
41             plt.annotate(word, alpha=0.5, xy=(x[i], y[i]),
42                         xytext=(5, 2), textcoords='offset points',
43                         ha='right', va='bottom', size=8)
44
45     plt.legend(loc=4)
46     plt.title(title)
47     plt.grid(True)
48     if filename:
49         plt.savefig(filename, format='png',
50                     dpi=150, bbox_inches='tight')
51     plt.show()
52
53 tsne_plot_similar_words('Similar Nodes', keys, embeddings_en_2d,
54                         word_clusters, 0.7, 'similar_nodes.png')

```

## 2.3 Application à notre base de données

Nous savons maintenant comment calculer les similitudes entre les noeuds de notre graphe de connaissance *Neo4J* vu au début de ce rapport. Le code ci-dessous permet d'exporter notre graphe au format *graphml* et ainsi d'utiliser la bibliothèque *Node2Vec* dessus.

```

1 CALL apoc.export.graphml.all("drug.graphml", {})
2
3 // Cette commande nécessite le plugin apoc qui peut
4 // être installer facilement depuis l'onglet Plugins
5 // de l'interface Neo4J Desktop

```

## Améliorations

L'important maintenant est de mettre dans la base de données des informations qui permettront ensuite à *Node2Vec* de déterminer les nouvelles applications des médicaments.

Actuellement la base de données ne permet pas vraiment de faire cela mais on observe déjà que le modèle reconnaît les molécules synonymes.

Une piste intéressante serait d'ajouter dans la base de données les applications des médicaments ainsi que les effets secondaires et toute autres propriétés qui pourraient permettre d'effectuer le repositionnement en observant directement les liens entre les médicaments, les maladies et les effets.

On peut également envisager de relier un médicament à des mots clés eux-même reliés entre eux. Pour cela on peut par exemple, à partir du champ *description* ne garder que les mots les plus rares dans la langue anglaise, ce qui enlèvera les mots comme "the", "a", "an" et conservera les mots comme "antalgic". Ces mots seraient donc reliés au médicament décrit puis à l'aide des modèles de plongements de mots on pourrait également relier ces mots clés entre eux. De plus le fait de relier les mots clés entre eux ne nécessite pas forcément de partir de données médicales, n'importe quel texte peut être donné à un modèle de plongement de mot pour reconnaître les mots similaires.

Le fait d'avoir un nuage de mot clé reliés entre eux et aux médicaments associés serait très bénéfique pour *Node2Vec* qui pourra alors facilement retrouver les médicaments similaires.

## Conclusion

Ce projet nous a permis de mettre en pratique les méthodes de plongement de mots que nous avons vu en cours. Cependant étant donné que le sujet global du projet était difficile à appréhender et pas instinctivement corrélé à l'intelligence artificielle nous n'avons pas encore obtenu de résultats montrant qu'il est possible de faire du repositionnement de médicaments de façon automatisé.

Malgré tout, les résultats obtenus à la fin du projet vis à vis du traitement des données de *Drugbank* et de leur mise en place dans un graphe de connaissance sont encourageants et laisse à penser qu'il est effectivement possible d'utiliser les algorithmes de *Graph embedding* afin de trouver les repositionnements.