# Project 2: Multi-source code search

Ardigò Susanna

December 11, 2020

# Contents

# 1    Data Extraction

## 1.1    Goal and Input parameter

This part of the project consists of extracting names and comments of Python classes, methods and functions and save them in a csv file.

This file takes as argument the path of the directory of the project that we want to analyze. For this project we use the project `tensorflow`.

## 1.2    Description of the code

To efficiently parse the files in the directory, we created a class named `Visitor`, which extends the `NodeVisitor` class of the standard library ast (which stands for Abstract Syntax Tree). This class holds the path of the file. There is a global variable `data` used throughout the execution to store all the information extracted.

The function `start(directory_path)` *'walks'* the given directory using the function `walk` which generates a 3-tuple of directory path, directory names and file names. We open and read all the python files, checked with the extension of the file, we create a Visitor object and start to visit. The class we created has two different visit methods which differ in if the node visiting is a definition of a class or a function. The method `visit_FunctionDef(self, node:  FunctionDef)` adds the node information to the array of data if the function or method is not a main or a test. Since this method is used both for functions and methods, we know if the node is a method by checking if the first argument is `self`. The method `visit_ClassDef(self, node:  ClassDef)` calls a generic visit (of the ast library) and, as the previous method, adds the node information to the array of data if the class is not a main or a test.

After the parsing is complete we create a pandas dataframe, feeding it as data the data array, and export it in a csv extension.

## 1.3    Results

Table 1 show the number of Python files, classes, methods and functions found while parsing the Tensorflow directory. The results can be found in the file `res/data.csv`.

| Type | # |
| --- | --- |
| Python files | 2817 |
| Classes | 1904 |
| Methods | 7271 |
| Functions | 4881 |

Table 1: Count of data found in Tensorflow

# 2 Training of search engines

## 2.1 Goal and Input parameter

This part of the project consists of representing code entities using the four embeddings frequency, TF-IDF, LSI and Doc2Vec.
This file takes as argument a query which will be fed to the four search engines.

## 2.2 Description of the code

The function `start(query)` loads the csv into a pandas dataframe and then computes the results. The first part of function `compute_results(query, dataframe)` creates the necessary data and normalizes the query that the second part needs to produce the results.

The first part of function `create_data(dataframe)` extracts the names and comments of the data extracted in the first part. to create a clean array of arrays of tokens and a dictionary with the frequencies of each token. In the second part we create the corpus by processing the tokens, we create a gensim dictionary and the bag of words. At the end of the creation, we save the corpus, dictionary and bag of words in external files to then load them in future runs. In the second part of function `compute_results(query, dataframe)` we create a dictionary that hold the results of the searches and a dictionary to save the embedding vectors.

The function `query_frequency(query, bow, dictionary)` creates a sparse matrix of the bag of words and returns an array with the similarity scores of each entity of the given csv file. This array is then filtered to extract only the top 5 scoring entities. Similarly, the function `query_tfidf(query, bow, dictionary)` creates a sparse matrix of the tfidf model of the bag of words and returns an array with the similarity scores which is then filtered. The function `query_lsi(query, bow, dictionary)` creates a lsi model based on the bag of words, a vector based on the model and the dictionary, the matrix of the similarities and the embedding vectors. The result of the matrix, as in the previous cases, is filtered to get only the top 5 scores. The function `query_doc2vec(query, bow, dictionary)` creates a doc2vec model which then feed the corpus to and train it. We create a vector infering it from the query, we create the similarity and take only the top 5 scores and the embedding vectors.

We save the trained models in external pickle files to load then load them in the next runs. This improves the running time of the function.

We create a dataframe with the information stored in the dictionary, we print the results and save them in a separate file.

## 2.3 Results

To show the results we run this part of the project with the query:

        *'AST Visitor that looks for specific API usage without editing anything'*        .

The correct document is `PastaAnalyzeVisitor` with path `../tensorflow/tensorflow/tools/compatibility/ast_edits.py`.

Figure 1 show the result of the given query.

As we can see in the image all search engine find the correct result. Frequency, TF-IDF and LSI find it as the first result. Doc2Vec finds the correct result as second result. The first result is a class in the same file but named `APIChangeSpec` and with a different comment. This result is found only by this search engine. Doc2Vec has different results compared to the other search engines. Frequency and TF-IDF have the most similar results.

| | name | file | line | type | comment | search |
|---|---|---|---|---|---|---|
| 1 | PastaAnalyzeVisitor | ../tensorflow/tensorflow/tools/compatibility/ast_edits.py | 815 | class | AST Visitor that looks for specific API usage without editing an… | FREQ |
| 2 | CompatV1ImportRepla… | ../tensorflow/tensorflow/tools/compatibility/tf_upgrade_v2.py | 72 | class | AST Visitor that replaces `import tensorflow.compat.v1 as tf`. | FREQ |
| 3 | CleanCopier | ../tensorflow/tensorflow/python/autograph/pyct/ast_util.py | 30 | class | NodeTransformer-like visitor that copies an AST. | FREQ |
| 4 | FunctionVisitor | ../tensorflow/tensorflow/python/autograph/pyct/static_analysis/type_infer… | 394 | class | AST visitor that applies type inference to each function separat… | FREQ |
| 5 | track_usage | ../tensorflow/tensorflow/python/platform/analytics.py | 21 | function | No usage tracking for external library. | FREQ |
| 6 | PastaAnalyzeVisitor | ../tensorflow/tensorflow/tools/compatibility/ast_edits.py | 815 | class | AST Visitor that looks for specific API usage without editing an… | TF-IDF |
| 7 | PublicAPIVisitor | ../tensorflow/tensorflow/tools/common/public_api.py | 29 | class | Visitor to use with `traverse` to visit exactly the public TF AP… | TF-IDF |
| 8 | CleanCopier | ../tensorflow/tensorflow/python/autograph/pyct/ast_util.py | 30 | class | NodeTransformer-like visitor that copies an AST. | TF-IDF |
| 9 | FunctionVisitor | ../tensorflow/tensorflow/python/autograph/pyct/static_analysis/type_infer… | 394 | class | AST visitor that applies type inference to each function separat… | TF-IDF |
| 10 | track_usage | ../tensorflow/tensorflow/python/platform/analytics.py | 21 | function | No usage tracking for external library. | TF-IDF |
| 11 | PastaAnalyzeVisitor | ../tensorflow/tensorflow/tools/compatibility/ast_edits.py | 815 | class | AST Visitor that looks for specific API usage without editing an… | LSI |
| 12 | CountingVisitor | ../tensorflow/tensorflow/python/autograph/pyct/cfg_test.py | 28 | class | <null> | LSI |
| 13 | matches | ../tensorflow/tensorflow/python/autograph/pyct/ast_util.py | 214 | function | Basic pattern matcher for AST. | LSI |
| 14 | concentration | ../tensorflow/tensorflow/python/ops/distributions/dirichlet.py | 213 | method | Concentration parameter; expected counts for that coordinate. | LSI |
| 15 | ProfileOptionBuilder | ../tensorflow/tensorflow/python/profiler/option_builder.py | 27 | class | Option Builder for Profiling API. | LSI |
| 16 | APIChangeSpec | ../tensorflow/tensorflow/tools/compatibility/ast_edits.py | 192 | class | This class defines the transformations that need to happen. | Doc2Vec |
| 17 | PastaAnalyzeVisitor | ../tensorflow/tensorflow/tools/compatibility/ast_edits.py | 815 | class | AST Visitor that looks for specific API usage without editing an… | Doc2Vec |
| 18 | Trackable | ../tensorflow/tensorflow/python/training/tracking/base.py | 537 | class | Base class for `Trackable` objects without automatic dependencie… | Doc2Vec |
| 19 | check_dtype | ../tensorflow/tensorflow/python/ops/linalg/linear_operator_util.py | 139 | function | Check that arg.dtype == self.dtype. | Doc2Vec |
| 20 | ProfileOptionBuilder | ../tensorflow/tensorflow/python/profiler/option_builder.py | 27 | class | Option Builder for Profiling API. | Doc2Vec |

Figure 1: Results of the given query

# 3  Evaluation of search engines

## 3.1  Goal and Input parameter

This part of the project consists of measuring the precision and recall given 10 queries along with their ground truth.

This file takes as argument the path of the ground truth file.

## 3.2  Description of the code

The function `start(path_ground_truth)` loads the csv of the data into a pandas dataframe, parses the ground truth and then computes the precision and recall.

To efficiently parse the ground truth file, we created a class named `Truth` which holds the name, path and query. We read the ground truth file and create an array with all the entries of the ground truth and the queries.

To compute precision and recall we get the data of the results and the embedding vectors from the previous part. We create a dictionary to save the scores of the queries and a dictionary for the vectors. We then compute the precision and recall, by comparing our results and the ground truth.

## 3.3   Results

Table 2 show the statistics of precision and recall compared to the unique ground truth. We can see that the precision is low for all engines. The engine with the highest precision is `LSI`, which is the only score higher than 0.4. The second highest precision is `TF-IDF`, followed by `Frequencies` and then **Doc2Vec**. We can say that almost all search engine have a similar precision.

The recall is higher than the precision. The `TF-IDF` engine has a recall equal to 1, which means that for each query the search engine has found the correct result in the top 5. The second highest recall is of `Frequencies` with score 0.9. Both `LSI` and `Doc2Vec` have a recall of 0.8.

| Engine | Precision | Recall |
|--------|-----------|--------|
| Frequencies | 0.332 | 0.9 |
| TD-IDF | 0.365 | 1.0 |
| LSI | 0.403 | 0.8 |
| Doc2Vec | 0.323 | 0.8 |

Table 2: Statistics of the search engines

# 4   Visualisation of query results

## 4.1   Goal and Input parameter

This part of the project consists of visualizing the embedding vectors of the queries and the top 5 answers in a 2D plot. This file takes as argument the ground truth file.

## 4.2   Description of the code

The first part of the execution is the same as the previous file. After the results are calculated, we use the embedding vectors, that we retrieved in the explanation above but we did not use. For vector we apply TSNE to produce 2D vectors composed of queries and the top 4 results. The plot is straight-forward: we create a dataframe with the information of x and y coordinates and print them of different hues. We use the library `seaborn` to create the charts and we then save them on disk.
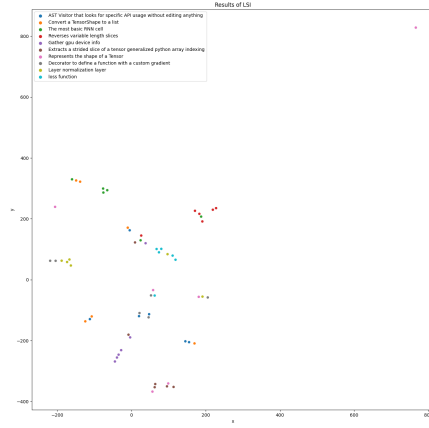
## 4.3   Results

Figure 2 shows the plots of the visualization of the queries. At first we notice that the LSI scatterplot tends to be more a bit compact, while the Doc2Vec scatterplot is more sparse. The optimal solution is to have defined clusters for each query. This does not happen in any of the two images.
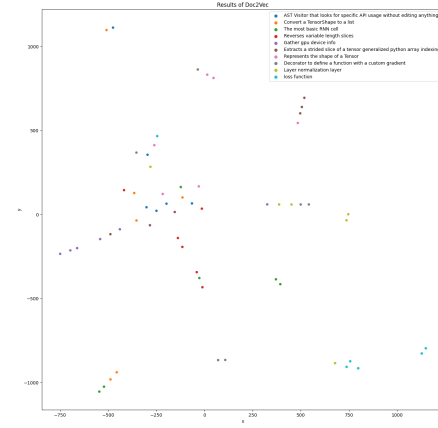
### 4.3.1 LSI

Analyzing the plot of LSI, shown in figure 2a, we can see that some results of the queries tend to stay close, but not completely. There are some clusters that are close. In most cases the cluster have different queries. The most well defined cluster is *'Gather gpu device info'*, colored in purple. The sparseness of the data reflects the low precision.

### 4.3.2 Doc2Vec

Analyzing the plot of Doc2Vec, shown in figure 2b, we can see that few results of the queries tend to stay close. In this image there are some cluster with more queries. The most well defined cluster is *'Gather gpu device info'*, as in the previous plot. The sparseness of the data reflects the low precision, which is lower than the LSI. This is reflected in the plot.



(a) Results of LSI                                      (b) Results of Doc2Vec

Figure 2: Visualization of the plots of the queries

# A  Python code

## A.1  Data Extraction

```python
from sys import argv, exit
from ast import *
from os import walk
import pandas as pd


class Visitor(NodeVisitor):
    def __init__(self, file_path, node):
        super().__init__()
        self.file_path = clean_file_path(file_path)
        self.visit(parse(node))

    def visit_ClassDef(self, node: ClassDef):
        self.generic_visit(node)
        if is_valid_entity(node.name):
            self.append_data(node, "class")

    def visit_FunctionDef(self, node: FunctionDef):
        if is_valid_entity(node.name):
            self.append_data(node, "method" if is_method(node) else "function")

    def append_data(self, node, def_type):
        comment = get_docstring(node)
        comment = comment.split('\n')[0] if comment is not None else ""
        data.append((node.name, self.file_path, node.lineno, def_type, comment))


def clean_file_path(path):
    directories = path.split('/')
    return '../' + '/'.join(directories[directories.index('tensorflow'):])


def is_valid_entity(name):
    return name[0] != '_' and name != "main" and "test" not in name.lower()


def is_method(function):
    return function.args and len(function.args.args) > 0 and 'self' in function.args.args[0].arg


def start(directory_path):
    if directory_path[-1] == '/':
        directory_path = directory_path[: -1]
    counter = 0
    for path, _, files in walk(directory_path):
        for file_name in files:
            if file_name.endswith('.py'):
                counter += 1
                file_path = path + '/' + file_name
                with open(file_path) as file:
                    Visitor(file_path, file.read())
```

```
52
53        dataframe = pd.DataFrame(data=data, columns=["name", "file", "line", "type", "comment"])
54        dataframe.to_csv('res/data.csv', index=False, encoding='utf-8')
55        print("files\t    " + str(counter))
56        print(dataframe["type"].value_counts())
57
58
59    if len(argv) < 2:
60        print("Please give as input the path of the directory to analyze")
61        exit(1)
62    data = []
63    start(argv[1])
```

## A.2   Training of search engines

```
1     from datetime import datetime
2     import string
3     from os import path
4     import pandas as pd
5     import pickle as pkl
6     from re import finditer
7     from sys import argv, exit
8     from collections import defaultdict
9     from gensim.corpora import Dictionary
10    from gensim.models.doc2vec import TaggedDocument
11    from gensim.utils import simple_preprocess
12    from gensim.models import TfidfModel, LsiModel, Doc2Vec
13    from gensim.similarities import MatrixSimilarity, SparseMatrixSimilarity
14
15
16    def start(query):
17        dataframe = pd.read_csv("res/data.csv").fillna(value="")
18        results_dictionary, _ = compute_results(query, dataframe)
19        results = pd.DataFrame(data=create_result_dataframe(results_dictionary, dataframe),
20                               columns=['name', "file", "line", "type", "comment", "search"])
21        pd.options.display.max_colwidth = 200
22        print_results(results)
23        results.to_latex('res/search_data.tex', index=False, encoding='utf-8')
24        results.to_csv('res/search_data.csv', index=False, encoding='utf-8')
25
26
27    def compute_results(query, dataframe):
28        processed_corpus, frequencies, bag_of_words = get_data(dataframe)
29        query_to_execute = normalize_query(query)
30        results = {
31            "FREQ": query_frequency(query_to_execute, bag_of_words, frequencies),
32            "TF-IDF": query_tfidf(query_to_execute, bag_of_words, frequencies)
33        }
34        vectors = dict()
35        results["LSI"], vectors["LSI"] = query_lsi(query_to_execute, bag_of_words, frequencies)
36        results["Doc2Vec"], vectors["Doc2Vec"] = query_doc2vec(query_to_execute, processed_corpus
              )
37        return results, vectors
38
39
```

8

```
40  def get_data(df):
41      return load_data_files() if exists_data_files() else create_data(df)
42
43
44  def create_data(df):
45      tokens = [filter_stopwords(normalize_tokens(handle_camel_case(split_underscore(
46          [row["name"]] + split_space(row["comment"]))))) for _, row in df.iterrows()]
47
48      frequency = defaultdict(int)
49      for token in tokens:
50          for word in token:
51              frequency[word] += 1
52
53      corpus = [[token for token in text if frequency[token] > 1] for text in tokens]
54      dictionary = Dictionary(corpus)
55      bow = [dictionary.doc2bow(text) for text in corpus]
56
57      save_data(corpus, 'corpus')
58      save_data(dictionary, 'dictionary')
59      save_data(bow, 'bow')
60      return corpus, dictionary, bow
61
62
63  def exists_data_files():
64      return exists_file('corpus') and exists_file('dictionary') and exists_file('bow')
65
66
67  def exists_file(name):
68      return path.exists('res/' + name + '.pkl')
69
70
71  def load_data_files():
72      return load_file('corpus'), load_file('dictionary'), load_file('bow')
73
74
75  def save_data(data, name):
76      pkl.dump(data, open('res/' + name + '.pkl', "wb"), protocol=pkl.HIGHEST_PROTOCOL)
77
78
79  def load_file(name):
80      return pkl.load(open('res/' + name + '.pkl', "rb"))
81
82
83  def split_space(text):
84      return text.translate(str.maketrans('', '', string.punctuation)).split(' ') if text != ""
               else []
85
86
87  def split_underscore(tokens):
88      return [word for token in tokens for word in token.split('_')]
89
90
91  def handle_camel_case(tokens):
92      words = []
93      for token in tokens:
94          matches = finditer('.+?(?:(?<=[a-z])(?=[A-Z])|(?<=[A-Z])(?=[A-Z][a-z])|$)', token)
95          words += [m.group(0) for m in matches]
```

9

```
96          return words
97
98
99   def normalize_tokens(tokens):
100         return [token.lower() for token in tokens]
101
102
103  def filter_stopwords(tokens):
104         for token in tokens:
105             if token in ['test', 'tests', 'main']:
106                 return []
107         return tokens
108
109
110  def normalize_query(query):
111         return query.strip().lower().split()
112
113
114  def save_model(model, name):
115         save_data(model, 'model_' + name)
116
117
118  def query_frequency(query, bow, dictionary):
119         return filter_results(get_freq_model(bow, dictionary)[dictionary.doc2bow(query)])
120
121
122  def get_freq_model(bow, dictionary):
123         return load_file('model_freq') if exists_file('model_freq') else create_freq_model(bow,
                 dictionary)
124
125
126  def create_freq_model(bow, dictionary):
127         model = SparseMatrixSimilarity(bow, num_features=len(dictionary.token2id))
128         save_model(model, 'freq')
129         return model
130
131
132  def query_tfidf(query, bow, dictionary):
133         model = get_tfidf_model(bow)
134         matrix = get_tfidf_matrix(model, bow, dictionary)
135         return filter_results(matrix[model[dictionary.doc2bow(query)]])
136
137
138  def get_tfidf_model(bow):
139         return load_file('model_tfidf') if exists_file('model_tfidf') else create_tfidf_model(bow
                 )
140
141
142  def create_tfidf_model(bow):
143         model = TfidfModel(bow)
144         save_model(model, 'tfidf')
145         return model
146
147
148  def get_tfidf_matrix(model, bow, dictionary):
149         return load_file('matrix_tfidf') if exists_file('matrix_tfidf') else create_tfidf_matrix(
                 model, bow, dictionary)
```

```
150
151
152 def create_tfidf_matrix(model, bow, dictionary):
153     matrix = SparseMatrixSimilarity(model[bow], num_features=len(dictionary.token2id))
154     save_data(matrix, 'matrix_tfidf')
155     return model
156
157
158 def query_lsi(query, bow, dictionary):
159     model = get_lsi_model(bow, dictionary)
160     vector = model[dictionary.doc2bow(query)]
161     result = abs(MatrixSimilarity(model[bow])[vector])
162     embedding = [[value for _, value in vector]] + [[value for _, value in model[bow][i]] for
                i, value in
163                                                     sorted(enumerate(result), key=lambda x: x
                                                         [1], reverse=True)[:5]]
164     return filter_results(result), embedding
165
166
167 def get_lsi_model(bow, dictionary):
168     return load_file('model_lsi') if exists_file('model_lsi') else create_lsi_model(bow,
            dictionary)
169
170
171 def create_lsi_model(bow, dictionary):
172     model = LsiModel(bow, id2word=dictionary, num_topics=300)
173     save_model(model, 'lsi')
174     return model
175
176
177 def filter_results(arrg):
178     return [i for i, v in sorted(enumerate(arrg), key=lambda x: x[1], reverse=True)[:5]]
179
180
181 def query_doc2vec(query, corpus):
182     model = get_doc2vec_model(get_doc2vec_corpus(corpus))
183     vector = model.infer_vector(query)
184     similar = model.docvecs.most_similar([vector], topn=5)
185     return [index for (index, _) in similar], \
186             [list(vector)] + [list(model.infer_vector(corpus[index])) for index, _ in similar]
187
188
189 def get_doc2vec_corpus(corpus):
190     return [TaggedDocument(simple_preprocess(' '.join(element)), [index])
191             for index, element in enumerate(corpus)]
192
193
194 def get_doc2vec_model(corpus):
195     return load_file('model_doc2vec') if exists_file('model_doc2vec') else
            create_doc2vec_model(corpus)
196
197
198 def create_doc2vec_model(corpus):
199     model = Doc2Vec(vector_size=300, min_count=2, epochs=77)
200     model.build_vocab(corpus)
201     model.train(corpus, total_examples=model.corpus_count, epochs=model.epochs)
202     save_model(model, 'doc2vec')
```

11

```
203        return model
204
205
206  def create_result_dataframe(queries_dictionary, df):
207      for key, values in queries_dictionary.items():
208          for index in sorted(values):
209              row = df.iloc[index]
210              yield [row["name"], row["file"], row["line"], row["type"], row["comment"], key]
211
212
213  def print_results(df):
214      grouped = df.groupby(['search'])
215      for key, item in grouped:
216          print(grouped.get_group(key), "\n\n")
217
218
219  if len(argv) < 2:
220      print("Please give as input the query")
221      exit(1)
222
223  start(argv[1])
```

## A.3   Evaluation of search engines and Visualisation of query results

```
1  import itertools
2  from datetime import datetime
3
4  import string
5  import pandas as pd
6  from os import path
7  import pickle as pkl
8  import seaborn as sns
9  from re import finditer
10  from sys import argv, exit
11  import matplotlib.pyplot as plt
12  from sklearn.manifold import TSNE
13  from collections import defaultdict
14  from gensim.corpora import Dictionary
15  from gensim.models.doc2vec import TaggedDocument
16  from gensim.utils import simple_preprocess
17  from gensim.models import TfidfModel, LsiModel, Doc2Vec
18  from gensim.similarities import MatrixSimilarity, SparseMatrixSimilarity
19
20  ##################
21  def get_results(query, dataframe):
22      results_dictionary, vectors = compute_results(query, dataframe)
23      return pd.DataFrame(data=create_result_dataframe(results_dictionary, dataframe),
24                          columns=['name', "file", "line", "type", "comment", "search"]),
                                 vectors
25
26
27  def compute_results(query, dataframe):
28      processed_corpus, frequencies, bag_of_words = get_data(dataframe)
29      query_to_execute = normalize_query(query)
30      results = {
```

```
31              "FREQ": query_frequency(query_to_execute, bag_of_words, frequencies),
32              "TF-IDF": query_tfidf(query_to_execute, bag_of_words, frequencies)
33          }
34      vectors = dict()
35      results["LSI"], vectors["LSI"] = query_lsi(query_to_execute, bag_of_words, frequencies)
36      results["Doc2Vec"], vectors["Doc2Vec"] = query_doc2vec(query_to_execute, processed_corpus
            )
37      return results, vectors
38
39
40  def get_data(df):
41      return load_data_files() if exists_data_files() else create_data(df)
42
43
44  def create_data(df):
45      tokens = [filter_stopwords(normalize_tokens(handle_camel_case(split_underscore(
46          [row["name"]] + split_space(row["comment"]))))) for _, row in df.iterrows()]
47
48      frequency = defaultdict(int)
49      for token in tokens:
50          for word in token:
51              frequency[word] += 1
52
53      corpus = [[token for token in text if frequency[token] > 1] for text in tokens]
54      dictionary = Dictionary(corpus)
55      bow = [dictionary.doc2bow(text) for text in corpus]
56
57      save_data(corpus, 'corpus')
58      save_data(dictionary, 'dictionary')
59      save_data(bow, 'bow')
60      return corpus, dictionary, bow
61
62
63  def exists_data_files():
64      return exists_file('corpus') and exists_file('dictionary') and exists_file('bow')
65
66
67  def exists_file(name):
68      return path.exists('res/' + name + '.pkl')
69
70
71  def load_data_files():
72      return load_file('corpus'), load_file('dictionary'), load_file('bow')
73
74
75  def save_data(data, name):
76      pkl.dump(data, open('res/' + name + '.pkl', "wb"), protocol=pkl.HIGHEST_PROTOCOL)
77
78
79  def load_file(name):
80      return pkl.load(open('res/' + name + '.pkl', "rb"))
81
82
83  def split_space(text):
84      return text.translate(str.maketrans('', '', string.punctuation)).split(' ') if text != ""
            else []
85
```

13

```
86
87  def split_underscore(tokens):
88      return [word for token in tokens for word in token.split('_')]
89
90
91  def handle_camel_case(tokens):
92      words = []
93      for token in tokens:
94          matches = finditer('.+?(?:(?<=[a-z])(?=[A-Z])|(?<=[A-Z])(?=[A-Z][a-z])|$)', token)
95          words += [m.group(0) for m in matches]
96      return words
97
98
99  def normalize_tokens(tokens):
100     return [token.lower() for token in tokens]
101
102
103 def filter_stopwords(tokens):
104     for token in tokens:
105         if token in ['test', 'tests', 'main']:
106             return []
107     return tokens
108
109
110 def normalize_query(query):
111     return query.strip().lower().split()
112
113
114 def save_model(model, name):
115     save_data(model, 'model_' + name)
116
117
118 def query_frequency(query, bow, dictionary):
119     return filter_results(get_freq_model(bow, dictionary)[dictionary.doc2bow(query)])
120
121
122 def get_freq_model(bow, dictionary):
123     return load_file('model_freq') if exists_file('model_freq') else create_freq_model(bow,
             dictionary)
124
125
126 def create_freq_model(bow, dictionary):
127     model = SparseMatrixSimilarity(bow, num_features=len(dictionary.token2id))
128     save_model(model, 'freq')
129     return model
130
131
132 def query_tfidf(query, bow, dictionary):
133     model = get_tfidf_model(bow)
134     matrix = get_tfidf_matrix(model, bow, dictionary)
135     return filter_results(matrix[model[dictionary.doc2bow(query)]])
136
137
138 def get_tfidf_model(bow):
139     return load_file('model_tfidf') if exists_file('model_tfidf') else create_tfidf_model(bow
             )
140
```

14

```
141
142   def create_tfidf_model(bow):
143       model = TfidfModel(bow)
144       save_model(model, 'tfidf')
145       return model
146
147
148   def get_tfidf_matrix(model, bow, dictionary):
149       return load_file('matrix_tfidf') if exists_file('matrix_tfidf') else create_tfidf_matrix(
                 model, bow, dictionary)
150
151
152   def create_tfidf_matrix(model, bow, dictionary):
153       matrix = SparseMatrixSimilarity(model[bow], num_features=len(dictionary.token2id))
154       save_data(matrix, 'matrix_tfidf')
155       return model
156
157
158   def query_lsi(query, bow, dictionary):
159       model = get_lsi_model(bow, dictionary)
160       vector = model[dictionary.doc2bow(query)]
161       result = abs(MatrixSimilarity(model[bow])[vector])
162       embedding = [[value for _, value in vector]] + [[value for _, value in model[bow][i]] for
                 i, value in
163                                                   sorted(enumerate(result), key=lambda x: x
                                                       [1], reverse=True)[:5]]
164       return filter_results(result), embedding
165
166
167   def get_lsi_model(bow, dictionary):
168       return load_file('model_lsi') if exists_file('model_lsi') else create_lsi_model(bow,
                 dictionary)
169
170
171   def create_lsi_model(bow, dictionary):
172       model = LsiModel(bow, id2word=dictionary, num_topics=300)
173       save_model(model, 'lsi')
174       return model
175
176
177   def filter_results(arrg):
178       return [i for i, v in sorted(enumerate(arrg), key=lambda x: x[1], reverse=True)[:5]]
179
180
181   def query_doc2vec(query, corpus):
182       model = get_doc2vec_model(get_doc2vec_corpus(corpus))
183       vector = model.infer_vector(query)
184       similar = model.docvecs.most_similar([vector], topn=5)
185       return [index for (index, _) in similar], \
186               [list(vector)] + [list(model.infer_vector(corpus[index])) for index, _ in similar]
187
188
189   def get_doc2vec_corpus(corpus):
190       return [TaggedDocument(simple_preprocess(' '.join(element)), [index])
191               for index, element in enumerate(corpus)]
192
193
```

```
194  def get_doc2vec_model(corpus):
195      return load_file('model_doc2vec') if exists_file('model_doc2vec') else
             create_doc2vec_model(corpus)
196
197
198  def create_doc2vec_model(corpus):
199      model = Doc2Vec(vector_size=300, min_count=2, epochs=77)
200      model.build_vocab(corpus)
201      model.train(corpus, total_examples=model.corpus_count, epochs=model.epochs)
202      save_model(model, 'doc2vec')
203      return model
204
205
206  def create_result_dataframe(queries_dictionary, df):
207      for key, values in queries_dictionary.items():
208          for index in sorted(values):
209              row = df.iloc[index]
210              yield [row["name"], row["file"], row["line"], row["type"], row["comment"], key]
211
212
213  #####################################
214
215  class Truth:
216      def __init__(self, query, name, path):
217          self.name = name
218          self.path = path
219          self.query = query.lower()
220
221
222  class Stat:
223      def __init__(self, precisions, recalls):
224          self.precisions = precisions
225          self.recalls = recalls
226
227
228  def start(path_ground_truth):
229      dataframe = pd.read_csv("res/data.csv").fillna(value="")
230      ground_truth, queries = parse_ground_truth(path_ground_truth)
231      scores, vectors = compute_precision_recall(ground_truth, dataframe)
232      plot_vectors(compute_tsne(vectors), queries)
233      print_scores(scores)
234
235
236  def parse_ground_truth(path_ground_truth):
237      classes, queries = [], []
238      for entry in open(path_ground_truth, "r").read().split("\n\n"):
239          data = entry.split("\n")
240          classes.append(Truth(data[0], data[1], data[2]))
241          queries.append(data[0])
242      return classes, queries
243
244
245  def compute_precision_recall(ground_truth, dataframe):
246      scores = {"FREQ": [], "TF-IDF": [], "LSI": [], "Doc2Vec": []}
247      vectors = {"LSI": [], "Doc2Vec": []}
248      for entry in ground_truth:
249          results, vectors_i = get_results(entry.query, dataframe)
```

```
250            vectors["LSI"] += vectors_i["LSI"]
251            vectors["Doc2Vec"] += vectors_i["Doc2Vec"]
252            for query_type in ["FREQ", "TF-IDF", "LSI", "Doc2Vec"]:
253                precision = compute_precision(entry, query_type, results)
254                scores[query_type].append(Stat(precision, compute_recall(precision)))
255        return scores, vectors
256
257
258    def compute_precision(truth, search_type, dataframe):
259        counter = 0
260        for _, row in dataframe[dataframe['search'] == search_type].iterrows():
261            counter += 1
262            if row["name"] == truth.name and row["file"] == truth.path:
263                return 1 / counter
264        return 0
265
266
267    def compute_recall(precision):
268        return 1 if precision > 0 else 0
269
270
271    def compute_tsne(dictionary):
272        results = {}
273        for key, values in dictionary.items():
274            tsne = TSNE(n_components=2, verbose=1, perplexity=2, n_iter=3000)
275            results[key] = tsne.fit_transform(values)
276        return results
277
278
279    def plot_vectors(dictionary, queries):
280        for key, values in dictionary.items():
281            dataframe = pd.DataFrame()
282            dataframe['x'] = values[:, 0]
283            dataframe['y'] = values[:, 1]
284            plt.figure(figsize=(16, 16))
285            plt.title("Results of " + key)
286
287            sns_plot = sns.scatterplot(
288                x="x",
289                y="y",
290                hue=queries + list(itertools.chain.from_iterable([query] * 5 for query in queries
                        )),
291                data=dataframe,
292                legend="full",
293                alpha=1.0
294            )
295            sns_plot.get_figure().savefig("res/plot_" + key.lower())
296
297
298    def print_scores(scores):
299        print("##### PRINT #####")
300        for key, values in scores.items():
301            print(key)
302            precision, recall = compute_mean(values)
303            print("\tprecision:\t" + precision + "\n\trecall:\t\t" + recall)
304
305
```

17

```
306  def compute_mean(stats):
307      precision, recall, counter = 0, 0, 0
308      for stat in stats:
309          precision += stat.precisions
310          recall += stat.recalls
311          counter += 1
312      return str(precision / counter), str(recall / counter)
313
314
315  if len(argv) < 1:
316      print("Please give as input ground truth file")
317      exit(1)
318
319
320  start(argv[1])
```

# B   Bash Code

```
1  #!/bin/bash
2
3  rm res/*.pkl
4  python3 src/extract_data.py $1
5  python3 src/search_data.py $2
6  python3 src/prec_recall.py res/data.csv res/ground-truth.txt
```