


```
def abc_score(seed_set, test_set, co_occur_mat):

    scores = {}

    for Ti in test:
        score_Ti = 0
        ranks = {}
        if Ti in co_occur_mat.columns.tolist():
            co_occur_mat_sorted_by_Ti = co_occur_mat.sort_values(by=[Ti], ascending=False)
            co_occur_mat_sorted_by_Ti = co_occur_mat_sorted_by_Ti.loc[co_occur_mat_sorted_by_Ti["drug_gene"] != ""]
            n = len(co_occur_mat_sorted_by_Ti["drug_gene"])
            j = 1
            for drug_gene in co_occur_mat_sorted_by_Ti["drug_gene"]:
                ranks[drug_gene] = j
                j+=1
            for Rij in co_occur_mat_sorted_by_Ti["drug_gene"]:
                if Rij in seed_set:
                    score_Ti += ranks[Rij]
            scores[Ti] = score_Ti

    return scores
```

Assess the quality of unsupervised clustering

The original paper generates randomly 1000 seed sets and test sets, resulting in 1000 different rankings. The rankings' quality is determined with the method below.

For each verified drug-gene pair ranking, we're supposed to determine if the ranking is true or false. Using this information, we can generate an ROC curve with the true-positive-rate on the y-axis and false-positive-rate on the x-axis. The AUC score, or the area under the ROC curve, provides us a measure of how likely it is that a positive element of the test set will be ranked higher than a negative element. The original paper considers any rankings with an AUC > 0.7 acceptable.

Notes:

- We are uncertain on how true or false is assigned to each ranking.
- The AUC value of 0.7 was chosen by the original author after observing the frequency of correct rankings based on their AUC value. The information can be found on page 125 in Dr. Percha's dissertation titled "Biomedical Text Mining from Context".

5. Scaling to Dask

In order to run all .xml files, we use Dask to scale up our operations. All core functions are kept the same. The main difference between our sequential computation and parallel computation is a conversion from .xml to .csv, and the mapping of the functions to the dask dataframe.

Convert all .xml files to .csv files

This step is necessary because Dask Dataframe requires a .csv input.

```
def convert_xml_to_csv(xml_file_directory):
    csv_file = pd.DataFrame.from_dict(pp.parse_medline_xml(xml_file))
    csv_file.to_csv(path_or_buf = directoryre.sub(r'.xml',r'.csv'),xml_file[len(directory):-3]), index=False)
    pass

def convert_all_xml_to_csv(directory='all_xml/'):

    """
    Convert xml to csv
    Change the range to the length of all the files listed to get the full range
    directory: folder where xml are located
    """

    file_list = os.listdir('all_xml/')
    pubmed_paths = [directory+pubmed2in+str(num).zfill(4)+'.xml.gz' for num in range(1,len(file_list))]
    for xml_file in pubmed_paths:
        convert_xml_to_csv(xml_file,directory)
```

Map all functions to Dask Dataframe

The mapping is applied from the beginning up to the obtain dependency matrix step. After that, parallel computing is not necessary. However, with the many loops in the rest of the code, we could parallel all the steps after the dependency matrix with more time.

```
import dask
import dask.dataframe as dd
import spacy
import pandas as pd
import regex as re
import pubmed_parser as pp
import numpy as np
from all_functions import *
import os

#read in all csv as dask dataframe
df = dd.read_csv('all_xml/pubmed1*.csv', sample=2500000,
                dtype={'title':'object','abstract':'object','journal':'object',
                        'authors':'object','pubdate':'object','pmid':'object',
                        'mesh_terms':'object','publication_types':'object','chemical_list':'object',
                        'keywords':'object','doi':'object','references':'object','delete':'object',
                        'affiliations':'object','pmc':'object','other_id':'object','medline_ta':'object',
                        'nlm_unique_id':'object','issn_linking':'object','country':'object'})

#take only the abstract column
abs_df = df[['abstract']]

#map filtering for all sentences containing drug-gene pair function
abs_filt = abs_df.map_partitions(lambda x: get_filtered_data_all_sent(x, drugs, genes))

#map obtaining matrix from sentences and drug-gene pairs
dependency_matrix = abs_filt.map_partitions(lambda x: get_dependency_matrix(x))

#compute all map function
dependency_matrix_pd = dependency_matrix.compute()
```

6. Creating Dendrogram with the Cluster Assignments

The dendrogram contains all the drug-gene pairs and how they are related. Imagine a family tree with each branch describing a relationship between family members. In our context, all the family members are drug-gene pairs and the relationships are within biomedical contexts (activator, inhibitor, etc). With the cluster assignments, the "closeness" of the pairs are determined through the Spearman correlation.

After this step, we should be able to see which pairs have already been discovered, and which ones are completely new. The new drug-gene pairs are potential testing targets in the lab.

The code for this portion is in R.

7. Final Notes

Final bottleneck

While the entire pipeline works well, with both our sequential and parallel computing fully functional, we weren't able to process entirely 777 different xml files due to our limited hardware. Our personal computers ran out of memory after a few days of processing. Our project would be finished completely with either better hardware or more time to clean up the data structure of the code.

Future improvements

Our project focused mainly on replicating Dr. Percha's paper. However, we have some ideas to improve the current process.

1. Use the new Stanford Parser in the form of the stanza package.
2. With the stanza package, we can use the AWS server to compute all the xml files, eliminating the final bottleneck.
3. Use all updated xml files.