

Project 2: Multi-source code search

Ardigò Susanna

December 7, 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 | 5 |
| A | Python code | 6 |
| A.1 | Data Extraction | 6 |
| A.2 | Training of search engines | 7 |
| A.3 | Evaluation of search engines and Visualisation of query results | 10 |
| B | Bash Code | 15 |

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 that is a method 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 I 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.

| 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.

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 normalize the query that the second part needs to produce the results. The first part of function `create_data(dataframe)` extracts the names and comments from the dataframes 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. 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 save the trained model in an external pickle file to load it in the next runs. This improves the running time of this function. 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 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: *'Optimizer that implements the Adadeltha algorithm'*.

Figure 1 show the result of the given query. As we can see in the image all results are classes. Almost all results have as comment the sentence *'Optimizer that implements the x algorithm'* with x being an algorithm name.

The most common document found is **Adadelta** with 3 findings, **AdadeltaOptimizer** with 3. Both documents have the same comments. The first, with path `../tensorflow/tensorflow/python/keras/optimizer_v2/adadelta.py`, was found by Frequencies, TD-IDF and LSI. The second, with path `../tensorflow/tensorflow/python/training/adadelta.py`, was found by Frequency, TF-IDF and Doc2Vec. Frequency and TF-IDF have the best results.

| # | name | file | Line | type | comment | search |
|----|-------------------|--|------|-------|---|---------|
| 1 | Nadam | ../tensorflow/tensorflow/python/keras/optimizer_v2/nadam.py | 34 | class | Optimizer that implements the Nadam algorithm. | FREQ |
| 2 | Adadelta | ../tensorflow/tensorflow/python/keras/optimizer_v2/adadelta.py | 32 | class | Optimizer that implements the Adadelta algorithm. | FREQ |
| 3 | FtrlOptimizer | ../tensorflow/tensorflow/python/training/ftrl.py | 29 | class | Optimizer that implements the FTRL algorithm. | FREQ |
| 4 | AdagradOptimizer | ../tensorflow/tensorflow/python/training/adagrad.py | 32 | class | Optimizer that implements the Adagrad algorithm. | FREQ |
| 5 | AdadeltaOptimizer | ../tensorflow/tensorflow/python/training/adadelta.py | 29 | class | Optimizer that implements the Adadelta algorithm. | FREQ |
| 6 | Adadelta | ../tensorflow/tensorflow/python/keras/optimizers.py | 383 | class | Adadelta optimizer. | TF-IDF |
| 7 | Nadam | ../tensorflow/tensorflow/python/keras/optimizer_v2/nadam.py | 34 | class | Optimizer that implements the Nadam algorithm. | TF-IDF |
| 8 | Adadelta | ../tensorflow/tensorflow/python/keras/optimizer_v2/adadelta.py | 32 | class | Optimizer that implements the Adadelta algorithm. | TF-IDF |
| 9 | AdamOptimizer | ../tensorflow/tensorflow/python/training/adam.py | 32 | class | Optimizer that implements the Adam algorithm. | TF-IDF |
| 10 | AdadeltaOptimizer | ../tensorflow/tensorflow/python/training/adadelta.py | 29 | class | Optimizer that implements the Adadelta algorithm. | TF-IDF |
| 11 | RMSprop | ../tensorflow/tensorflow/python/keras/optimizer_v2/rmsprop.py | 35 | class | Optimizer that implements the RMSprop algorithm. | LSI |
| 12 | Adamax | ../tensorflow/tensorflow/python/keras/optimizer_v2/adamax.py | 33 | class | Optimizer that implements the Adamax algorithm. | LSI |
| 13 | Ftrl | ../tensorflow/tensorflow/python/keras/optimizer_v2/ftrl.py | 30 | class | Optimizer that implements the FTRL algorithm. | LSI |
| 14 | Adam | ../tensorflow/tensorflow/python/keras/optimizer_v2/adam.py | 34 | class | Optimizer that implements the Adam algorithm. | LSI |
| 15 | Adadelta | ../tensorflow/tensorflow/python/keras/optimizer_v2/adadelta.py | 32 | class | Optimizer that implements the Adadelta algorithm. | LSI |
| 16 | Adam | ../tensorflow/tensorflow/python/keras/optimizer_v2/adam.py | 34 | class | Optimizer that implements the Adam algorithm. | Doc2Vec |
| 17 | AdagradOptimizer | ../tensorflow/tensorflow/python/training/adagrad.py | 32 | class | Optimizer that implements the Adagrad algorithm. | Doc2Vec |
| 18 | AdamOptimizer | ../tensorflow/tensorflow/python/training/adam.py | 32 | class | Optimizer that implements the Adam algorithm. | Doc2Vec |
| 19 | MomentumOptimizer | ../tensorflow/tensorflow/python/training/momentum.py | 29 | class | Optimizer that implements the Momentum algorithm. | Doc2Vec |
| 20 | AdadeltaOptimizer | ../tensorflow/tensorflow/python/training/adadelta.py | 29 | class | Optimizer that implements the Adadelta algorithm. | 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 ground truth. We can see that the precision is low for all engines. The engine with the highest precision is **Doc2Vec**. The recall is higher than the precision. The engine **TF-IDF** has a recall equal to 1. 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 | 0.417 | 0.8 |

Table 2: Statistics of the search engines

4 Visualisation of query results

4.1 Goal and Input parameter

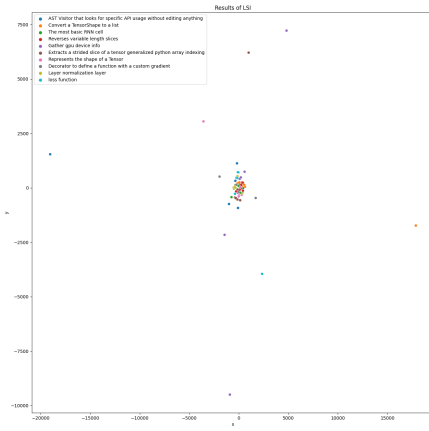
This part of the project consists of visualizing the embedding vectors of the queries and the top 5 answers in a 2D plot. This file takes as argument the ground truth file.

4.2 Description of the code

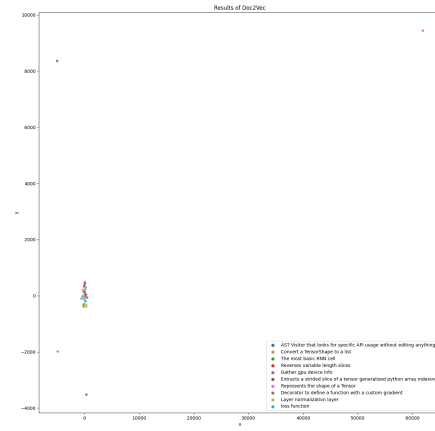
The first part of the execution is the same as the previous file. After the results are calculated, we plot the TSNE of the embedding vectors, that we retrieved in the explanation above but we did not use. The plot is straight-forward: we create a dataframe with the information of x and y coordinates and print them of different hues.

4.3 Results

Figure 2 shows the plots of the visualization of the queries.



(a) Results of LSI



(b) Results of Doc2Vec

Figure 2: Visualization of the plots of the 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 = 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 is_valid_entity(name):
29 |     return name[0] != '_' and name != "main" and "test" not in name.lower()
30 |
31 |
32 | def is_method(function):
33 |     return function.args and len(function.args.args) > 0 and 'self' in function.args.args[0].
        arg
34 |
35 |
36 | def start(directory_path):
37 |     if directory_path[-1] == '/':
38 |         directory_path = directory_path[: -1]
39 |     counter = 0
40 |     for path, _, files in walk(directory_path):
41 |         for file_name in files:
42 |             if file_name.endswith('.py'):
43 |                 counter += 1
44 |                 file_path = path + '/' + file_name
45 |                 with open(file_path) as file:
46 |                     Visitor(file_path, file.read())
47 |
48 |     dataframe = pd.DataFrame(data=data, columns=["name", "file", "line", "type", "comment"])
49 |     dataframe.to_csv('res/data.csv', index=False, encoding='utf-8')
50 |     print("files\t      " + str(counter))
51 |     print(dataframe["type"].value_counts())
52 |
53 |
54 | if len(argv) < 2:
55 |     print("Please give as input the path of the directory to analyze")
56 |     exit(1)
57 | data = []
58 | 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 = create_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 create_data(df):
41     tokens = [filter_stopwords(normalize_tokens(handle_camel_case(split_underscore(
42         [row["name"]] + split_space(row["comment"])))) for _, row in df.iterrows())]
43
44     frequency = defaultdict(int)
45     for token in tokens:
46         for word in token:
47             frequency[word] += 1
48
49     processed = [[token for token in text if frequency[token] > 1] for text in tokens]
50     dictionary = Dictionary(processed)
51     bow = [dictionary.doc2bow(text) for text in processed]
52
53     return processed, dictionary, bow
54
55
56 def split_space(text):
57     return text.translate(str.maketrans('', '', string.punctuation)).split(' ') if text != ""
58     else []
59
60 def split_underscore(tokens):
61     return [word for token in tokens for word in token.split('_')]
62
63
64 def handle_camel_case(tokens):
65     words = []
66     for token in tokens:
67         matches = finditer('.+?(?:(?<=[a-z]) (?=[A-Z]) | (?<=[A-Z]) (?=[A-Z][a-z]) | $)', token)
68         words += [m.group(0) for m in matches]
69     return words
70
71
72 def normalize_tokens(tokens):
73     return [token.lower() for token in tokens]
74

```

```

75
76 def filter_stopwords(tokens):
77     for token in tokens:
78         if token in ['test', 'tests', 'main']:
79             return []
80     return tokens
81
82
83 def normalize_query(query):
84     return query.strip().lower().split()
85
86
87 def query_frequency(query, bow, dictionary):
88     return filter_results(SparseMatrixSimilarity(bow, num_features=len(dictionary.token2id))[
89         dictionary.doc2bow(query)])
90
91 def query_tfidf(query, bow, dictionary):
92     model = TfidfModel(bow)
93     return filter_results(SparseMatrixSimilarity(model[bow], num_features=len(dictionary.
94         token2id))[model[dictionary.doc2bow(query)]])
95
96 def query_lsi(query, bow, dictionary):
97     model = LsiModel(bow, id2word=dictionary, num_topics=300)
98     vector = model[dictionary.doc2bow(query)]
99     result = abs(MatrixSimilarity(model[bow])[vector])
100     embedding = [[value for _, value in vector]] + [[value for _, value in model[bow][i]] for
101         i, value in
102             sorted(enumerate(result), key=lambda x: x
103                 [1], reverse=True)[:5]]
104
105     return filter_results(result), embedding
106
107
108 def filter_results(arr):
109     return [i for i, v in sorted(enumerate(arr), key=lambda x: x[1], reverse=True)[:5]]
110
111
112 def query_doc2vec(query, corpus):
113     model = get_doc2vec_model(get_doc2vec_corpus(corpus))
114     vector = model.infer_vector(query)
115     similar = model.docvecs.most_similar([vector], topn=5)
116     return [index for (index, _) in similar], \
117         [list(vector)] + [list(model.infer_vector(corpus[index])) for index, _ in similar]
118
119
120 def get_doc2vec_corpus(corpus):
121     return [TaggedDocument(simple_preprocess(' '.join(element)), [index])
122         for index, element in enumerate(corpus)]
123
124
125 def get_doc2vec_model(corpus):
126     return pickle.load(open('res/doc2vec.pkl', "rb")) if path.exists('res/doc2vec.pkl') else
127         create_doc2vec_model(corpus)
128
129
130 def create_doc2vec_model(corpus):

```

```

127     model = Doc2Vec(vector_size=300, min_count=2, epochs=77)
128     model.build_vocab(corpus)
129     model.train(corpus, total_examples=model.corpus_count, epochs=model.epochs)
130     pickle.dump(model, open('res/doc2vec.pkl', "wb"), protocol=pickle.HIGHEST_PROTOCOL)
131     return model
132
133
134 def create_result_dataframe(queries_dictionary, df):
135     for key, values in queries_dictionary.items():
136         for index in sorted(values):
137             row = df.iloc[index]
138             yield [row["name"], row["file"], row["line"], row["type"], row["comment"], key]
139
140
141 def print_results(df):
142     grouped = df.groupby(['search'])
143     for key, item in grouped:
144         print(grouped.get_group(key), "\n\n")
145
146
147 if len(argv) < 2:
148     print("Please give as input the query")
149     exit(1)
150
151 begin_time = datetime.now()
152 start(argv[1])
153 print(datetime.now() - begin_time)

```

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)

```



```

80     return tokens
81
82
83 def normalize_query(query):
84     return query.strip().lower().split()
85
86
87 def query_frequency(query, bow, dictionary):
88     return SparseMatrixSimilarity(bow, num_features=len(dictionary.token2id))[dictionary.doc2bow(query)]
89
90
91 def query_tfidf(query, bow, dictionary):
92     model = TfidfModel(bow)
93     return SparseMatrixSimilarity(model[bow], num_features=len(dictionary.token2id))[model[dictionary.doc2bow(query)]]
94
95
96 def query_lsi(query, bow, dictionary):
97     model = LsiModel(bow, id2word=dictionary, num_topics=300)
98     vector = model[dictionary.doc2bow(query)]
99     result = abs(MatrixSimilarity(model[bow])[vector])
100    embedding = [[value for _, value in vector]] + [[value for _, value in model[bow][i]] for
101                                                    i, value in
102                                                         sorted(enumerate(result), key=lambda x: x
103                                                           [1], reverse=True)[:5]]
104
105    return filter_results(result), embedding
106
107
108 def filter_results(arr):
109     return [i for i, v in sorted(enumerate(arr), key=lambda x: x[1], reverse=True)[:5]]
110
111
112 def query_doc2vec(query, corpus):
113     model = get_doc2vec_model(get_doc2vec_corpus(corpus))
114     vector = model.infer_vector(query)
115     similar = model.docvecs.most_similar([vector], topn=5)
116     return [index for (index, _) in similar], \
117            [list(vector)] + [list(model.infer_vector(corpus[index])) for index, _ in similar]
118
119
120 def get_doc2vec_corpus(corpus):
121     return [TaggedDocument(simple_preprocess(' '.join(element)), [index])
122             for index, element in enumerate(corpus)]
123
124
125 def get_doc2vec_model(corpus):
126     return pickle.load(open('res/doc2vec.pkl', "rb")) if path.exists('res/doc2vec.pkl') else
127            create_doc2vec_model(corpus)
128
129
130 def create_doc2vec_model(corpus):
131     model = Doc2Vec(vector_size=300, min_count=2, epochs=77)
132     model.build_vocab(corpus)
133     model.train(corpus, total_examples=model.corpus_count, epochs=model.epochs)
134     pickle.dump(model, open('res/doc2vec.pkl', "wb"), protocol=pickle.HIGHEST_PROTOCOL)
135     return model

```

```

132
133
134 def create_result_dataframe(queries_dictionary, df):
135     for key, values in queries_dictionary.items():
136         for index in sorted(values):
137             row = df.iloc[index]
138             yield [row["name"], row["file"], row["line"], row["type"], row["comment"], key]
139
140
141 #####
142
143 class Truth:
144     def __init__(self, query, name, path):
145         self.name = name
146         self.path = path
147         self.query = query.lower()
148
149
150 class Stat:
151     def __init__(self, precisions, recalls):
152         self.precisions = precisions
153         self.recalls = recalls
154
155
156 def start(path_ground_truth):
157     dataframe = pd.read_csv("res/data.csv").fillna(value="")
158     ground_truth, queries = parse_ground_truth(path_ground_truth)
159     scores, vectors = compute_precision_recall(ground_truth, dataframe)
160     plot_vectors(compute_tsne(vectors), queries)
161     print_scores(scores)
162
163
164 def parse_ground_truth(path_ground_truth):
165     classes, queries = [], []
166     for entry in open(path_ground_truth, "r").read().split("\n\n"):
167         data = entry.split("\n")
168         classes.append(Truth(data[0], data[1], data[2]))
169         queries.append(data[0])
170     return classes, queries
171
172
173 def compute_precision_recall(ground_truth, dataframe):
174     scores = {"FREQ": [], "TF-IDF": [], "LSI": [], "Doc2Vec": []}
175     vectors = {"LSI": [], "Doc2Vec": []}
176     for entry in ground_truth:
177         results, vectors_i = get_results(entry.query, dataframe)
178         vectors["LSI"] += vectors_i["LSI"]
179         vectors["Doc2Vec"] += vectors_i["Doc2Vec"]
180         for query_type in ["FREQ", "TF-IDF", "LSI", "Doc2Vec"]:
181             precision = compute_precision(entry, query_type, results)
182             scores[query_type].append(Stat(precision, compute_recall(precision)))
183     return scores, vectors
184
185
186 def compute_precision(truth, search_type, dataframe):
187     precision, counter = 0, 0
188     for _, row in dataframe[dataframe['search'] == search_type].iterrows():

```

```

189         if row["name"] == truth.name and row["file"] == truth.path:
190             return 1 / (counter + 1)
191         counter += 1
192     return precision
193
194
195 def compute_recall(precision):
196     return 1 if precision > 0 else 0
197
198
199 def compute_tsne(dictionary):
200     results = {}
201     for key, values in dictionary.items():
202         tsne = TSNE(n_components=2, verbose=1, perplexity=2, n_iter=3000)
203         results[key] = tsne.fit_transform(values)
204     return results
205
206
207 def plot_vectors(dictionary, queries):
208     for key, values in dictionary.items():
209         dataframe = pd.DataFrame()
210         dataframe['x'] = values[:, 0]
211         dataframe['y'] = values[:, 1]
212         plt.figure(figsize=(16, 16))
213         plt.title("Results of " + key)
214
215         sns_plot = sns.scatterplot(
216             x="x",
217             y="y",
218             hue=queries + list(itertools.chain.from_iterable([query] * 5 for query in queries
219
220                 data=dataframe,
221                 legend="full",
222                 alpha=1.0
223             )
224             sns_plot.get_figure().savefig("res/plot_" + key.lower())
225
226 def print_scores(scores):
227     print("#### PRINT ####")
228     for key, values in scores.items():
229         print(key)
230         precision, recall = compute_mean(values)
231         print("\tprecision:\t" + precision)
232         print("\trecall:\t\t" + recall)
233
234
235 def compute_mean(stats):
236     precision, recall, counter = 0, 0, 0
237     for stat in stats:
238         precision += stat.precisions
239         recall += stat.recalls
240         counter += 1
241     return str(precision / counter), str(recall / counter)
242
243
244 if len(argv) < 1:

```

```
245     print("Please give as input ground truth file")
246     exit(1)
247
248
249 begin_time = datetime.now()
250 start(argv[1])
251 print(datetime.now() - begin_time)
```

B Bash Code

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