# Project 2: Multi-source code search

Ardigò Susanna

December 6, 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 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 Tensorlow 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 create a vector infering 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 Adadelta algorithm*’. Figure 1 show the result of the given query.

| | 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.

| Engine | Precision | Recall |
|--------|-----------|--------|
| Frequencies | 0.332 | 0.9 |
| TD-IDF | 0.365 | 1.0 |
| LSI | 0.403 | 0.8 |
| Doc2Vec | 0.508 | 0.8 |

Table 2: Statistics of the search engines
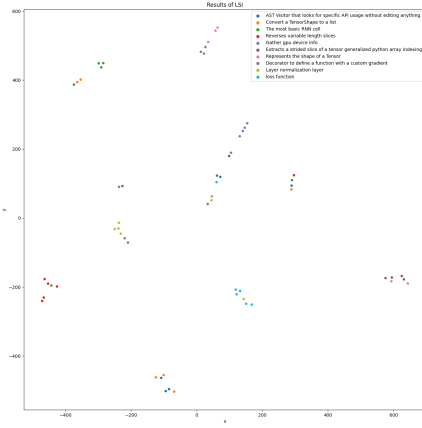
# 4 Visualisation of query results

## 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.
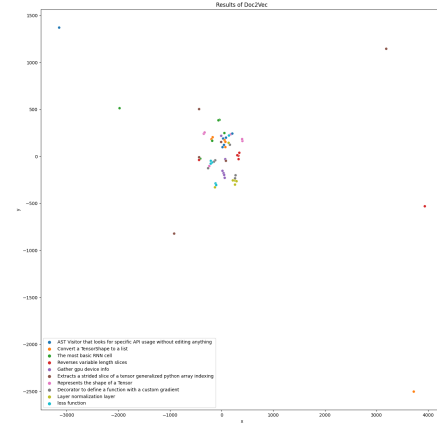
## 4.2 Description of the code

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

(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 = 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))
```

```
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   import pandas as pd
4   from re import finditer
5   from sys import argv, exit
6   from collections import defaultdict
7   from gensim.corpora import Dictionary
8   from gensim.models.doc2vec import TaggedDocument
9   from gensim.utils import simple_preprocess
10  from gensim.models import TfidfModel, LsiModel, Doc2Vec
11  from gensim.similarities import MatrixSimilarity, SparseMatrixSimilarity
12
13
14  def start(query):
15      dataframe = load_csv("res/data.csv")
16      results_dictionary, _ = compute_results(query, dataframe)
17      results = pd.DataFrame(data=create_result_dataframe(results_dictionary, dataframe),
18                             columns=['name', "file", "line", "type", "comment", "search"])
19      pd.options.display.max_colwidth = 200
```

```
20      print_results(results)
21      results.to_latex('res/search_data.tex', index=False, encoding='utf-8')
22      results.to_csv('res/search_data.csv', index=False, encoding='utf-8')
23
24
25  def compute_results(query, dataframe):
26      processed_corpus, frequencies, bag_of_words = create_data(dataframe)
27      query_to_execute = normalize_query(query)
28      results = {
29          "FREQ": query_frequency(query_to_execute, bag_of_words, frequencies),
30          "TF-IDF": query_tfidf(query_to_execute, bag_of_words, frequencies)
31      }
32      vectors = dict()
33      results["LSI"], vectors["LSI"] = query_lsi(query_to_execute, bag_of_words, frequencies)
34      results["Doc2Vec"], vectors["Doc2Vec"] = query_doc2vec(query_to_execute, processed_corpus
            )
35      return results, vectors
36
37
38  def load_csv(path):
39      return pd.read_csv(path).fillna(value="")
40
41
42  def create_data(df):
43      tokens = [filter_stopwords(normalize_tokens(handle_camel_case(split_underscore(
44          [row["name"]] + split_space(row["comment"]))))) for _, row in df.iterrows()]
45
46      frequency = defaultdict(int)
47      for token in tokens:
48          for word in token:
49              frequency[word] += 1
50
51      processed = [[token for token in text if frequency[token] > 1] for text in tokens]
52      dictionary = Dictionary(processed)
53      bow = [dictionary.doc2bow(text) for text in processed]
54
55      return processed, dictionary, bow
56
57
58  def split_space(text):
59      return text.translate(str.maketrans('', '', string.punctuation)).split(' ') if text != ""
            else []
60
61
62  def split_underscore(tokens):
63      return [word for token in tokens for word in token.split('_')]
64
65
66  def handle_camel_case(tokens):
67      words = []
68      for token in tokens:
69          matches = finditer('.+?(?:(?<=[a-z])(?=[A-Z])|(?<=[A-Z])(?=[A-Z][a-z])|$)', token)
70          words += [m.group(0) for m in matches]
71      return words
72
73
74  def normalize_tokens(tokens):
```

```
75        return [token.lower() for token in tokens]
76
77
78   def filter_stopwords(tokens):
79        for token in tokens:
80            if token in ['test', 'tests', 'main']:
81                return []
82        return tokens
83
84
85   def normalize_query(query):
86        return query.strip().lower().split()
87
88
89   def query_frequency(query, bow, dictionary):
90        return filter_results(SparseMatrixSimilarity(bow, num_features=len(dictionary.token2id))[
91            dictionary.doc2bow(query)])
92
93   def query_tfidf(query, bow, dictionary):
94        model = TfidfModel(bow)
95        return filter_results(SparseMatrixSimilarity(model[bow], num_features=len(dictionary.
                token2id))[model[dictionary.doc2bow(query)]])
96
97
98   def query_lsi(query, bow, dictionary):
99        model = LsiModel(bow, id2word=dictionary, num_topics=300)
100       vector = model[dictionary.doc2bow(query)]
101       result = abs(MatrixSimilarity(model[bow])[vector])
102       embedding = [[value for _, value in vector]] + [[value for _, value in model[bow][i]] for
                i, value in
103                                                       sorted(enumerate(result), key=lambda x: x
                                                           [1], reverse=True)[:5]]
104       return filter_results(result), embedding
105
106
107  def filter_results(arrg):
108       return [i for i, v in sorted(enumerate(arrg), key=lambda x: x[1], reverse=True)[:5]]
109
110
111  def query_doc2vec(query, corpus):
112       model = get_doc2vec_model(get_doc2vec_corpus(corpus))
113       vector = model.infer_vector(query)
114       similar = model.docvecs.most_similar([vector], topn=5)
115       return [index for (index, _) in similar], \
116            [list(vector)] + [list(model.infer_vector(corpus[index])) for index, _ in similar]
117
118
119  def get_doc2vec_corpus(corpus):
120       return [TaggedDocument(simple_preprocess(' '.join(element)), [index])
121            for index, element in enumerate(corpus)]
122
123
124  def get_doc2vec_model(corpus):
125       model = Doc2Vec(vector_size=300, min_count=2, epochs=77)
126       model.build_vocab(corpus)
127       model.train(corpus, total_examples=model.corpus_count, epochs=model.epochs)
```

```
128        return model
129
130
131  def create_result_dataframe(queries_dictionary, df):
132      for key, values in queries_dictionary.items():
133          for index in sorted(values):
134              row = df.iloc[index]
135              yield [row["name"], row["file"], row["line"], row["type"], row["comment"], key]
136
137
138  def print_results(df):
139      grouped = df.groupby(['search'])
140      for key, item in grouped:
141          print(grouped.get_group(key), "\n\n")
142
143
144  if len(argv) < 2:
145      print("Please give as input the query")
146      exit(1)
147
148  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   import seaborn as sns
7   from re import finditer
8   from sys import argv, exit
9   import matplotlib.pyplot as plt
10  from sklearn.manifold import TSNE
11  from collections import defaultdict
12  from gensim.corpora import Dictionary
13  from gensim.models.doc2vec import TaggedDocument
14  from gensim.utils import simple_preprocess
15  from gensim.models import TfidfModel, LsiModel, Doc2Vec
16  from gensim.similarities import MatrixSimilarity, SparseMatrixSimilarity
17
18  ###################
19  def get_results(query, dataframe):
20      results_dictionary, vectors = compute_results(query, dataframe)
21      return pd.DataFrame(data=create_result_dataframe(results_dictionary, dataframe),
22                          columns=['name', "file", "line", "type", "comment", "search"]),
                                 vectors
23
24
25  def compute_results(query, dataframe):
26      processed_corpus, frequencies, bag_of_words = create_data(dataframe)
27      query_to_execute = normalize_query(query)
28      results = {
29          "FREQ": filter_results(query_frequency(query_to_execute, bag_of_words, frequencies)),
30          "TF-IDF": filter_results(query_tfidf(query_to_execute, bag_of_words, frequencies))
```

```
31         }
32     vectors = dict()
33     results["LSI"], vectors["LSI"] = query_lsi(query_to_execute, bag_of_words, frequencies)
34     results["Doc2Vec"], vectors["Doc2Vec"] = query_doc2vec(query_to_execute, processed_corpus
               )
35     return results, vectors
36
37
38 def load_csv(path):
39     return pd.read_csv(path).fillna(value="")
40
41
42 def create_data(df):
43     tokens = [filter_stopwords(normalize_tokens(handle_camel_case(split_underscore(
44         [row["name"]] + split_space(row["comment"]))))) for _, row in df.iterrows()]
45
46     frequency = defaultdict(int)
47     for token in tokens:
48         for word in token:
49             frequency[word] += 1
50
51     processed = [[token for token in text if frequency[token] > 1] for text in tokens]
52     dictionary = Dictionary(processed)
53     bow = [dictionary.doc2bow(text) for text in processed]
54
55     return processed, dictionary, bow
56
57
58 def split_space(text):
59     return text.translate(str.maketrans('', '', string.punctuation)).split(' ') if text != ""
               else []
60
61
62 def split_underscore(tokens):
63     return [word for token in tokens for word in token.split('_')]
64
65
66 def handle_camel_case(tokens):
67     words = []
68     for token in tokens:
69         matches = finditer('.+?(?:(?<=[a-z])(?=[A-Z])|(?<=[A-Z])(?=[A-Z][a-z])|$)', token)
70         words += [m.group(0) for m in matches]
71     return words
72
73
74 def normalize_tokens(tokens):
75     return [token.lower() for token in tokens]
76
77
78 def filter_stopwords(tokens):
79     for token in tokens:
80         if token in ['test', 'tests', 'main']:
81             return []
82     return tokens
83
84
85 def normalize_query(query):
```

11

```
86        return query.strip().lower().split()
87
88
89  def query_frequency(query, bow, dictionary):
90      return SparseMatrixSimilarity(bow, num_features=len(dictionary.token2id))[dictionary.
            doc2bow(query)]
91
92
93  def query_tfidf(query, bow, dictionary):
94      model = TfidfModel(bow)
95      return SparseMatrixSimilarity(model[bow], num_features=len(dictionary.token2id))[model[
            dictionary.doc2bow(query)]]
96
97
98  def query_lsi(query, bow, dictionary):
99      model = LsiModel(bow, id2word=dictionary, num_topics=300)
100     vector = model[dictionary.doc2bow(query)]
101     result = abs(MatrixSimilarity(model[bow])[vector])
102     embedding = [[value for _, value in vector]] + [[value for _, value in model[bow][i]] for
            i, value in
103                                               sorted(enumerate(result), key=lambda x: x
                                                    [1], reverse=True)[:5]]
104     return filter_results(result), embedding
105
106
107 def filter_results(arrg):
108     return [i for i, v in sorted(enumerate(arrg), key=lambda x: x[1], reverse=True)[:5]]
109
110
111 def query_doc2vec(query, corpus):
112     model = get_doc2vec_model(get_doc2vec_corpus(corpus))
113     vector = model.infer_vector(query)
114     similar = model.docvecs.most_similar([vector], topn=5)
115     return [index for (index, _) in similar], \
116             [list(vector)] + [list(model.infer_vector(corpus[index])) for index, _ in similar]
117
118
119 def get_doc2vec_corpus(corpus):
120     return [TaggedDocument(simple_preprocess(' '.join(element)), [index])
121             for index, element in enumerate(corpus)]
122
123
124 def get_doc2vec_model(corpus):
125     model = Doc2Vec(vector_size=300, min_count=2, epochs=77)
126     model.build_vocab(corpus)
127     model.train(corpus, total_examples=model.corpus_count, epochs=model.epochs)
128     return model
129
130
131 def create_result_dataframe(queries_dictionary, df):
132     for key, values in queries_dictionary.items():
133         for index in sorted(values):
134             row = df.iloc[index]
135             yield [row["name"], row["file"], row["line"], row["type"], row["comment"], key]
136
137
138 ####################################
```

```
139
140   class Truth:
141       def __init__(self, query, name, path):
142           self.name = name
143           self.path = path
144           self.query = query.lower()
145
146
147   class Stat:
148       def __init__(self, precisions, recalls):
149           self.precisions = precisions
150           self.recalls = recalls
151
152
153   def start(path_ground_truth):
154       dataframe = pd.read_csv("res/data.csv").fillna(value="")
155       ground_truth, queries = parse_ground_truth(path_ground_truth)
156       scores, vectors = compute_precision_recall(ground_truth, dataframe)
157       plot_vectors(compute_tsne(vectors), queries)
158       print_scores(scores)
159
160
161   def parse_ground_truth(path_ground_truth):
162       print("##### GROUND TRUTH #####")
163       classes, queries = [], []
164       for entry in open(path_ground_truth, "r").read().split("\n\n"):
165           data = entry.split("\n")
166           classes.append(Truth(data[0], data[1], data[2]))
167           queries.append(data[0])
168       return classes, queries
169
170
171   def compute_precision_recall(ground_truth, dataframe):
172       scores = {"FREQ": [], "TF-IDF": [], "LSI": [], "Doc2Vec": []}
173       vectors = {"LSI": [], "Doc2Vec": []}
174       for entry in ground_truth:
175           results, vectors_i = get_results(entry.query, dataframe)
176           vectors["LSI"] += vectors_i["LSI"]
177           vectors["Doc2Vec"] += vectors_i["Doc2Vec"]
178           for query_type in ["FREQ", "TF-IDF", "LSI", "Doc2Vec"]:
179               precision = compute_precision(entry, query_type, results)
180               scores[query_type].append(Stat(precision, compute_recall(precision)))
181       return scores, vectors
182
183
184   def compute_precision(truth, search_type, dataframe):
185       precision, counter = 0, 0
186       for _, row in dataframe[dataframe['search'] == search_type].iterrows():
187           if row["name"] == truth.name and row["file"] == truth.path:
188               return 1 / (counter + 1)
189           counter += 1
190       return precision
191
192
193   def compute_recall(precision):
194       return 1 if precision > 0 else 0
195
```

```python
196
197  def compute_tsne(dictionary):
198      results = {}
199      for key, values in dictionary.items():
200          tsne = TSNE(n_components=2, verbose=1, perplexity=2, n_iter=3000)
201          results[key] = tsne.fit_transform(values)
202      return results
203
204
205  def plot_vectors(dictionary, queries):
206      for key, values in dictionary.items():
207          dataframe = pd.DataFrame()
208          dataframe['x'] = values[:, 0]
209          dataframe['y'] = values[:, 1]
210          plt.figure(figsize=(16, 16))
211          plt.title("Results of " + key)
212
213          sns_plot = sns.scatterplot(
214              x="x",
215              y="y",
216              hue=queries + list(itertools.chain.from_iterable([query] * 5 for query in queries
                      )),
217              data=dataframe,
218              legend="full",
219              alpha=1.0
220          )
221          sns_plot.get_figure().savefig("res/plot_" + key.lower())
222
223
224  def print_scores(scores):
225      print("##### PRINT #####")
226      for key, values in scores.items():
227          print(key)
228          precision, recall = compute_mean(values)
229          print("\tprecision:\t" + precision)
230          print("\trecall:\t\t" + recall)
231
232
233  def compute_mean(stats):
234      precision, recall, counter = 0, 0, 0
235      for stat in stats:
236          precision += stat.precisions
237          recall += stat.recalls
238          counter += 1
239      return str(precision / counter), str(recall / counter)
240
241
242  if len(argv) < 1:
243      print("Please give as input ground truth file")
244      exit(1)
245
246
247  begin_time = datetime.now()
248  start(argv[1])
249  print(datetime.now() - begin_time)
```

# B    Bash Code

```bash
#!/bin/bash

python3 src/extract_data.py $1
python3 src/search_data.py $2
python3 src/prec_recall.py res/data.csv res/ground-truth.txt
```