

Project 2: Multi-source code search

Ardigò Susanna

December 16, 2020

https://github.com/SusyPinkBash/multi_source_code_search

Contents

1	Data Extraction	2
1.1	Goal and Input parameter	2
1.2	Description of the code	2
1.3	Results	2
2	Training of search engines	3
2.1	Goal and Input parameter	3
2.2	Description of the code	3
2.3	Results	3
3	Evaluation of search engines	4
3.1	Goal and Input parameter	4
3.2	Description of the code	4
3.3	Results	5
4	Visualisation of query results	5
4.1	Goal and Input parameter	5
4.2	Description of the code	5
4.3	Results	6
4.3.1	LSI	6
4.3.2	Doc2Vec	7
A	Python code	8
A.1	Data Extraction	8
A.2	Training of search engines	9
A.3	Evaluation of search engines and Visualisation of query results	13
B	Bash Code	20

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 inferring 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 document as a first result. We can see that Frequency and TF-IDF have three results in common, if we don't consider the correct result, but at different positions.

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 anything.	FREQ
2 CleanCopier	../tensorflow/tensorflow/python/autograph/pyct/ast_util.py	38	class	NodeTransformer-like visitor that copies an AST.	FREQ
3 CompatVilImportReplacer	../tensorflow/tensorflow/tools/compatibility/tf_upgrade_v2.py	72	class	AST Visitor that replaces 'import tensorflow.compat.v1 as tf'.	FREQ
4 FunctionVisitor	../tensorflow/tensorflow/python/autograph/pyct/static_analysis/type_inference.py	394	class	AST visitor that applies type inference to each function separately.	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 anything.	TF-IDF
7 FunctionVisitor	../tensorflow/tensorflow/python/autograph/pyct/static_analysis/type_inference.py	394	class	AST visitor that applies type inference to each function separately.	TF-IDF
8 track_usage	../tensorflow/tensorflow/python/platform/analytics.py	21	function	No usage tracking for external library.	TF-IDF
9 CleanCopier	../tensorflow/tensorflow/python/autograph/pyct/ast_util.py	38	class	NodeTransformer-like visitor that copies an AST.	TF-IDF
10 PublicAPIVisitor	../tensorflow/tensorflow/tools/common/public_api.py	29	class	Visitor to use with 'traverse' to visit exactly the public TF API.	TF-IDF
11 PastaAnalyzeVisitor	../tensorflow/tensorflow/tools/compatibility/ast_edits.py	815	class	AST Visitor that looks for specific API usage without editing anything.	LSI
12 matches	../tensorflow/tensorflow/python/autograph/pyct/ast_util.py	214	function	Basic pattern matcher for AST.	LSI
13 CountingVisitor	../tensorflow/tensorflow/python/autograph/pyct/cfg_test.py	28	class	<null>	LSI
14 ProfileOptionBuilder	../tensorflow/tensorflow/python/profiler/option_builder.py	27	class	Option Builder for Profiling API.	LSI
15 concentration	../tensorflow/tensorflow/python/ops/distributions/dirichlet.py	213	method	Concentration parameter; expected counts for that coordinate.	LSI
16 PastaAnalyzeVisitor	../tensorflow/tensorflow/tools/compatibility/ast_edits.py	815	class	AST Visitor that looks for specific API usage without editing anything.	Doc2Vec
17 APICChangeSpec	../tensorflow/tensorflow/tools/compatibility/ast_edits.py	192	class	This class defines the transformations that need to happen.	Doc2Vec
18 Profile	../tensorflow/tensorflow/python/profiler/profiler_v2.py	181	class	Context-manager profile API.	Doc2Vec
19 TFAPICChangeSpec	../tensorflow/tensorflow/tools/compatibility/tf_upgrade_v2_safety.py	26	class	List of maps that describe what changed in the API.	Doc2Vec
20 Profiler	../tensorflow/tensorflow/python/eager/profiler.py	176	class	Context-manager eager profiler 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 high for all engines except for one. The engine with the highest precision is **TF-IDF**, with score 0.85. The second highest precision is **Frequencies**, with score 0.83. The third highest precision is **LSI**, with score 0.80. The least precise search engine is **Doc2Vec** with score 0.55, which is the only score lower than 0.8

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. It is followed by **LSI** with a score of 0.9. At last **Doc2Vec** has a recall of 0.6.

We can say that almost all search engine have similar scores for precision and recalls. The only exception is **Doc2Vec**.

Engine	Precision	Recall
Frequencies	0.83	0.90
TD-IDF	0.85	1.00
LSI	0.80	0.80
Doc2Vec	0.55	0.60

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

Figures 2 and 3 shows the plots of the visualization of the queries. At first we notice that the LSI scatterplot tends to have better query clusters. The optimal solution is to have defined clusters for each query. This does not happen in any of the two images. There are queries that create clusters but this does not happen for all queries.

4.3.1 LSI

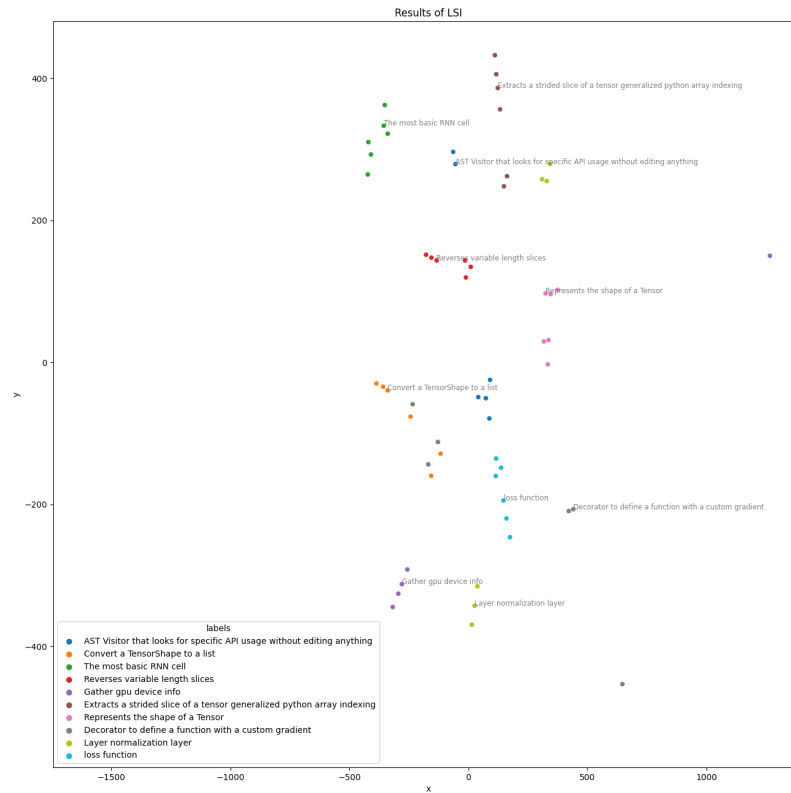


Figure 2: Visualization of the queries using LSI

Analyzing the plot of LSI, shown in figure 2, we can see that some results of the queries tend to stay close, but not completely. There are some well defined clusters composed of the same queries.

Some examples are *loss function*, *The most basic RNN cell* and *Reverses variable length slices*. The queries of *Convert a TensorShape to a list* and *Decorator to define a function with a custom gradient* that tend to have some results that are somewhat close.

4.3.2 Doc2Vec

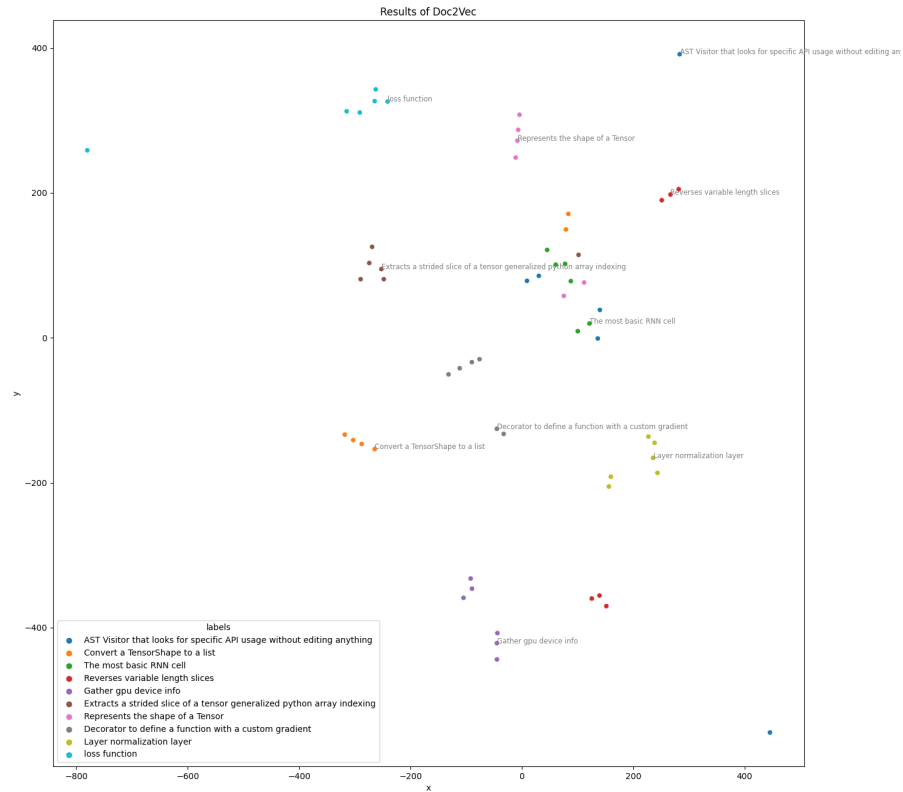


Figure 3: Visualization of the queries using Doc2Vec

Analyzing the plot of Doc2Vec, shown in figure 3, we can see that there are results of the queries tend to stay closer than the previous plot. In this image there are some cluster that are well defined and are composed only of one query. Some queries are divided into two clusters close to each other. The most well defined cluster is '*loss function*', which was also well defined in the previous plot. There is a cluster which is composed of five different queries.

A Python code

A.1 Data Extraction

```

1  from sys import argv, exit
2  from ast import *
3  from os import walk
4  import pandas as pd
5
6
7  class Visitor(NodeVisitor):
8      def __init__(self, file_path, node):
9          super().__init__()
10         self.file_path = clean_file_path(file_path)
11         self.visit(parse(node))
12
13     def visit_ClassDef(self, node: ClassDef):
14         self.generic_visit(node)
15         if is_valid_entity(node.name):
16             self.append_data(node, "class")
17
18     def visit_FunctionDef(self, node: FunctionDef):
19         if is_valid_entity(node.name):
20             self.append_data(node, "method" if is_method(node) else "function")
21
22     def append_data(self, node, def_type):
23         comment = get_docstring(node)
24         comment = comment.split('\n')[0] if comment is not None else ""
25         data.append((node.name, self.file_path, node.lineno, def_type, comment))
26
27
28 def clean_file_path(path):
29     directories = path.split('/')
30     return '..../' + '/'.join(directories[directories.index('tensorflow'):])
31
32
33 def is_valid_entity(name):
34     return name[0] != '_' and name != "main" and "test" not in name.lower()
35
36
37 def is_method(function):
38     return function.args and len(function.args.args) > 0 and 'self' in function.args.args[0].arg
39
40
41 def start(directory_path):
42     if directory_path[-1] == '/':
43         directory_path = directory_path[:-1]
44     counter = 0
45     for path, _, files in walk(directory_path):
46         for file_name in files:
47             if file_name.endswith('.py'):
48                 counter += 1
49                 file_path = path + '/' + file_name
50                 with open(file_path) as file:
51                     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

```

```

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/pickle/' + 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     pickle.dump(data, open('res/pickle/' + name + '.pkl', "wb"), protocol=pickle.HIGHEST_PROTOCOL)
77
78
79 def load_file(name):
80     return pickle.load(open('res/pickle/' + name + '.pkl', "rb"))
81
82
83 def split_space(text):
84     return text.translate(str.maketrans(' ', ' ', string.punctuation)).split(' ') if text != ""
85     else []
86
87
88 def split_underscore(tokens):
89     return [word for token in tokens for word in token.split('_')]
90
91
92 def handle_camel_case(tokens):
93     words = []
94     for token in tokens:
95         matches = finditer('.+?(?:(?:<=[a-z])|(?:<=[A-Z])|(?:<=[A-Z])|(?:<=[A-Z][a-z])|$', token)
96         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 query_frequency(query, bow, dictionary):
115     return filter_results(create_top_5_result_tuples(get_freq_model(bow, dictionary)[
116         dictionary.doc2bow(query)]))
117
118
119 def get_freq_model(bow, dictionary):
120     return load_file('model_freq') if exists_file('model_freq') else create_freq_model(bow,
121         dictionary)
122
123
124 def create_freq_model(bow, dictionary):
125     model = SparseMatrixSimilarity(bow, num_features=len(dictionary.token2id))
126     save_data(model, 'model_freq')
127     return model
128
129
130 def query_tfidf(query, bow, dictionary):
131     model = get_tfidf_model(bow)
132     matrix = get_tfidf_matrix(model, bow, dictionary)
133     return filter_results(create_top_5_result_tuples(matrix[model[dictionary.doc2bow(query)
134         ]]))
135
136
137 def get_tfidf_model(bow):
138     return load_file('model_tfidf') if exists_file('model_tfidf') else create_tfidf_model(bow
139         )
140
141
142 def create_tfidf_model(bow):
143     model = TfidfModel(bow)
144     save_data(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(
150         model, bow, dictionary)

```

```

148 def create_tfidf_matrix(model, bow, dictionary):
149     matrix = SparseMatrixSimilarity(model[bow], num_features=len(dictionary.token2id))
150     save_data(matrix, 'matrix_tfidf')
151     return model
152
153
154 def query_lsi(query, bow, dictionary):
155     model = get_lsi_model(bow, dictionary)
156     matrix = get_lsi_matrix(model, bow)
157     vector = model[dictionary.doc2bow(query)]
158     result = abs(matrix[vector])
159     embedding = [[value for _, value in vector]] + [[value for _, value in model[bow][i]] for
160                                                     i, value in
161                                                         sorted(enumerate(result), key=lambda x: x
162                                                             [1], reverse=True)[:5]]
163
164     return filter_results(create_top_5_result_tuples(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,
169     dictionary)
170
171
172 def create_lsi_model(bow, dictionary):
173     model = LsiModel(bow, id2word=dictionary, num_topics=300)
174     save_data(model, 'model_lsi')
175     return model
176
177
178 def get_lsi_matrix(model, bow):
179     return load_file('matrix_lsi') if exists_file('matrix_lsi') else create_lsi_matrix(model,
180     bow)
181
182
183 def create_lsi_matrix(model, bow):
184     matrix = MatrixSimilarity(model[bow])
185     save_data(matrix, 'matrix_lsi')
186     return matrix
187
188
189 def create_top_5_result_tuples(arr):
190     return sorted(enumerate(arr), key=lambda x: x[1], reverse=True)[:5]
191
192
193 def filter_results(tuples):
194     return [i for i, v in tuples]
195
196
197 def query_doc2vec(query, corpus):
198     model = get_doc2vec_model(get_doc2vec_corpus(corpus))
199     vector = model.infer_vector(query)
200     similar = model.docvecs.most_similar([vector], topn=5)
201     return [index for (index, _) in similar], \
202         [list(vector)] + [list(model.infer_vector(corpus[index])) for index, _ in similar]
203
204
205 def get_doc2vec_corpus(corpus):

```

```

201     return [TaggedDocument(simple_preprocess(' '.join(element)), [index])
202             for index, element in enumerate(corpus)]
203
204
205 def get_doc2vec_model(corpus):
206     return load_file('model_doc2vec') if exists_file('model_doc2vec') else
        create_doc2vec_model(corpus)
207
208
209 def create_doc2vec_model(corpus):
210     model = Doc2Vec(vector_size=300, min_count=2, epochs=77)
211     model.build_vocab(corpus)
212     model.train(corpus, total_examples=model.corpus_count, epochs=model.epochs)
213     save_data(model, 'model_doc2vec')
214     return model
215
216
217 def create_result_dataframe(queries_dictionary, df):
218     for key, values in queries_dictionary.items():
219         for index in values:
220             row = df.iloc[index]
221             yield [row["name"], row["file"], row["line"], row["type"], row["comment"], key]
222
223
224 def print_results(df):
225     grouped = df.groupby(['search'])
226     for key, item in grouped:
227         print(grouped.get_group(key), "\n\n")
228
229
230 if len(argv) < 2:
231     print("Please give as input the query")
232     exit(1)
233
234 print("NOPE")
235 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'],
25                             vectors)
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/pickle/' + 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     pickle.dump(data, open('res/pickle/' + name + '.pkl', "wb"), protocol=pickle.HIGHEST_PROTOCOL)
77
78
79 def load_file(name):
80     return pickle.load(open('res/pickle/' + name + '.pkl', "rb"))
81
82
83 def split_space(text):
84     return text.translate(str.maketrans('', '', string.punctuation)).split(' ') if text != ""
85         else []
86
87
88 def split_underscore(tokens):
89     return [word for token in tokens for word in token.split('_')]
90
91
92 def handle_camel_case(tokens):
93     words = []
94     for token in tokens:
95         matches = finditer('.+?(?:(?<=[a-z])(?=[A-Z])|(?<=[A-Z])(?=[A-Z][a-z])|$)', token)
96         words += [m.group(0) for m in matches]
97     return words
98
99
100 def normalize_tokens(tokens):
101     return [token.lower() for token in tokens]
102
103
104 def filter_stopwords(tokens):
105     for token in tokens:
106         if token in ['test', 'tests', 'main']:
107             return []
108     return tokens
109
110
111 def normalize_query(query):
112     return query.strip().lower().split()
113
114
115 def query_frequency(query, bow, dictionary):
116     return filter_results(get_freq_model(bow, dictionary)[dictionary.doc2bow(query)])
117
118
119 def get_freq_model(bow, dictionary):
120     return load_file('model_freq') if exists_file('model_freq') else create_freq_model(bow,
121         dictionary)
122
123
124 def create_freq_model(bow, dictionary):
125     model = SparseMatrixSimilarity(bow, num_features=len(dictionary.token2id))
126     save_data(model, 'model_freq')
127     return model

```



```

127
128 def query_tfidf(query, bow, dictionary):
129     model = get_tfidf_model(bow)
130     matrix = get_tfidf_matrix(model, bow, dictionary)
131     return filter_results(matrix[model[dictionary.doc2bow(query)]])
132
133
134 def get_tfidf_model(bow):
135     return load_file('model_tfidf') if exists_file('model_tfidf') else create_tfidf_model(bow
136     )
137
138 def create_tfidf_model(bow):
139     model = TfidfModel(bow)
140     save_data(model, 'model_tfidf')
141     return model
142
143
144 def get_tfidf_matrix(model, bow, dictionary):
145     return load_file('matrix_tfidf') if exists_file('matrix_tfidf') else create_tfidf_matrix(
146         model, bow, dictionary)
147
148 def create_tfidf_matrix(model, bow, dictionary):
149     matrix = SparseMatrixSimilarity(model[bow], num_features=len(dictionary.token2id))
150     save_data(matrix, 'matrix_tfidf')
151     return model
152
153
154 def query_lsi(query, bow, dictionary):
155     model = get_lsi_model(bow, dictionary)
156     matrix = get_lsi_matrix(model, bow)
157     vector = model[dictionary.doc2bow(query)]
158     result = abs(matrix[vector])
159     embedding = [[value for _, value in vector]] + [[value for _, value in model[bow][i]] for
160                                                         i, value in
                                                         sorted(enumerate(result), key=lambda x: x
                                                         [1], reverse=True)[:5]]
161     return filter_results(result), embedding
162
163
164 def get_lsi_model(bow, dictionary):
165     return load_file('model_lsi') if exists_file('model_lsi') else create_lsi_model(bow,
166         dictionary)
167
168 def create_lsi_model(bow, dictionary):
169     model = LsiModel(bow, id2word=dictionary, num_topics=300)
170     save_data(model, 'model_lsi')
171     return model
172
173
174 def get_lsi_matrix(model, bow):
175     return load_file('matrix_lsi') if exists_file('matrix_lsi') else create_lsi_matrix(model,
176         bow)
177

```

```

178 def create_lsi_matrix(model, bow):
179     matrix = MatrixSimilarity(model[bow])
180     save_data(matrix, 'matrix_lsi')
181     return matrix
182
183
184 def filter_results(arrg):
185     return [i for i, v in sorted(enumerate(arrg), key=lambda x: x[1], reverse=True)[:5]]
186
187
188 def query_doc2vec(query, corpus):
189     model = get_doc2vec_model(get_doc2vec_corpus(corpus))
190     vector = model.infer_vector(query)
191     similar = model.docvecs.most_similar([vector], topn=5)
192     return [index for (index, _) in similar], \
193           [list(vector)] + [list(model.infer_vector(corpus[index])) for index, _ in similar]
194
195
196 def get_doc2vec_corpus(corpus):
197     return [TaggedDocument(simple_preprocess(' '.join(element)), [index])
198           for index, element in enumerate(corpus)]
199
200
201 def get_doc2vec_model(corpus):
202     return load_file('model_doc2vec') if exists_file('model_doc2vec') else
203           create_doc2vec_model(corpus)
204
205
206 def create_doc2vec_model(corpus):
207     model = Doc2Vec(vector_size=300, min_count=2, epochs=77)
208     model.build_vocab(corpus)
209     model.train(corpus, total_examples=model.corpus_count, epochs=model.epochs)
210     save_data(model, 'model_doc2vec')
211     return model
212
213
214 def create_result_dataframe(queries_dictionary, df):
215     for key, values in queries_dictionary.items():
216         for index in values:
217             row = df.iloc[index]
218             yield [row["name"], row["file"], row["line"], row["type"], row["comment"], key]
219
220 #####
221
222 class Truth:
223     def __init__(self, query, name, path):
224         self.name = name
225         self.path = path
226         self.query = query.lower()
227
228
229 class Stat:
230     def __init__(self, precisions, recalls):
231         self.precisions = precisions
232         self.recalls = recalls
233

```

```

234
235 def start(path_ground_truth):
236     dataframe = pd.read_csv("res/data.csv").fillna(value="")
237     ground_truth, labels = parse_ground_truth(path_ground_truth)
238     scores, vectors = compute_precision_recall(ground_truth, dataframe)
239     plot_vectors(compute_tsne(vectors), labels)
240     print_scores(scores)
241
242
243 def parse_ground_truth(path_ground_truth):
244     classes, labels = [], []
245     for entry in open(path_ground_truth, "r").read().split("\n\n"):
246         data = entry.split("\n")
247         query = data[0]
248         classes.append(Truth(query, data[1], data[2]))
249         labels += [query]*6
250     return classes, labels
251
252
253 def compute_precision_recall(ground_truth, dataframe):
254     scores = {"FREQ": [], "TF-IDF": [], "LSI": [], "Doc2Vec": []}
255     vectors = {"LSI": [], "Doc2Vec": []}
256     for entry in ground_truth:
257         results, vectors_i = get_results(entry.query, dataframe)
258         vectors["LSI"] += vectors_i["LSI"]
259         vectors["Doc2Vec"] += vectors_i["Doc2Vec"]
260         for query_type in ["FREQ", "TF-IDF", "LSI", "Doc2Vec"]:
261             precision = compute_precision(entry, query_type, results)
262             scores[query_type].append(Stat(precision, compute_recall(precision)))
263     return scores, vectors
264
265
266 def compute_precision(truth, search_type, dataframe):
267     counter = 0
268     for _, row in dataframe[dataframe['search'] == search_type].iterrows():
269         counter += 1
270         if row["name"] == truth.name and row["file"] == truth.path:
271             return 1 / counter
272     return 0
273
274
275 def compute_recall(precision):
276     return 1 if precision > 0 else 0
277
278
279 def compute_tsne(dictionary):
280     results = {}
281     for key, values in dictionary.items():
282         tsne = TSNE(n_components=2, verbose=1, perplexity=2, n_iter=3000)
283         results[key] = tsne.fit_transform(values)
284     return results
285
286
287 def plot_vectors(dictionary, labels):
288     for key, values in dictionary.items():
289         dataframe = pd.DataFrame()
290         dataframe['x'] = values[:, 0]

```

```

291     dataframe['y'] = values[:, 1]
292     dataframe['labels'] = labels
293     plt.figure(figsize=(16, 16))
294     plt.title("Results of " + key)
295
296     plot = sns.scatterplot(
297         x="x",
298         y="y",
299         hue='labels',
300         data=dataframe,
301         legend="full",
302         alpha=1.0
303     )
304
305     for label in range(0, len(labels), 6):
306         plot.text(
307             dataframe['x'][label],
308             dataframe['y'][label],
309             dataframe['labels'][label],
310             horizontalalignment='left', size='small', color='gray'
311         )
312
313     plot.get_figure().savefig("res/plot_" + key.lower())
314
315
316 def print_scores(scores):
317     print("#### PRINT ####")
318     for key, values in scores.items():
319         print(key)
320         precision, recall = compute_mean(values)
321         print("\tprecision:\t" + precision + "\n\trecall:\t\t" + recall)
322
323
324 def compute_mean(stats):
325     precision, recall, counter = 0, 0, 0
326     for stat in stats:
327         precision += stat.precisions
328         recall += stat.recalls
329         counter += 1
330     return str(precision / counter), str(recall / counter)
331
332
333 if len(argv) < 1:
334     print("Please give as input ground truth file")
335     exit(1)
336
337
338 start(argv[1])

```

B Bash Code

```
1 #!/bin/bash
2
3 rm -rf res/pickle
4 mkdir res/pickle
5 python3 src/extract_data.py $1
6 python3 src/search_data.py $2
7 python3 src/prec_recall.py res/ground-truth-unique.txt
```