

Final Project

Personality Recognition Based on Text Input

Introduction

In human-computer interaction, the need for adaptive systems is becoming increasingly important. This is especially true when interactive systems fill more meaningful roles. For example, in systems where users interact socially with a virtual agent (chatbot, embodied conversational agent, avatar, etc.), the system needs to be able to both recognize the user to understand the input (text, speech, gestures, etc.) and also generate relevant responses (in the form of text, speech, gestures, etc.).

One approach to understanding the user and generating relevant responses is to recognize the user's personality type. A common personality type classification is the Meyers-Briggs Type Indicator (MBTI), which measures personality on 4 dimensions found in Jungian psychology:

- *extraversion/introversion*, which is known commonly as social energy,
- *sensing/intuition*, which describes how an individual focuses on information,
- *feeling/thinking*, which describes how people make decisions, and,
- *perceiving/judging*, which describes how people deal with the outside world.

Previous Work

Previously, researchers at [Stanford](#) successfully implemented a personality type recognizer model that achieved up to 80% accuracy. However, a shortfall of their work lies in the predictor models utilized. In their work, the authors presented several predictor models, though notably the logistic regression and support vector machine models produced better results. In this work, I show how the use of a more appropriate predictor model can achieve better results with less data pre-processing.

The authors at [Stanford](#) performed a lot of data pre-processing. In fact, they (strangely) used both stemming and lemmatization to train the Doc2Vec model for embedding representations of the articles (or documents) in the dataset. They also performed stopwords removal, assigned all letters to lowercase versions, and removed all documents with URLs. In total, before pre-processing, there are 422,844 articles:

- 37463 INTJ samples
- 44456 INTP samples
- 51999 INFJ samples
- 64612 INFP samples
- 6726 ISTJ samples
- 11580 ISTP samples
- 5807 ISFJ samples
- 9029 ISFP samples
- 7774 ENTJ samples
- 24997 ENTP samples
- 6647 ENFJ samples
- 22976 ENFP samples
- 1495 ESTJ samples
- 3328 ESTP samples
- 1489 ESFJ samples
- 1531 ESFP samples

After their pre-processing, the data was reduced by 6% to contain 397,613 samples:

- 35256 INTJ samples
- 41684 INTP samples
- 48860 INFJ samples
- 60329 INFP samples
- 6369 ISTJ samples
- 10794 ISTP samples
- 5446 ISFJ samples
- 8265 ISFP samples
- 7432 ENTJ samples
- 23900 ENTP samples
- 6288 ENFJ samples
- 21956 ENFP samples
- 1461 ESTJ samples
- 3131 ESTP samples
- 1460 ESFJ samples
- 1459 ESFP samples

Approach

Similar to the previous work, my reported models are individual models for each of the personality dichotomies (E/I, S/N, F/T, P/J). Therefore, I trained one model for each binary class (4 models). Each predictor model did use the same Doc2Vec trained model based on the training data input (the text articles). I chose the use of a Doc2Vec model over a Word2Vec model as the model should be making predictions on a document-level rather than word-level basis. Contrary to the baseline model, I did not remove those articles which contained URLs, instead I removed the URL links themselves from the articles. So, in total, there are still 422,844 articles, and I split

the data (both the baseline and my pre-processed) into a 75/25 ratio for training and test data, respectively.

Results

My results (average of 10 test runs) are reported in the following tables, for both the baseline (previous work logistic regression) model, as well as my two Naive-Bayes models (Gaussian and Multinomial). Table 1 shows the results when training a Doc2Vec model with *vector_size* = 15, and Table 2 shows the results when training a Doc2Vec model with *vector_size* = 30.

I ran several trials of various combinations of parameters for the sklearn models. For both the Gaussian and Multinomial NB models, the *var_smoothing* (Gaussian) and *alpha* (Multinomial) parameters produced better results when they were set to 0.8. The probabilities of the class (*priors* in Gaussian NB and *class_priors* in the Multinomial NB) were found to produce better results (higher accuracy) when the probabilities were not balanced, but in the following ranges:

- E/I, in $[0.3, 0.4]$,
- S/N, in $[0.4, 0.5]$,
- F/T, in $[0.6, 0.7]$,
- P/J, in $[0.5, 0.6]$.

I attribute the results of my model (and the class probabilities) to the disproportionate amount of I/E, S/N, F/T, and P/J articles in the dataset.

Model	Accuracy (%)
Baseline	80.1
Gaussian NB	78.7
Multinomial NB	82.3

Table 1: Results: Embedding length=15

Model	Accuracy (%)
Baseline	81.6
Gaussian NB	83.6
Multinomial NB	84.4

Table 2: Results: Embedding length=30

Discussion

The results presented could be more improved had I been able to test more data on my machine/environment. In my system, I had to modify the page file sizing to allow for larger allocation for the PyCharm program to build on combined set sizes larger than 1500 (though at the expense of time).

One area for improvement in my solution is to report one F1-score measures to determine just how much my models were really learning about the data. In this way, I could analyze the results for their precision and recall. If the recall was low and the precision was high, then that is decent, but it would be ideal if both measures were high: my model should (ideally) learn more about introversion/extraversion (or the others) when making its predictions. For this, I'd utilize the `sklearn.metrics` library to report on these results for each of the classes.

Finally, another improvement to this work would be to develop Naive-Bayes predictor models on more-accepted personality indicators (the Big Five model, for example). Another avenue for future work would be to generate the responses that would match the personality type predicted. This would be especially useful in on-the-fly interactive systems for more immersive and human-like representations of virtual agent. For example, in an interactive system with an embodied conversational agent, it would be useful to identify a human interactant by their personality to be more assistive in learning environments, casual interactions, or medical scenarios.

Footnote

The source code of the described program is available at https://github.com/mahdafr/19w_cs5361-labs.

Previous work: <http://web.stanford.edu/class/archive/cs/cs221/cs221.1192/2018/restricted/posters/dkedia/poster.pdf>

```

1  import numpy as np
2  import pandas as pd
3  import os, math
4
5  # fixme: change to dataset location
6  d = "D:\\Google Drive\\skool\\CS 5361\\datasets\\project\\"
7  f = "mbti.csv"
8
9  class Dataset(object):
10     def __init__(self, first_time=False, lo=0.71, hi=0.77, to_load='', chop=0.0005):
11         self.X = None; self.Y = None
12         self.x = None; self.y = None
13         self.lo = lo; self.hi = hi
14         self.chop = chop
15         self.rand = np.random.randint(1,100000)
16         if first_time:
17             self._load()      # read from csv (for first run only)
18             self._save()      # save the data as npy for faster future runs
19         else:
20             self._read(to_load)    # read from npy
21
22     """ FIRST RUN ONLY Read the data to be saved for faster runs """
23     def _load(self):
24         data = np.asarray(pd.read_csv(d + f, sep=',', header=0))    # load from file
25         x, y = self._break(data)    # get the data into samples and target values
26         self._split(x, y)    # get training and test sets
27
28     """ Read the data to be stored """
29     def _read(self, to_load):
30         # find the npy files
31         lst = os.listdir(d)
32         # and ('baseline' not in s) and ('main' not in s) for non-spec use
33         file = [i for i, s in enumerate(lst) if (to_load+'npy' in s)]
34         # todo randomly load from list of test/train sets
35         self.X = np.load(d+lst[file[1]], allow_pickle=True)
36         self.Y = np.load(d+lst[file[3]], allow_pickle=True)
37         self.x = np.load(d+lst[file[0]], allow_pickle=True)
38         self.y = np.load(d+lst[file[2]], allow_pickle=True)
39         print('Loaded',self.X.shape[0],'train, ',self.x.shape[0],'test samples.')
40
41     """ Break the data into target values and training samples """
42     def _break(self, data):
43         X = []; Y = []
44         for i in range(data.shape[0]):
45             for s in np.array(data[i][1].split('|||')):
46                 X.append(s)    # sample
47                 Y.append(data[i][0])    # target value
48         return np.asarray(X), np.asarray(Y)
49
50     """ Split the data into training and test sets """
51     def _split(self, x, y):
52         # randomly choosing a value between lo%-hi% of dataset for train/test split
53         breakpoint = np.random.randint(self.lo*x.shape[0], self.hi*x.shape[0])
54         self.X = x[:breakpoint]
55         self.Y = y[:breakpoint]
56         # saving the test data
57         self.x = x[breakpoint+1:]
58         self.y = y[breakpoint+1:]
59         print('Using',self.X.shape[0],'train, ',self.x.shape[0],'test samples.')
60
61     """ Saves the training and test data into npy for efficient access """
62     def _save(self, name=''):
63         np.save(d+ 'X_' +str(self.X.shape[0]) + name,self.X)
64         np.save(d+ 'Y_' +str(self.X.shape[0]) + name,self.Y)
65         np.save(d+ 'x_' +str(self.x.shape[0]) + name,self.x)
66         np.save(d+ 'y_' +str(self.x.shape[0]) + name,self.y)
67         print('Saved npy files to ' +d+name)

```

```

68 """ Getter methods: testing and training data """
69
70 def train(self):
71     print('Using',math.ceil(self.chop*self.X.shape[0]),'training samples.')
72     self.Y = self._representation(self.Y)
73     self.X = self.X[self.rand+math.ceil(self.chop*self.X.shape[0]):]
74     self.Y = self.Y[self.rand+math.ceil(self.chop*self.Y.shape[0]):]
75     return self.X, self.Y
76
77 def test(self):
78     print('Using',math.floor(0.25*self.chop*self.X.shape[0]),'test samples.')
79     self.y = self._representation(self.y)
80     self.x = self.x[self.rand:self.rand+math.floor(0.25*self.chop*self.x.shape[0])]
81     self.y = self.y[self.rand:self.rand+math.floor(0.25*self.chop*self.y.shape[0])]
82     return self.x, self.y
83
84 """ Change the representation of target values. """
85 def _representation(self, y, one_hot=True):
86     if one_hot:
87         y = list(s.replace('E', '0') for s in y)
88         y = list(s.replace('I', '1') for s in y)
89         y = list(s.replace('S', '0') for s in y)
90         y = list(s.replace('N', '1') for s in y)
91         y = list(s.replace('F', '0') for s in y)
92         y = list(s.replace('T', '1') for s in y)
93         y = list(s.replace('P', '0') for s in y)
94         y = list(s.replace('J', '1') for s in y)
95     return np.array(y)
96
97 """ Setter method: testing and training data after pre-processing """
98 def set(self, X, Y, x, y, name=''):
99     self.X = X
100     self.Y = Y
101     self.x = x
102     self.y = y
103     self._save(name)
104
105 """ Get sample count of each personality/class type """
106 def of_each(self, in_strings=False):
107     if in_strings:
108         classes = ['INTJ', 'INTP', 'INFJ', 'INFP', 'ISTJ', 'ISTP', 'ISFJ',
109                   'ISFP', 'ENTJ', 'ENTP', 'ENFJ', 'ENFP', 'ESTJ', 'ESTP',
110                   'ESFJ', 'ESFP']
111     else:
112         classes = ['1111', '1110', '1101', '1100', '1011', '1010', '1001',
113                   '1000', '0111', '0110', '0101', '0100', '0011', '0010',
114                   '0001', '0000']
115     for c in classes:
116         print(str(np.sum(self.Y == c)), c, 'samples')
117
118 """ Get the target value for c """
119 def get_train_target(self, c):
120     return np.array([x[c] for x in self.Y])
121
122 def get_test_target(self, c):
123     return np.array([x[c] for x in self.y])
124

```

```

1  from gensim.models.doc2vec import Doc2Vec, TaggedDocument
2  import numpy as np
3
4  # fixme: change to dataset location
5  dr = "D:\\Google Drive\\skool\\CS 5361\\datasets\\project\\"
6
7  """ Get the embeddings for each training sample """
8  def train(d, title, first_time=False):
9      X, Y = d.train()
10     x, y = d.test()
11
12     if not first_time:
13         print('Using ' + title + ' model')
14         model = Doc2Vec.load(dr+title+"d2v.model", )
15         return vector(model, X), Y, vector(model, x), y
16
17     train = [TaggedDocument(words=list(X), tags=list(Y))]
18     np.save(dr + title + 'train', np.array(train))
19     test = [TaggedDocument(words=list(x), tags=list(y))]
20     np.save(dr + title + 'test', np.array(test))
21     data = [TaggedDocument(words=list(np.append(X,x)),
22                             tags=list(np.append(Y,y)))]
23     np.save(dr + title + 'data', np.array(data))
24     model = _model(title, data, epochs=2)
25     return vector(model, X), Y, vector(model, x), y
26
27  """ Train the model for a Doc2Vec embedding of input data """
28  def _model(title, tag, epochs=100, v=15, alpha=0.025):
29      # https://medium.com/@mishra.thedeepak/doc2vec-simple
30      # -implementation-example-df2afbbfbad5
31      model = Doc2Vec(size=v, alpha=alpha, min_alpha=0.00025,
32                      min_count=1, dm=1)
33      model.build_vocab(tag)
34
35      # print('Training ' + title + ' doc2vec model')
36      for epoch in range(epochs):
37          model.train(tag,
38                      total_examples=model.corpus_count,
39                      epochs=model.iter)
40          model.alpha -= 0.0002 # decrease the learning rate
41          model.min_alpha = model.alpha # fix to no decay
42          # print('Completed epoch', epoch)
43
44      print('Trained ' + title + ' doc2vec model with epochs=', str(epochs), 'vector_size=' +
45            str(v))
46      model.save(dr + title + "d2v_v=" + str(v) + ".model")
47      return model
48
49  """ Get the feature vector for the classifier """
50  def vector(model, doc):
51      # https://towardsdatascience.com/multi-class-text-classification
52      # -with-doc2vec-logistic-regression-9da9947b43f4
53      # y, x = zip(*[(d.tags[0],
54                     #     model.infer_vector(d.words, steps=20))
55                     # for d in doc])
56      # return x, y
57      ret = []
58      for i in range(len(doc)):
59          ret.append(model.infer_vector(list(doc[i])))
60      return np.array(ret)

```

```

1  import numpy as np
2  from nltk import WordNetLemmatizer
3  from nltk.stem import PorterStemmer
4  from nltk.corpus import stopwords
5  from nltk.tokenize import word_tokenize
6  from sklearn.linear_model import LogisticRegression
7  from project import dataset, doc2vec as emb
8
9  # fixme: change to dataset location
10 d = "D:\\Google Drive\\skool\\CS 5361\\datasets\\project\\"
11 title = 'baseline'
12
13 """ To preprocess the data """
14 def second_run(data):
15     data.of_each()
16     X, Y, x, y = _preprocess(*data.train(), *data.test())
17     data.set(X,Y,x,y,name='_'+title)
18     data.of_each()
19
20 """ One-time use: preprocess the data """
21 def _preprocess(X,Y,x,y):
22     X, Y = _urls_lower_lem_stem_stop(X,Y)
23     x, y = _urls_lower_lem_stem_stop(x,y)
24     return X, Y, x, y
25
26 """ Removes samples with URLs, lemmatizes, stems, then removes stopwords """
27 def _urls_lower_lem_stem_stop(X, Y):
28     nwX = []; nwY = []
29     ps = PorterStemmer()
30     lem = WordNetLemmatizer()
31     stop_words = set(stopwords.words('english'))
32
33     # for each sample
34     for i in range(X.shape[0]):
35         if 'http' not in X[i]:
36             filt = ps.stem(lem.lemmatize(X[i].lower()))
37             nwX.append(str([w for w in word_tokenize(filt) if not w in stop_words]))
38             nwY.append(Y[i])
39     return np.array(nwX), np.array(nwY)
40
41 """ Logistic regression """
42 def build_model(data):
43     X, Y, x, y = emb.train(data, title+'_', first_time=True)
44     model = LogisticRegression(solver='lbfgs', max_iter=100, multi_class='multinomial')
45     print('Classifier:\tLogistic Regression')
46     score = []
47     for i in range(4):
48         score.append(model.fit(X,data.get_train_target(i)).score(x,data.get_test_target(
49             i)))
50         print("\tTarget="+str(i), "Score:\t%f" % score[i])
51     print("Overall:",str(np.average(score)))
52
53 if __name__ == "__main__":
54     data = dataset.Dataset(to_load=title+'.', chop=0.05)
55     # second_run(data)
56     build_model(data)

```



```

1  import numpy as np
2  import re
3  from nltk.stem import PorterStemmer
4  from nltk.tokenize import word_tokenize
5  from sklearn.naive_bayes import GaussianNB, MultinomialNB
6  from project import dataset, doc2vec as emb
7  from sklearn.preprocessing import MinMaxScaler
8
9  # fixme: change to dataset location
10 d = "D:\\Google Drive\\skool\\CS 5361\\datasets\\project\\"
11 title = 'main'
12
13 """ To preprocess the data """
14 def second_run(data):
15     data.of_each()
16     X, Y, x, y = _preprocess(*data.train(), *data.test())
17     data.set(X,Y,x,y,name='_main')
18     data.of_each()
19
20 """ One-time use: preprocess the data """
21 def _preprocess(X,Y,x,y):
22     X, Y = _urls_lower_stem(X,Y)
23     x, y = _urls_lower_stem(x,y)
24     return X, Y, x, y
25
26 """ Removes samples with URLs, and makes samples lowercase """
27 def _urls_lower_stem(X, Y):
28     nwX = []; nwY = []
29     ps = PorterStemmer()
30     for i in range(X.shape[0]):
31         if 'http' in X[i]: # chop the URL from the string instead
32             X[i] = re.sub(
33                 r"(https?:\\/\\/)(\\s)*(www\\.)?(\\s)*((\\w|\\s)+\\.)*([\\w\\-\\s]+\\/)*([\\w\\-]+)((\\?)[\\w\\s]*=\\s*[\\w%&]*)*", '', X[i], re.MULTILINE)
34             nwX.append(str([w for w in word_tokenize(ps.stem(X[i].lower()))]))
35             nwY.append(Y[i])
36     return np.array(nwX), np.array(nwY)
37
38 """ Naive Bayes' Models """
39 def build_models(data):
40     X, Y, x, y = emb.train(data, title+'_', first_time=True)
41     # todo change the parameters for tests
42     pred(X, x, data, GaussianNB(), name='Gaussian')
43     pred(*rescale(X,x), data, MultinomialNB(class_prior=0.6), name='Multinomial')
44
45 """ Train and report on results """
46 def pred(X, x, data, model, name=''):
47     print('Classifier:', name)
48     score = []
49     for i in range(4):
50         score.append(model.fit(X,data.get_train_target(i)).score(x,data.get_test_target(i)))
51         print("\tTarget="+str(i),"Score:\t%f" % score[i])
52     print("Overall:",str(np.average(score)))
53
54 """ Rescales the features for MNB """
55 def rescale(X, x):
56     s = MinMaxScaler()
57     s.fit(X)
58     return s.transform(X), s.transform(x)
59
60 if __name__=="__main__":
61     data = dataset.Dataset(to_load=title+'.', chop=0.05)
62     # second_run(data)
63     build_models(data)

```