

Myers Briggs Personality Prediction using Machine Learning Techniques

Nishita Vaddem
Department of Comp Sc.
PES University
Bangalore, India

Pooja Agarwal
Department of Comp Sc.
PES University
Bangalore, India

ABSTRACT

In natural language processing and in the scientific realm of psychology, automatic personality analysis from social media is gaining growing interest. Currently, the Myers Briggs Type Indicator (MBTI) is deemed to be one of the most regularly used and reliable forms of personality recognition. The dataset used in this research is derived from Myers Briggs Forum on personalitycafe.com, a medium hitherto ignored for prediction of personality. This dataset is named as Myers-Briggs Type Indicators (MBTI) Personality Type and is available on Kaggle. The aim of this work is to predict the personality type of an individual linked to their posts and to explore the relevance of the test in the study and categorization of human behavior using Learning models.

Keywords

Myers-Briggs Type Indicators (MBTI)

1. INTRODUCTION

The notion of personality is considered an important yet imprecisely established construct in the world of psychology. Therefore, psychologists would benefit greatly from creating more specific, objective tests of current personality models. Most research on personality prediction has been extensively concentrated on the MBTI or Big Five personality models, which are the most regularly observed and commonly encountered personality recognition models in the world. The Big Five personality model can be described as a set of five dis-tinct dimensions, namely (1) extraversion, (2) agreeableness, (3) conscientiousness, (4) neuroticism and (5) openness [1]. MBTI is an introspective self-report test designed to show specific psychological patterns about how people view and make choices about the environment around them. This model recognizes 16 forms of personality, spanning four dimensions:

1. Introversion/Extraversion (how one gains energy),
2. Sensing/Intuition (how one takes in information),
3. Thinking/Feeling (how one builds decisions), and
4. Judging/Perceiving (how one presents herself or himself to the outside world) [2].

The Fig 1. shows the 16 specific personality types which is a combination of the groups which were spanned from the above four dimensions. For example, ISTJ is the personality type formed by combining Introversion, Sensing, Thinking, Judging. Fig 2. shows the 8 different types of personality utilized in the Myers–Briggs Type Indicator.



Fig 1. Types of Personality

Research indicates that the MBTI model has further applications, particularly in business and for self-exploration of personality types, despite disagreement over the reliability as well as the feasibility of these two models [3]. The goal of this work is to use machine learning to construct a classifier which will accept text as input (e.g. a social media post) and generate a prediction of the author's MBTI personality type as output.

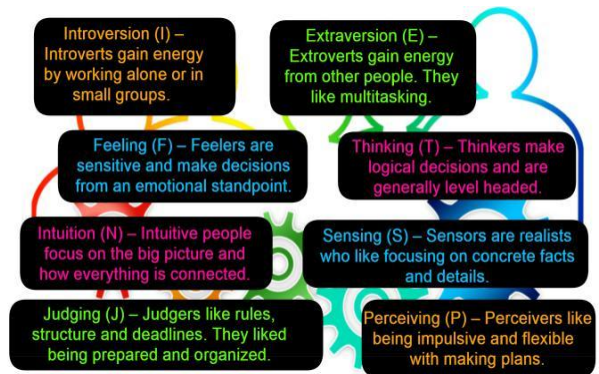


Fig 2. Four Dimensions of MBTI Personality

2. RELATED WORK

Research into prediction of personality styles from textual data is limited. However, huge strides have been taken through machine learning in this endeavor. Classic neural networks and machine learning techniques have been used for text classification, paraphrase detection, and predicting MBTI personality types. In a research by Komisin and Guinn [6], Naïve Bayes and Support Vector Machine (SVM) methods were made use of to predict the personality type of a single person depending on their choice of word.

Their database was constructed on the grounds of samples of in-class writing taken from 40 graduate students and their MBTI personality form. Later, they juxtaposed these two techniques, and established that on their limited dataset, the Naïve Bayes technique performs better because it treat them as independent from a "features" point of view, whereas SVM

looks at the interactions between them to some degree.

After two years, Wan et al. [5] made use of correlation analysis and principal component analysis to predict the Big Five personality type of individuals through their posts in a Chinese social network called inWeibo, and they were in a position to predict the personality type of those people successfully. They used two machine learning algorithms and the mean precision of prediction of five personalities was 0.707.

Li, Wan and Wang [7] made use of the gray prediction model, the multiple regression model and the multi-tasking model in order to predict the user's personality type using the Big Five model and their text samples. The gray prediction model is one of the most important predictive methods of the time series, which is used to solve problems of complexity with limited data. Multiple regression is an extension of simple linear regression. It is used when we want to estimate a variable's value based on two or more variables. In the multi-tasking model, multiple learning problems are solved, while the similarities are exploited across problems. The authors found that among these three models, the gray prediction model demonstrated the overall effect of the prediction between 0.8 and 0.9, the general accuracy of good prediction.

In a different research by Tandra et al. [9], the Big Five personality model and a deep learning architecture was utilized in order to predict a user's personality using the individual's information on their Facebook accounts. Their model achieved an average accuracy of 74.17%.

Additionally, in a separate study, Hernandez and Knight [8] utilized several kinds of recurrent neural networks (RNNs) namely simple RNN, gated recurrent unit (GRU), long short-term memory (LSTM) and Bidirectional LSTM in order to construct a classifier which is competent of predicting individual's MBTI personality trait depending on the text samples from their social media accounts. A simple RNN is a class of neural networks that allows previous outputs to be used as inputs while having hidden states. Gated recurrent unit (GRU) is an enhanced variant of the typical recurrent neural networks, using the concept of update gate and reset gate. These gates determine which information should be proceeded to the output. A special feature about this model is that it can be accustomed to retain the data from the distant past, without eliminating knowledge that is irrelevant to the prediction. Long short-term memory (LSTM) is an artificial recurrent neural network architecture which is made use of in the deep learning domain. They also compared the outputs and discovered that LSTM delivered the best results.

In another research by Gjurkovic´ and Snajder [10] utilized Logistic Regression, Multilayer Perceptron (MLP) and SVM to predict a person's MBTI personality type. This was performed using a dataset derived from Reddit. They found that MLP performed better with an overall accuracy of around 42%.

3. PROPOSED METHODOLOGY AND IMPLEMENTATION

3.1 Dataset

In this research, the Myers Briggs personality type dataset from Kaggle was utilized which is publicly available. The dataset includes 8675 observations where each finding shows an author's Myers Briggs personality type (a four-letter code) and raw text containing the author's last 50 posts where each post is separated by three pipe characters. Fig 3. shows two

columns labelled type and posts. The type column includes the 16 personality types, and the posts column includes raw text.

	type	posts
0	INFJ	'http://www.youtube.com/watch?v=qsXHcwe3krw ...
1	ENTP	'I'm finding the lack of me in these posts ver...
2	INTP	'Good one _____ https://www.youtube.com/wat...
3	INTJ	'Dear INTP, I enjoyed our conversation the o...
4	ENTJ	'You're fired. That's another silly misconce...

Fig 3. Sample dataset

3.2 Proportionality

Matplotlib, a Python plotting library and seaborn, a Python data visualization library have been used to preview data and value counts () method is used to get the number of occurrences of each personality kind. Fig 4. portrays the number of occurrences for each MBTI personality type. Since, INFP has more occurrences than ESTJ, this is an unbalanced dataset, as the classifier is biased to predict as INFP.

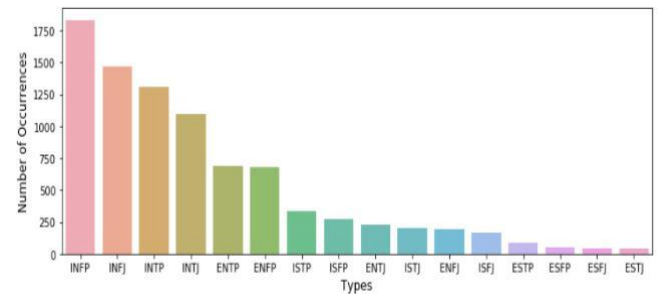


Fig 4. Distribution of dataset

3.3 Pre-processing

Since the data was collected from an online forum where the authors strictly convey their thoughts through text, some word removal was required. The most important explanation for this was the non-uniform distribution of MBTI types in the sample does not correspond with the true proportions of MBTI types. Finally, it has been concluded that this is due to the knowledge obtained from an Internet forum developed exclusively for debate about the type of personality and the MBTI characteristics were repeated in the posts too often. This might also impact the model's accuracy. As a result, all the characters were kept at lowercase. The following data were removed:

1. URLs,
2. Phrases that are not English letters (e.g., +, -, etc.),
3. Stop words (e.g., frequent words like "a", "an", "the", etc.) using the NLTK library,
4. MBTI profile strings (e.g., INFP, ESTJ, etc.) so that the model is not manipulated to identify MBTI mentions by name.

Ultimately, the text was lemmatized so as to achieve in making the dataset more consequential and pertinent using the NLTK library.

3.4 Model

The classification algorithms used in this step are: Logistic Regression (LR), Support Vector Machine (SVM) and XG-Boost. It belongs to the Supervised Learning category. The

classification task is divided into sixteen classes and further-
more into four binary classification tasks. This is because the
MBTI type is composed of four binary groups. Consequently,
four prominent binary classifiers have been trained, each
specializing in one of the personality aspects.

Using the model X and Y are obtained. X is the cleaned posts
after applying CountVectorizer and Term Frequency–Inverse
Document Frequency (TFIDF) Transformer. Y is the 4
columns 'I/E', 'N/S', 'F/T', 'J/P' in a binarized matrix. This X
and Y will be passed onto a Classification algorithm.
CountVectorizer essentially counts all the unique words in the
document and creates a matrix where occurrence of each word
is denoted by a 1, and non-occurrence of course, is denoted by
a 0.

```
Out[17]: [(0, 'ability'),
(1, 'able'),
(2, 'absolutely'),
(3, 'accept'),
(4, 'accurate'),
(5, 'across'),
(6, 'act'),
(7, 'action'),
(8, 'actual'),
(9, 'actually'),
(10, 'add'),
(11, 'admit'),
(12, 'advice'),
(13, 'afraid'),
(14, 'age'),
(15, 'ago'),
(16, 'agree'),
```

Fig 5. Sample of unique words

TFIDF Vectorizer converts the CountVectorizer output which
is a matrix of counts of each word in the document and
applies the Term-Frequency-Inverse Document Frequency
(TFIDF) formula to convert it into another matrix which will
be used as input by the machine learning model.

By examining the count vectorized matrix, 791 unique words
have been found. Fig 5. shows the list of unique words
obtained from the dataset. Using the proposed model, first the
prediction of individual personality types is done, for any
post, as Introversion (I) or Extroversion (E). Then, it will
attempt to predict Intuition (N) or Sensing (S) and so on.

MBTI type indicators were trained independently, and then
the data was separated into testing and training datasets
making use of the sklearn library's, train test split() feature.
Ultimately, 75% of the data was utilized for the training set
and 25% of the data was utilized for the test set. Model
according to the training data and the test data predictions was
designed.

4. RESULTS AND DISCUSSION

Though, XGBoost comparatively gave a lower accuracy as
opposed to the other two models (LR and SVM) but it could
figure out the features (words) in the document having highest
importance, using the plot importance feature. In this work,
four XGBoost models are applied to the dataset, assigning
max num features=10 to each model and then top 10
important features are obtained. These four XGBoost models
are a collection of 40 features. Finally the model could obtain
the most important feature(word) from each of those 10
important features(words) (i.e., one per each XGBoost model)
for I/E, N/S, F/T, J/P namely ne, si, feeling and ni
respectively. This is due to lemmatization not being accurate
enough. Fig 6. shows feature importance of words. X-axis
represents the top 10 features (i.e., words) in a model, whereas
the Y-axis represents the number of occurrences of each

feature.

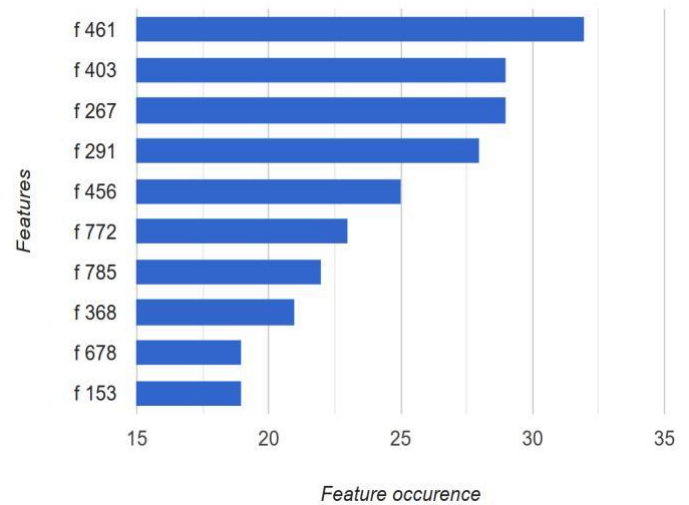


Fig 6. Feature Importance

The most important words identified by the model are non-
words, so the phrases: ne, si and ni are removed from the
training dataset. We attempt multi-class classification only to
get an accuracy around 31.17%. While the multi-class
precision provided by the XGBoost model was low, we can
apply binary classifications individually for each of the 4
types to generate higher accuracy. The accuracy for the binary
classification tasks after configuration is mentioned in the
following tables 1, 2 and 3.

Table 1: Using LR

MBTI Personality Type	Accuracy
Introversion (I)/ Extroversion (E)	79.07%
Intuition (N)/ Sensing (S)	86.68%
Feeling (F)/ Thinking (T)	77.87%
Judging (J)/ Perceiving (P)	67.17%

Table 2: Using SVM

MBTI Personality Type	Accuracy
Introversion (I)/ Extroversion (E)	76.39%
Intuition (N)/ Sensing (S)	86.58%
Feeling (F)/ Thinking (T)	76.76%
Judging (J)/ Perceiving (P)	68.42%

Table 3: Using XGBoost

MBTI Personality Type	Accuracy
Introversion (I)/ Extroversion (E)	76.21%
Intuition (N)/ Sensing (S)	85.62%
Feeling (F)/ Thinking (T)	75.01%
Judging (J)/ Perceiving (P)	63.76%

5. CONCLUSION

Various Python libraries like Pandas, Numpy, NLTK,
Seaborn, Matplot Lib and Sklearn etc were used in the process
of building a machine learning classifier to automate the
process of predicting MBTI personality types. In this work,
three algorithms namely, Logistic Regression, SVM and
XGBoost were used for classification. Although, the multi-

class accuracy was around 31.17% produced by the XGBoost model, we can apply binary classifications for each of the 4 types individually to generate higher accuracy.

6. REFERENCES

- [1] Soto, C.J. Big Five personality traits. In *The SAGE Encyclopedia of Lifespan Human Development*; Borstein, M.H., Arterberry, M.E., Fingerman, K.L., Lansford, J.E., Eds.; SAGE Publications: Thousand Oaks, CA, USA, 2018; pp. 240–241.
- [2] Myers, I.B.; McCaulley, M. *Manual: A Guide to the Development and Use of the Myers-Briggs Type Indicator*, 15th ed.; Consulting Psychologists Press: Santa Clara, CA, USA, 1989.
- [3] John, E.; Barbuto, J.R. A critique of the Myers-Briggs Type indicator and its operationalization of Carl Jung's Psychological types. *Psychol. Rep.* 1997, 80, 611–625.
- [4] Tieger, P.D.; Barron-Tieger, B. *Do What You Are: Discover the Perfect Career for You through the Secrets of Personality Type*, 4th ed.; Sphere: London, UK, 2007.
- [5] Wan, D.; Zhang, C.; Wu, M.; An, Z. Personality prediction based on all characters of user social media information. In *Proceedings of the Chinese National Conference on Social Media Processing*, Beijing, China, 1–2 November 2014; pp. 220–230.
- [6] Komisin, M.; Guinn, C. Identifying personality types using document classification methods. In *Proceedings of the 25th International Florida Artificial Intelligence Research Society Conference*, Marco Island, FL, USA, 23–25 May 2012; pp. 232–237.
- [7] Li, C.; Wan, J.; Wang, B. Personality Prediction of Social Network Users. In *Proceedings of the 16th International Symposium on Distributed Computing and Applications to Business, Engineering and Science*, Anyang, China, 13–16 October 2017.
- [8] Hernandez, R.; Knight, I.S. Predicting Myers-Briggs Type Indicator with text classification. In *Proceedings of the 31st Conference on Neural Information Processing Systems*, Long Beach, CA, USA, 4–9 December 2017.
- [9] Tandra, T.; Suhartono, D.; Wongso, R.; Prasetyo, Y. Personality prediction system from Facebook users. In *Proceedings of the 2nd International Conference on Computer Science and Computational Intelligence*, Bali, Indonesia, 13–14 October 2017.
- [10] Gjurkovic, Matej & Snajder, Jan. (2018). *Reddit: A Gold Mine for Personality Prediction*. 87-97. 10.18653/v1/W18-1112.