# Contents

# 1 README.txt

```
1  ##########################################################################
2  #                                                                        #
3  #                   README FILE for IML 2019 Hackaton task               #
4  #                                    Tweets                              #
5  #                                                                        #
6  #                 Tal Porezky,   tal.poreky, 311322499                   #
7  #                 Ilya Merkulov, iliamark,   317559631                   #
8  #                 Yoav Cohen,    cyoavc,     307944017                   #
9  #                 Amit Weiss,    amit.weiss, 304873334                   #
10 #                                06.06.2019                              #
11 #                                                                        #
12 ##########################################################################
13
14
15 classifier.py - Contains our implementation of the classifier.
16 requirements.txt - Contains the requirements for the virtualenv.
17 countries.txt - contains a list of all the countries in the universe.
18 README - THIS VERY FILE!
19 USERS.txt - we!
20 project.pdf - The description of our choices in design.
21
22 The following files are binary files for the classifier:
23 vcTags.pickle
24 vcQuestionMark.pickle
25 vcDotsMark.pickle
26 vc.pickle
27 train_avg_len.pickle
28 vcExclamationMark.pickle
29 best_classifier_ever.pickle
30
```

# 2 USERS.txt

```
1    tal.porezky, 311322499
2    iliamark, 317559631
3    cyoavc, 307944017
4    amit.weiss, 304873334
```

# IML Hackaton

Cohen Yoav
Porezky Tal
Merkulov Ilya
Weiss Amit

June 7, 2019

Shortly after we started working on our classifier we reasoned that our main problem is the feature selection. Both choosing them and deciding whether they are binary or not. The choice was often against basic logic or intuition, was made to the appeasing of the test data.

- Our first step was to use the bag of words to maximum results by making use of the inherent methods of the CountVectorizer class such as stop words, both English and Portuguese.

- It is worth noting, that the first step already achieved 70% success rate and from there, we managed to pull a few minor tweaks to raise the success rate by approximately 10 additional percents. Among these 'tweaks':

    - Length of each tweet
    - Tags - We counted the number of appearances of the symbol '@' representing the amount of tags a tweet is containing and the tendency of a user to tag other people in posts.
    - Countries - We examined whether a word from the file countries.txt, a text file containing every country, appear in a tweet.
    - Marks - Various marks and their combinations were examined. Among them: '!' ',' '?' and the regex ...* which detected a continuous use of 2 or more periods.

- Binary - Some features, much to our surprise, proved to be more useful in a binary score function. Describing whether a country was mentioned in a tweet instead of counting the number of countries, raised the success rate by 3%. However, in others e.g tag, no improvement was achieved.

- Correlation - We have also tried adding features describing correlation between length and mark or tag. These proved to be inconsequential.

- We used various classifiers of sklearn and found the best results with MultinomialNB classifier. We have tried to use the nltk library to process the effects of noun adjectives etcetra and emotions core but with no avail!

1

# 4 task1/src/classifier.py

```python
#!/usr/bin/env python
# coding: utf-8

import numpy as np
from sklearn.feature_extraction.text import TfidfTransformer, CountVectorizer
import pandas as pd
import os
from scipy.sparse import hstack
from sklearn.naive_bayes import MultinomialNB
import pickle

# The name of the folder were the tweets are held.
TWEETS_FOLDER = "tweets"
COUNTRIES_FILE = "countries.txt"


def parse_data(data_folder):
    """
    Read all the csv files in the given folder and merge them to a single
    pandas dataframe. Note that this function drops the 'user' column and
    keeps the tweets owner id number at the head of each column.
    :param data_folder:
    :return:
    """
    tweets = list()
    for filename in os.listdir(data_folder):
        # Only choose files with CSV ending and ignore the test file.
        if filename.endswith(".csv") and filename != 'tweets_test_demo.csv':
            current_tweet = pd.read_csv(TWEETS_FOLDER + '/' + filename)
            tweets.append(current_tweet)
            continue
        else:
            continue

    # Join all the tweets to a single pandas array.
    tweets = pd.concat(tweets, sort=True, axis=0)
    return tweets


def len_feature(tweets):
    feature_vec = np.zeros((len(tweets), 1))
    for i, tweet in enumerate(tweets):
        feature_vec[i] = len(tweet.split())
    return feature_vec


def country_parse():
    """
    Generates an np.array of all the countries.
    :param countries_file: The file of countries location.
    :return: An array of all countries in lower case!
    """
    f = open(COUNTRIES_FILE, "r")
    countries = list()
    for country in f:
        countries.append(country.replace('\n', '').lower())
    np.array(countries)
    return countries
```

```python
60
61  def countries_tweets_count(tweets):
62      """
63      Receives an array of tweets and counts the appearances of countries
64      :param tweets: an array of tweets
65      :return: an array of integers
66      """
67      country_count = np.zeros((len(tweets), 1))
68      for i, tweet in enumerate(tweets):
69          tweet_country_count = 0
70          for country in country_parse():
71              tweet_country_count = tweet_country_count + tweet.lower().count(country)
72          country_count[i] = tweet_country_count
73      for i in range(len(country_count)):
74          if country_count[i] != 0:
75              country_count[i] = 1
76      return country_count
77
78
79  def find_tags(tweets):
80      """
81      Receives an array of tweets and counts the appearances of countries
82      :param tweets: an array of tweets
83      :return: an array of integers
84      """
85      tag_count = np.zeros((len(tweets), 1))
86      for i, tweet in enumerate(tweets):
87          if tweet.lower().count('@') > 2:
88              tag_count[i] = 1
89      return tag_count
90
91
92  def classify(tweets):
93      """
94      Classifies tweets to various twitter stars.
95      :param tweets:
96      :return:
97      """
98
99      vc_file = open("vc.pickle", "rb")
100     vc = pickle.load(vc_file)
101     vc_file.close()
102
103     vcTags_file = open("vcTags.pickle", "rb")
104     vcTags = pickle.load(vcTags_file)
105     vcTags_file.close()
106
107     vcExclamationMark_file = open("vcExclamationMark.pickle", "rb")
108     vcExclamationMark = pickle.load(vcExclamationMark_file)
109     vcExclamationMark_file.close()
110
111     vcQuestionMark_file = open("vcQuestionMark.pickle", "rb")
112     vcQuestionMark = pickle.load(vcQuestionMark_file)
113     vcQuestionMark_file.close()
114
115     vcDotsMark_file = open("vcDotsMark.pickle", "rb")
116     vcDotsMark = pickle.load(vcDotsMark_file)
117     vcDotsMark_file.close()
118
119     train_avg_len_file = open("train_avg_len.pickle", "rb")
120     train_avg_len = pickle.load(train_avg_len_file)
121     train_avg_len_file.close()
122
123
124     print("HEY")
125
126     tweets = tweets['tweet']
127
```

```python
128          print(tweets)
129          print(tweets.shape)
130
131          # VC transform
132          hist = vc.transform(tweets)
133          print("hist shape is " + str(hist.shape))
134
135          # Lengths
136          test_tweets_lengths = len_feature(tweets)
137          test_tweets_lengths = test_tweets_lengths / train_avg_len
138
139          # Find all tags
140          tags = vcTags.transform(tweets)
141
142          # Find all marks
143          exclamationMarks = np.sum(vcExclamationMark.transform(tweets), axis=1)
144
145          questionMarks = np.sum(vcQuestionMark.transform(tweets), axis=1)
146
147          dotsMarks = np.sum(vcDotsMark.transform(tweets), axis=1)
148
149          countries = countries_tweets_count(tweets)
150
151          tweets = hstack([hist, tags, exclamationMarks, questionMarks, dotsMarks])
152          print(test_tweets_lengths)
153          print(tweets.shape)
154          print(test_tweets_lengths.shape)
155          tweets = hstack([tweets, test_tweets_lengths])
156          tweets = hstack([tweets, countries])
157
158
159
160          classifier_file = open("best_classifier_ever.pickle", "rb")
161          classifier = pickle.load(classifier_file)
162          classifier_file.close()
163
164          print(classifier.predict(tweets))
165
166          return classifier.predict(tweets)
167
168
169
170   def train_it():
171          # generate train and test data
172          data = parse_data(TWEETS_FOLDER)
173          test_size = 1
174          test_data = data.sample(test_size)
175          train_data = data.drop(test_data.index)
176
177          #####################
178          # Generate Features #
179          #####################
180
181          # Transform all tweets to matrix of features per tweet
182          vc = CountVectorizer(ngram_range=(1, 2), stop_words='english')
183          hist = vc.fit_transform(train_data['tweet'])
184
185          # Generate X length feature and join it with the previous features
186          train_tweets_lengths = len_feature(train_data['tweet'])
187          train_avg_len = sum(train_tweets_lengths) / len(train_tweets_lengths)
188          train_tweets_lengths = train_tweets_lengths / train_avg_len
189
190          # Generate the length feature of the train set
191          test_tweets_lengths = len_feature(test_data['tweet'])
192          test_tweets_lengths = test_tweets_lengths / train_avg_len
193
194          # Find all tags
195          vcTags = CountVectorizer(token_pattern=r'\b@\b')
```

```python
196        # tags = find_tags(train_data['tweet'])
197        tags = vcTags.fit_transform(train_data['tweet'])
198
199        # Find all marks
200        vcExclamationMark = CountVectorizer(token_pattern=r'(?u)(?:[!]+)')
201        exclamationMarks = np.sum(vcExclamationMark.fit_transform(train_data['tweet']), axis=1)
202
203        vcQuestionMark = CountVectorizer(token_pattern=r'(?u)(?:[?]+)')
204        questionMarks = np.sum(vcQuestionMark.fit_transform(train_data['tweet']), axis=1)
205
206        vcDotsMark = CountVectorizer(token_pattern=r'(?u)(?:[.]{2,})')
207        dotsMarks = np.sum(vcDotsMark.fit_transform(train_data['tweet']), axis=1)
208
209        # Find Countries
210        countries = countries_tweets_count(train_data['tweet'])
211
212        features = hstack([hist, tags, exclamationMarks, questionMarks, dotsMarks])
213        features = hstack([features, train_tweets_lengths])
214        features = hstack([features, countries])
215
216        # Generate Test Features
217        hist_test = vc.transform(test_data['tweet'])
218        # tags_test = find_tags(test_data['tweet'])
219        tags_test = vcTags.transform(test_data['tweet'])
220        exclamation_marks_test = np.sum(vcExclamationMark.transform(test_data['tweet']), axis=1)
221        question_marks_test = np.sum(vcQuestionMark.transform(test_data['tweet']), axis=1)
222        dots_marks = np.sum(vcDotsMark.transform(test_data['tweet']), axis=1)
223        countries_test = countries_tweets_count(test_data['tweet'])
224
225        test_features = hstack([hist_test, tags_test, exclamation_marks_test, question_marks_test, dots_marks])
226        test_features = hstack([test_features, test_tweets_lengths])
227        test_features = hstack([test_features, countries_test])
228        # test_features = hstack([test_features, tags_test])
229
230
231        # Transform features
232        transformer = TfidfTransformer()
233        freq_all = transformer.fit_transform(features)
234        freq_test = transformer.fit_transform(test_features)
235
236        clf2 = MultinomialNB().fit(freq_all, train_data['user'])
237        # predicted2 = clf2.predict(freq_test)
238
239        # Success of prediction on test_data.
240        # print(len(predicted2[predicted2 == test_data['user']]) / test_size)
241
242
243        save_classifier = open("best_classifier_ever.pickle", "wb")
244        pickle.dump(clf2, save_classifier)
245        save_classifier.close()
246
247        save_vc = open("vc.pickle", "wb")
248        pickle.dump(vc, save_vc)
249        save_vc.close()
250
251        save_vcTags = open("vcTags.pickle", "wb")
252        pickle.dump(vcTags, save_vcTags)
253        save_vcTags.close()
254
255        save_vcExclamationMark = open("vcExclamationMark.pickle", "wb")
256        pickle.dump(vcExclamationMark, save_vcExclamationMark)
257        save_vcExclamationMark.close()
258
259        save_vcQuestionMark = open("vcQuestionMark.pickle", "wb")
260        pickle.dump(vcQuestionMark, save_vcQuestionMark)
261        save_vcQuestionMark.close()
262
263        save_vcDotsMark = open("vcDotsMark.pickle", "wb")
```

```python
264         pickle.dump(vcDotsMark, save_vcDotsMark)
265         save_vcDotsMark.close()
266
267         save_train_avg = open("train_avg_len.pickle", "wb")
268         pickle.dump(train_avg_len, save_train_avg)
269         save_train_avg.close()
270
271
272  if __name__ == '__main__':
273
274
275         file ="tweets_test_demo.csv"
276         current_tweet = pd.read_csv(file)
277
278         classify(current_tweet)
```

# 5 task1/src/countries.txt

```
1   Afghanistan
2   Albania
3   Algeria
4   Andorra
5   Angola
6   Antigua
7   Argentina
8   Armenia
9   Australia
10  Austria
11  Azerbaijan
12  Bahamas
13  Bahrain
14  Bangladesh
15  Barbados
16  Belarus
17  Belgium
18  Belize
19  Benin
20  Bhutan
21  Bolivia
22  Bosnia Herzegovina
23  Botswana
24  Brazil
25  Brunei
26  Bulgaria
27  Burkina
28  Burundi
29  Cambodia
30  Cameroon
31  Canada
32  Cape Verde
33  Central African Rep
34  Chad
35  Chile
36  China
37  Colombia
38  Comoros
39  Congo
40  Costa Rica
41  Croatia
42  Cuba
43  Cyprus
44  Czech Republic
45  Denmark
46  Djibouti
47  Dominica
48  Dominican Republic
49  East Timor
50  Ecuador
51  Egypt
52  El Salvador
53  Equatorial Guinea
54  Eritrea
55  Estonia
56  Ethiopia
57  Fiji
58  Finland
59  France
```

| | |
|---|---|
| 60 | Gabon |
| 61 | Gambia |
| 62 | Georgia |
| 63 | Germany |
| 64 | Ghana |
| 65 | Greece |
| 66 | Grenada |
| 67 | Guatemala |
| 68 | Guinea |
| 69 | Guinea-Bissau |
| 70 | Guyana |
| 71 | Haiti |
| 72 | Honduras |
| 73 | Hungary |
| 74 | Iceland |
| 75 | India |
| 76 | Indonesia |
| 77 | Iran |
| 78 | Iraq |
| 79 | Ireland |
| 80 | Israel |
| 81 | Italy |
| 82 | Ivory Coast |
| 83 | Jamaica |
| 84 | Japan |
| 85 | Jordan |
| 86 | Kazakhstan |
| 87 | Kenya |
| 88 | Kiribati |
| 89 | Korea North |
| 90 | Korea South |
| 91 | Kosovo |
| 92 | Kuwait |
| 93 | Kyrgyzstan |
| 94 | Laos |
| 95 | Latvia |
| 96 | Lebanon |
| 97 | Lesotho |
| 98 | Liberia |
| 99 | Libya |
| 100 | Liechtenstein |
| 101 | Lithuania |
| 102 | Luxembourg |
| 103 | Macedonia |
| 104 | Madagascar |
| 105 | Malawi |
| 106 | Malaysia |
| 107 | Maldives |
| 108 | Mali |
| 109 | Malta |
| 110 | Marshall Islands |
| 111 | Mauritania |
| 112 | Mauritius |
| 113 | Mexico |
| 114 | Micronesia |
| 115 | Moldova |
| 116 | Monaco |
| 117 | Mongolia |
| 118 | Montenegro |
| 119 | Morocco |
| 120 | Mozambique |
| 121 | Myanmar |
| 122 | Namibia |
| 123 | Nauru |
| 124 | Nepal |
| 125 | Netherlands |
| 126 | New Zealand |
| 127 | Nicaragua |

```
128    Niger
129    Nigeria
130    Norway
131    Oman
132    Pakistan
133    Palau
134    Panama
135    Papua New Guinea
136    Paraguay
137    Peru
138    Philippines
139    Poland
140    Portugal
141    Qatar
142    Romania
143    Russian Federation
144    Rwanda
145    St Kitts & Nevis
146    St Lucia
147    Saint Vincent & the Grenadines
148    Samoa
149    San Marino
150    Sao Tome & Principe
151    Saudi Arabia
152    Senegal
153    Serbia
154    Seychelles
155    Sierra Leone
156    Singapore
157    Slovakia
158    Slovenia
159    Solomon Islands
160    Somalia
161    South Africa
162    South Sudan
163    Spain
164    Sri Lanka
165    Sudan
166    Suriname
167    Swaziland
168    Sweden
169    Switzerland
170    Syria
171    Taiwan
172    Tajikistan
173    Tanzania
174    Thailand
175    Togo
176    Tonga
177    Trinidad & Tobago
178    Tunisia
179    Turkey
180    Turkmenistan
181    Tuvalu
182    Uganda
183    Ukraine
184    United Arab Emirates
185    United Kingdom
186    United States
187    Uruguay
188    Uzbekistan
189    Vanuatu
190    Vatican City
191    Venezuela
192    Vietnam
193    Yemen
194    Zambia
195    Zimbabwe
```

# 6 task1/src/requirements.txt

```
1   alabaster==0.7.8
2   antlr==2.7.5rc1
3   apsw==3.16.2.post1
4   apt-xapian-index==0.49
5   arandr==0.1.9
6   asn1crypto==0.24.0
7   astroid==1.6.5
8   atomicwrites==1.1.5
9   attrs==17.4.0
10  autobahn==17.10.1
11  Automat==0.6.0
12  avro==1.8.2
13  Axiom==0.7.5
14  Babel==2.4.0
15  backports-abc==0.5
16  backports.functools-lru-cache==1.5
17  backports.shutil-get-terminal-size==1.0.0
18  backports.ssl-match-hostname==3.5.0.1
19  basemap==1.1.0
20  bcrypt==3.1.4
21  Beaker==1.9.1
22  BeautifulSoup==3.2.1
23  beautifulsoup4==4.6.1
24  biom-format==2.1.6
25  biopython==1.72
26  bleach==2.1.3
27  blinker==1.4
28  bottle==0.12.13
29  brial==1.0.1
30  Brlapi==0.6.7
31  bx-python==0.8.1
32  bz2file==0.98
33  bzr==2.8.0.dev1
34  BzrTools==2.6.0
35  CacheControl==0.11.7
36  cajarename==18.1.10
37  cbor==1.0.0
38  ccsm==0.9.13.1
39  certifi==2018.4.16
40  cffi==1.11.5
41  Chameleon==2.24
42  chardet==3.0.4
43  Cheetah3==3.1.0
44  CherryPy==8.9.1
45  click==6.7
46  cloudpickle==0.5.2
47  Coherence==0.6.6.2
48  colorama==0.3.7
49  compizconfig-python==0.9.13.1
50  configobj==5.0.6
51  configparser==3.5.0
52  constantly==15.1.0
53  cookies==2.2.1
54  cryptography==2.3
55  cssselect==1.0.3
56  cssutils==1.0.2
57  cutadapt==1.18
58  cvxopt==1.1.9
59  cwltool==1.0.20180302231433
```

```
60   cycler==0.10.0
61   cypari2==1.1.4
62   Cython==0.28.4
63   dap==2.2.6.7
64   dblatex==0.3.10
65   decorator==4.3.0
66   defer==1.0.6
67   defusedxml==0.5.0
68   deluge==1.3.15
69   devscripts==2.11.2
70   dissy==9
71   distro==1.0.1
72   Django==1.11.15
73   django-ajax-selects==1.7.0
74   django-app-plugins==0.1.1
75   django-auth-ldap==1.4.0
76   django-dajax==0.9.2
77   django-dajaxice==0.7
78   django-extensions==1.8.1
79   django-piston==0.2.3
80   django-websocket==0.3.0
81   Djblets==0.7a0.dev0
82   dnspython==1.15.0
83   docutils==0.14
84   drmaa==0.5
85   dulwich==0.19.6
86   elib.intl==0.0.3.dev0
87   entrypoints==0.2.3.post2
88   enum34==1.1.6
89   Epsilon==0.7.1
90   et-xmlfile==1.0.1
91   fastimport==0.9.8
92   feedparser==5.2.1
93   Flask==1.0.2
94   Flask-AutoIndex==0.6.1
95   Flask-Babel==0.11.2
96   Flask-OldSessions==0.10
97   Flask-OpenID==1.2.5
98   Flask-Silk==0.2
99   foolscap==0.13.1
100  FormEncode==1.3.0
101  fpconst==0.7.2
102  fpylll==0.3.0.dev0
103  funcsigs==1.0.2
104  functools32==3.2.3.post2
105  future==0.15.2
106  futures==3.2.0
107  GDAL==2.3.1
108  gdata==2.0.18
109  gdmodule==0.59
110  glpk==0.4.52
111  gmpy==1.17
112  GnuPGInterface==0.3.2
113  gyp==0.1
114  h5py==2.8.0
115  html5-parser==0.4.4
116  html5lib==1.0.1
117  HTSeq===0.6.1p1
118  httplib2==0.9.2
119  hupper==1.2
120  hyperlink==17.3.1
121  idna==2.6
122  imageio==2.2.0
123  imagesize==1.0.0
124  incremental==16.10.1
125  ipaddr==2.2.0
126  ipaddress==1.0.17
127  IPy==0.83
```

```
128  ipykernel==4.8.2
129  ipython==5.5.0
130  ipython-genutils==0.2.0
131  ipywidgets==6.0.0
132  isodate==0.6.0
133  isort==4.3.4
134  itsdangerous==0.24
135  jdcal==1.0
136  jedi==0.12.0
137  Jinja2==2.10
138  joblib==0.12.1
139  jsonschema==2.6.0
140  jupyter-client==5.2.3
141  jupyter-core==4.4.0
142  keyring==13.1.0
143  keyrings.alt==3.0
144  kiwisolver==1.0.1
145  launchpadlib==1.10.6
146  lazr.restfulclient==0.14.0
147  lazr.uri==1.0.3
148  lazy-object-proxy==1.3.1
149  leveldb==0.1
150  libvirt-python==4.5.0
151  lockfile==0.12.2
152  logilab-common==1.4.1
153  Louie==1.1
154  louis==3.6.0
155  lxml==4.2.3
156  lz4==1.1.0
157  M2Crypto==0.27.0
158  Mako==1.0.7
159  Markdown==2.6.9
160  MarkupSafe==1.0
161  mate-menu==18.4.3
162  matplotlib==2.2.2
163  mccabe==0.6.1
164  mechanize==0.2.5
165  mercurial==4.7
166  mistune==0.8.3
167  mock==2.0.0
168  mod-python==3.3.1
169  moin==1.9.9
170  more-itertools==4.2.0
171  mpi4py==2.0.0
172  mpmath==1.0.0
173  msgpack==0.5.6
174  mysql-connector-python==8.0.11
175  mysql-utilities==1.6.4
176  mysqlclient==1.3.10
177  nbconvert==5.3.1
178  nbformat==4.4.0
179  netifaces==0.10.4
180  networkx==2.1
181  Nevow==0.14.2
182  nltk==3.3
183  nose==1.3.7
184  notebook==5.4.1
185  numexpr==2.6.5
186  numpy==1.16.4
187  numpydoc==0.7.0
188  oauth==1.0.1
189  oauthlib==2.0.6
190  olefile==0.45.1
191  openpyxl==2.4.9
192  packaging==17.1
193  PAM==0.4.2
194  pandas==0.23.3
195  pandocfilters==1.4.2
```

```
196    parameterized==0.6.1
197    paramiko==2.4.0
198    parsedatetime==2.4
199    parso==0.2.1
200    passlib==1.7.1
201    Paste==2.0.3
202    PasteDeploy==1.5.2
203    PasteScript==2.0.2
204    path.py==11.0.1
205    pathlib2==2.3.2
206    patsy==0.4.1+dev
207    pbr==3.1.1
208    pdfrw==0.4
209    pexpect==4.6.0
210    pickleshare==0.7.4
211    Pillow==5.2.0
212    pkgconfig==1.3.1
213    plaster==1.0
214    plaster-pastedeploy==0.5
215    pluggy==0.6.0
216    ply==3.11
217    Pmw==1.3.2
218    progressbar==2.3
219    prompt-toolkit==1.0.15
220    protobuf==3.0.0
221    psphere==0.5.2
222    psutil==5.4.6
223    psycopg2==2.7.5
224    py==1.5.4
225    py-ubjson==0.8.5
226    pyasn1==0.4.2
227    pyasn1-modules==0.2.1
228    PyAudio==0.2.11
229    pycairo==1.16.2
230    pycodestyle==2.3.1
231    pycparser==2.18
232    pycrypto==2.6.1
233    pycups==1.9.73
234    pycurl==7.43.0.1
235    pydot==1.2.4
236    pyenchant==2.0.0
237    pyExcelerator==0.6.4.1
238    pyflakes==1.6.0
239    pyfribidi==0.11.0
240    pygame===1.9.1release
241    pyglet==1.3.0
242    Pygments==2.2.0
243    pygobject==3.28.2
244    pygpgme==0.3
245    pygpu==0.6.9
246    pygraphviz==1.4rc1
247    PyGreSQL==5.0.5
248    PyICU==2.0.3
249    pyinotify==0.9.6
250    PyJWT==1.6.4
251    pyliblzma==0.5.3
252    pylibssh2==1.0.0
253    pylint==1.9.2
254    pymol==2.1.0
255    PyMySQL==0.8.1
256    PyNaCl==1.2.1
257    pyodbc==4.0.22
258    PyOpenGL==3.1.0
259    pyOpenSSL==18.0.0
260    pyparsing==2.2.0
261    pyproj==1.9.5.1
262    pyqi==0.3.2
263    pyquery==1.2.9
```

```
264    pyramid==1.9.1
265    PyRRD==0.1.0
266    pysam==0.14
267    pysaml2==4.0.2
268    pyserial==3.4
269    pyshp==1.2.12
270    pysqlite==2.7.0
271    Pyste==0.9.10
272    pytest==3.6.4
273    python-apt==1.6.2
274    python-dateutil==2.6.1
275    python-debian==0.1.32
276    python-debianbts==2.7.2
277    python-distutils-extra==2.39
278    python-gflags==1.5.1
279    python-gnutls==3.0.0
280    python-ldap==3.1.0
281    python-libtorrent==1.1.5
282    python-lzo==1.12
283    python-magic==0.4.16
284    python-memcached==1.53
285    python-openid==2.2.5
286    python-snappy==0.5.3
287    python-xlib==0.23
288    PyTrie==0.2
289    pytz==2018.5
290    PyWavelets==0.5.1
291    pyxdg==0.25
292    PyYAML==3.12
293    pyzmq==17.1.0
294    qcli==0.1.1
295    qrcode==6.0
296    QtAwesome==0.4.4
297    qtconsole==4.3.1
298    QtPy==1.3.1
299    rabbitvcs==0.16.0
300    radiotray==0.7.3
301    rdflib==4.2.2
302    rdflib-jsonld==0.4.0
303    recaptcha-client==1.0.6
304    regex==2018.6.9
305    reportbug==6.6.6
306    reportlab==3.5.2
307    repoze.lru==0.7
308    repoze.who==2.2
309    requests==2.18.4
310    responses==0.9.0
311    roman==2.0.0
312    rope==0.10.5
313    Routes==2.4.1
314    rpm==4.14.1
315    rpy2==2.8.6
316    rrdtool==0.1.10
317    ## FIXME: could not find svn URL in dependency_links for this package:
318    rst2pdf===0.93.dev-r0
319    ruamel.yaml==0.15.34
320    rubber==1.4
321    sage==8.2
322    sagenb==1.0.1
323    sagenb-export==3.2
324    scandir==1.8
325    scgi==1.13
326    schema-salad==2.6.20171201034858
327    scikit-image==0.14.0
328    scikit-learn==0.20.3
329    scipy==1.2.1
330    SecretStorage==2.3.1
331    Send2Trash==1.4.2
```

```
332   service-identity==16.0.0
333   setproctitle==1.1.10
334   setuptools-scm==3.0.6
335   shellescape==3.4.1
336   simplegeneric==0.8.1
337   simplejson==3.15.0
338   singledispatch==3.4.0.3
339   six==1.11.0
340   sklearn==0.0
341   snowballstemmer==1.2.1
342   SOAPpy==0.12.22
343   SPARQLWrapper==1.7.6
344   Sphinx==1.7.6
345   sphinx-rtd-theme==0.4.0
346   spyder==3.2.6
347   SQLAlchemy==1.2.8
348   sqlparse==0.2.4
349   statsmodels==0.8.0
350   stevedore==1.28.0
351   subprocess32==3.5.2
352   suds-jurko==0.7.dev0
353   sympy==1.2
354   tables==3.4.4
355   tabulate==0.8.2
356   tagpy==2013.1
357   Tempita==0.5.2
358   termcolor==1.1.0
359   terminado==0.8.1
360   testpath==0.3.1
361   Theano==1.0.2
362   tornado==5.0.2
363   tqdm==4.23.4
364   traitlets==4.3.2
365   translationstring==1.3
366   trollius==2.0.1
367   Twisted==18.4.0
368   Twisted-Conch==11.0.0
369   Twisted-Core==11.0.0
370   Twisted-Names==11.0.0
371   Twisted-Web==11.0.0
372   txaio==2.8.1
373   typing==3.6.4
374   u-msgpack-python==2.1
375   urlgrabber==3.10.2
376   urllib3==1.22
377   uTidylib==0.3
378   vcversioner==2.16.0.0
379   venusian==1.1.0
380   virtualenv==15.1.0
381   virtualenv-clone==0.3.0
382   virtualenvwrapper==4.3.1
383   wadllib==1.3.3
384   waitress==1.1.0
385   wcwidth==0.1.7
386   webencodings==0.5
387   WebOb==1.7.3
388   WebTest==2.0.28
389   Werkzeug==0.14.1
390   wrapt==1.9.0
391   wsaccel==0.6.2
392   wstools==0.4.3
393   wxPython==2.8.12.1
394   wxPython-common==3.0.2.0
395   xappy==0.5
396   xlrd==1.1.0
397   XlsxWriter==0.9.6
398   xlwt==1.3.0
399   xmltodict==0.11.0
```

```
400    xopen==0.3.3
401    yum-metadata-parser==1.1.4
402    zope.component==4.3.0
403    zope.configuration==4.0.3
404    zope.deprecation==4.1.2
405    zope.event==4.2.0
406    zope.hookable==4.0.4
407    zope.i18nmessageid==4.0.3
408    zope.interface==4.3.2
409    zope.schema==4.4.2
410    zope.testing==4.6.2
```

# Index of comments

4.1    Seems like you've done some deep research to decide which features to extract - well done!

It could be nice to see some of the work you described in a more visualized fashion, and also to extend some your explanation about the model you have decided to use.

You showed creativity and understanding of different methods you have studied in the course - good job!