

Capstone project 2 Milestone Report- Personality Prediction MBIT

Problem statement

Have you ever been curious about the personality type of yourself and people around you? Is there any simple test to help us find the answer? The answer is yes. You can figure out your personality type by simply answering these four questions:

- Are you outwardly or inwardly focused ?
- How do you prefer to take in information?
- How do you prefer to make decisions?
- How do you prefer to live your outer life?

For each question, there are two options to choose from. You choose the side that seems most natural for you. The options are: Introversion (**I**) or Extraversion (**E**), Sensing (**S**) or Intuition (**N**), Thinking (**T**) or Feeling (**F**) and Judging (**J**) or Perception (**P**) respectively. The introspective self-report questionnaire we introduced above is known as Myers-Briggs -Type- Indicator (MBTI). It is a personality type system that divides everyone into 16 distinct personalities based on their answers to the previous questions. For example, someone who chooses extraversion, sensing, thinking and judging would be labeled an ESTP in the MBTI system.

This project focuses on using machine learning algorithms in supervised machine learning to predict the personality of a person from the type of posts they put on social media.

Data Source

The data was collected through the PersonalityCare forum and is available on Kaggle [Personality Prediction Dataset](#). This dataset consists of over 8675 rows representing 8675 different people and 2 columns representing a person's MBTI personality type and the things they have posted.

Exploratory Data Analysis

We first load the data to a DataFrame using `pd.read_csv` and look at the structure of the data set.

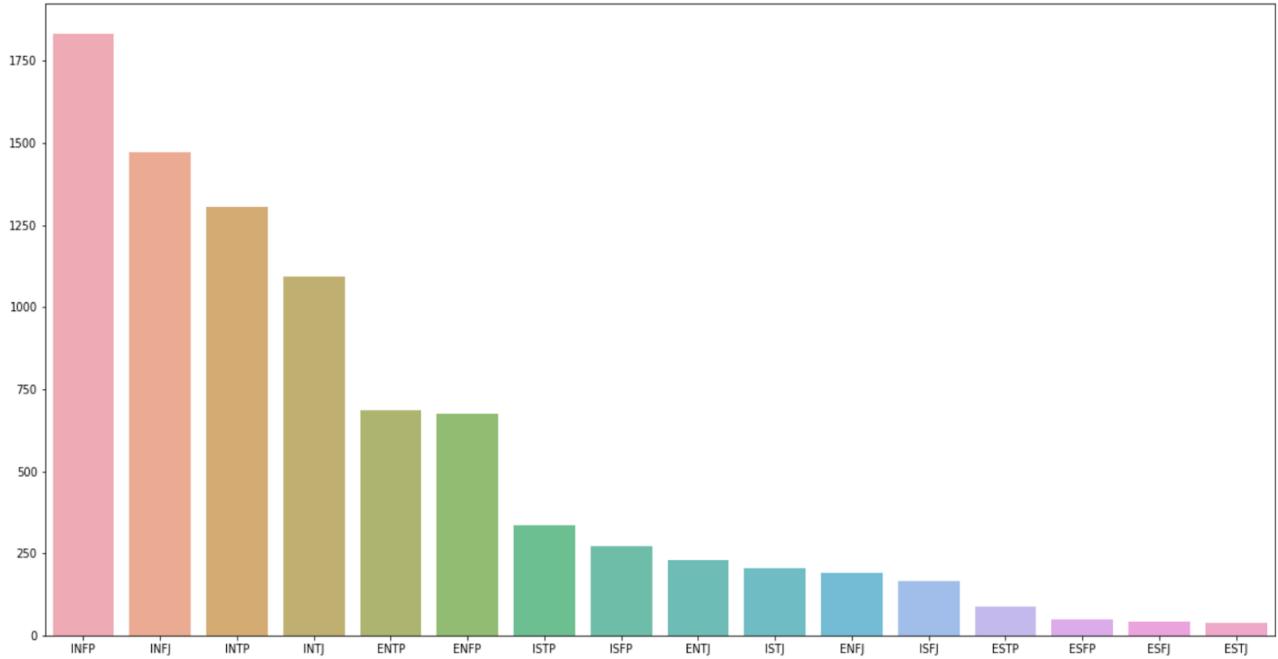
| | type | posts |
|---|------|---|
| 0 | INFJ | 'http://www.youtube.com/watch?v=qsXHcwe3krw ... |
| 1 | ENTP | 'I'm finding the lack of me in these posts ver... |
| 2 | INTP | 'Good one ____ https://www.youtube.com/wat... |
| 3 | INTJ | 'Dear INTP, I enjoyed our conversation the o... |
| 4 | ENTJ |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

The data has two columns: **type** and **posts** and no missing data. The two columns represent the personality type of 8675 different people and their corresponding posts on social media. The personality type is expressed by 4 letters representing the answers for 4 questions above. The first letter can either be **I**(Introversion) or **E**(Extraversion). The second one can either be **S**(sensing) or **N**(Intuition). The third one can either be **T**(Thinking) or **F**(Feeling). The last letter can either be **J**(Judging) or **P**(Perception). The first person in our table has personality type **INFJ** implying Introversion-Intuition-Feeling-Judging.

We are going to analyze the data set using common sense. It is natural to think that the length of a post, number of question marks, exclamation, ellipsis or images in a post somehow reveal the personality of the owner. We are going to find these numbers and add columns **post_length**, **question_per_post**, **excl_per_post**, **ellipsis_per_post**, **img_per_post** to our data df.

| | type | posts | post_length | question_per_post | img_per_post | excl_per_post | ellipsis_per_post |
|---|------|---|-------------|-------------------|--------------|---------------|-------------------|
| 0 | INFJ | 'http://www.youtube.com/watch?v=qsXHcwe3krw ... | 556 | 18 | 6 | 3 | 15 |
| 1 | ENTP | 'I'm finding the lack of me in these posts ver... | 1170 | 5 | 1 | 0 | 19 |
| 2 | INTP | 'Good one ____ https://www.youtube.com/wat... | 836 | 12 | 0 | 4 | 13 |
| 3 | INTJ | 'Dear INTP, I enjoyed our conversation the o... | 1064 | 11 | 0 | 3 | 26 |
| 4 | ENTJ | 'You're fired. That's another silly misconce... | 967 | 10 | 2 | 1 | 21 |

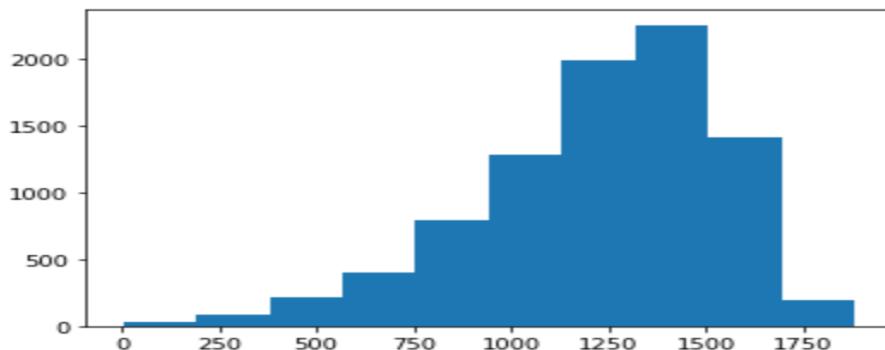
There are 16 personality types. How are they distributed in our data? We first want to explore the distribution of different personality types.



It is obvious that the distribution is unbalanced between different classes. The group **INFP(Introversion-Intuition-Feeling-Perception)** has the highest frequency while the group **ESTJ (Extraversion-Sensing-Thinking-Judging)** has the lowest frequency.

The largest group is **INFP (Introversion-Intuition-Feeling-Perception)** with 1832 posts, while the smallest group is **ESTJ (Extraversion-Sensing-Thinking-Judging)** with only 39 posts. We observe that the majority of posts (5697) in our data belong to **IN** (Introversion-Intuition) people while the **ES(Extraversion_Sensing)** people contribute only 218 posts. The difference in frequencies is really significant. The group with highest frequency is almost triple the group with 5th highest frequency, and 46 times more than the group with the lowest frequency.

We also want to look at the distribution of the length of posts.



The majority of posts in our data are more than 500 words. The models based on this data will work best for posts or texts with more than 500 words.

We now look at the relations between the type of personality and the length of post, as well as the number of question marks, exclamations, ellipsis and images per post.

| post_length | question_per_post | excl_per_post |
|-------------------------|--------------------------|-----------------------|
| type | type | type |
| ESFJ 1290.476190 | ESTP 12.292135 | ENFP 16.939259 |
| ENFJ 1286.584211 | ENTJ 12.000000 | ENFJ 13.800000 |
| INFJ 1278.431973 | ESFP 11.854167 | ESFJ 11.523810 |
| ENFP 1260.770370 | ENFP 11.362963 | ESFP 10.958333 |
| INFP 1244.552948 | ISTP 11.080119 | ISFJ 10.475904 |
| ISFJ 1241.295181 | INTP 11.078988 | ISFP 10.398524 |
| ESTJ 1229.538462 | ENTP 11.048175 | INFP 9.322052 |
| ENTJ 1218.086580 | ISFP 10.819188 | INFJ 9.077551 |
| ISTJ 1213.224390 | ISTJ 10.746341 | ENTJ 8.298701 |
| ENTP 1205.995620 | INTJ 10.742438 | ESTJ 8.205128 |
| INTP 1197.763037 | ENFJ 10.652632 | ESTP 8.000000 |
| INTJ 1194.577452 | ESTJ 10.487179 | ISTJ 7.980488 |
| ISTP 1165.563798 | INFJ 10.448299 | ENTP 7.129927 |
| ESTP 1162.595506 | INFP 10.126638 | INTJ 5.262145 |
| ISFP 1136.383764 | ISFJ 10.048193 | INTP 5.004601 |
| ESFP 1022.125000 | ESFJ 8.809524 | ISTP 4.629080 |

The **ESFJ** (Extraversion-Sensing-Feeling-Judging) tends to have the longest post(average 1290 words per post) while the **ESFP** (Extraversion-Sensing-Feeling-Perception) tends to have the shortest post (average 1022 words per post).

People from the **ESTP** (Extraversion-Sensing-Thinking-Perception) tend to use more question marks in their post (average 12.3 question marks per post) while people in the **ESFJ** (Extraversion-Sensing-Feeling-Judging) tend to use less question marks in their posts (average 8.8 question marks per post). The difference in question marks between groups is not significant compared to the length of these posts.

The group **ENFP** (Extraversion-Intuition-Feeling-Perception) use the most exclamation (Average 16.9 exclamation per post) while the group **ISTP** (Introversion-Sensing-Thinking-Perception) use the least exclamation (average 4.6 exclamation per post).

Data Preprocessing

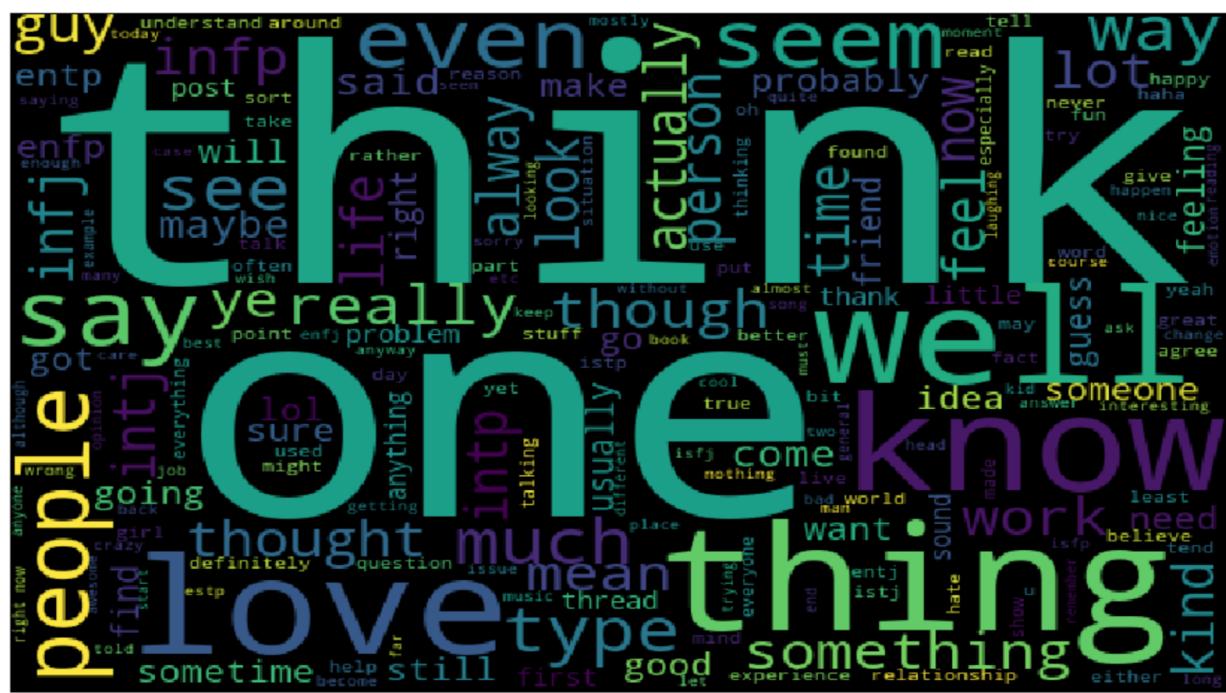
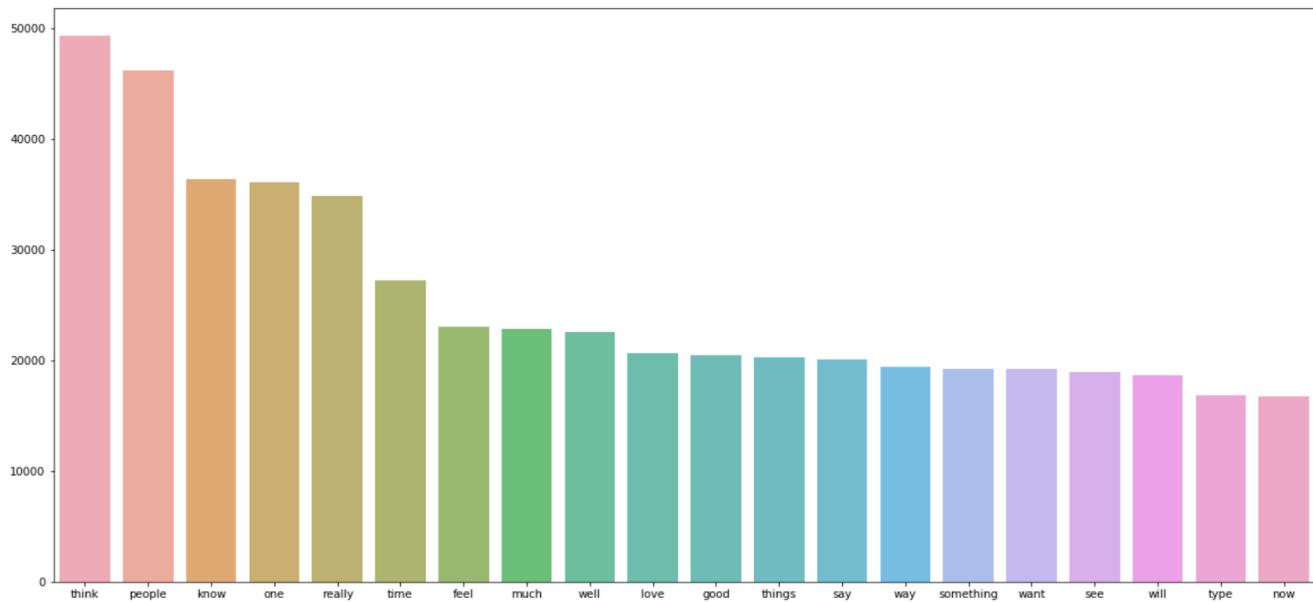
We are going to create 3 functions which play an important role in preprocessing text. They are : **clean_text**, **combine_text** and **text_preprocessing**.

- **clean_text** function takes the argument text, the make the text lowercase, remove text in square brackets, remove links, remove punctuation and remove words containing numbers.
- **text_preprocessing** function takes text as an argument. It first applies the function **clean_text** on the argument, tokenizes the result, removes all stop words before combining all tokens.
- **combine_text** function combine single words into text.

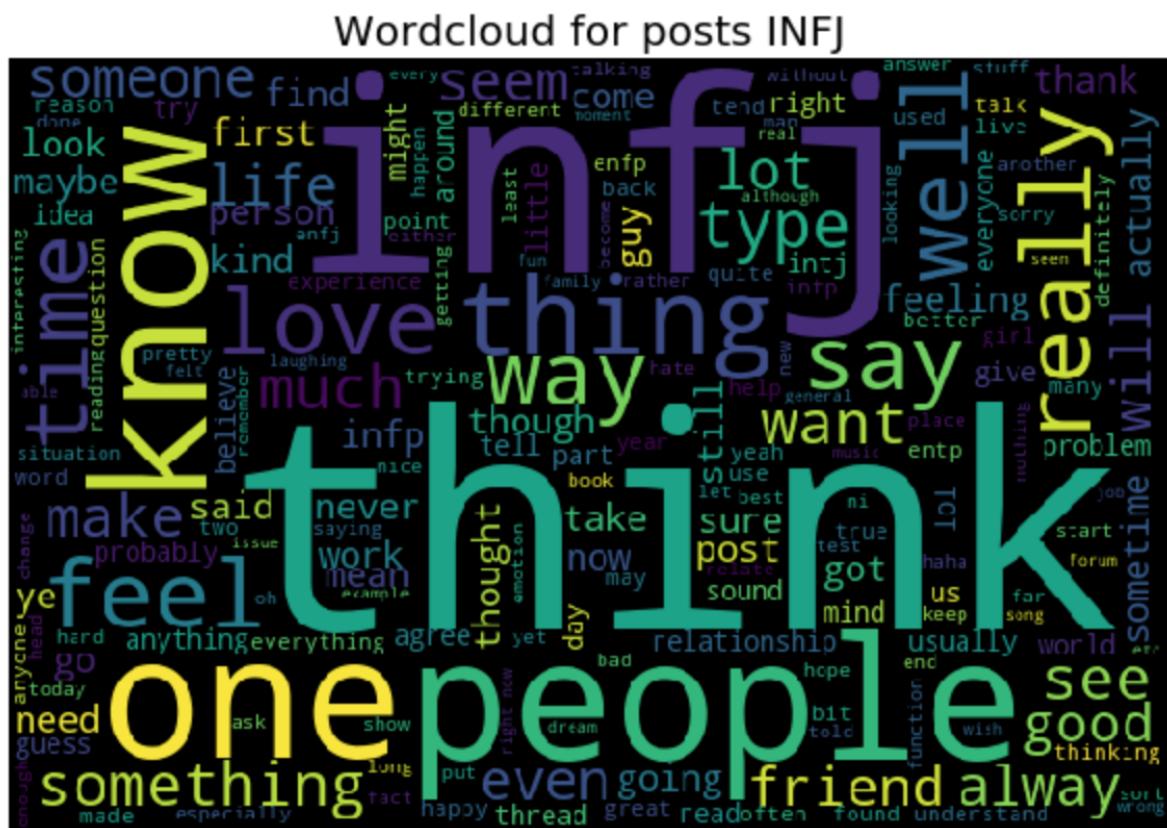
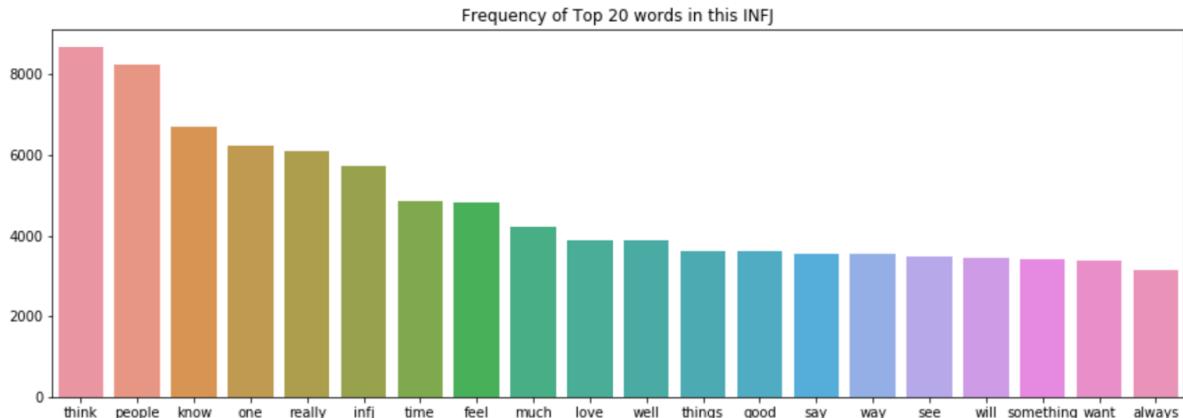
The **clean_posts** column is added which contains the corresponding posts which have been preprocessed.

Most Used Words and Word Cloud.

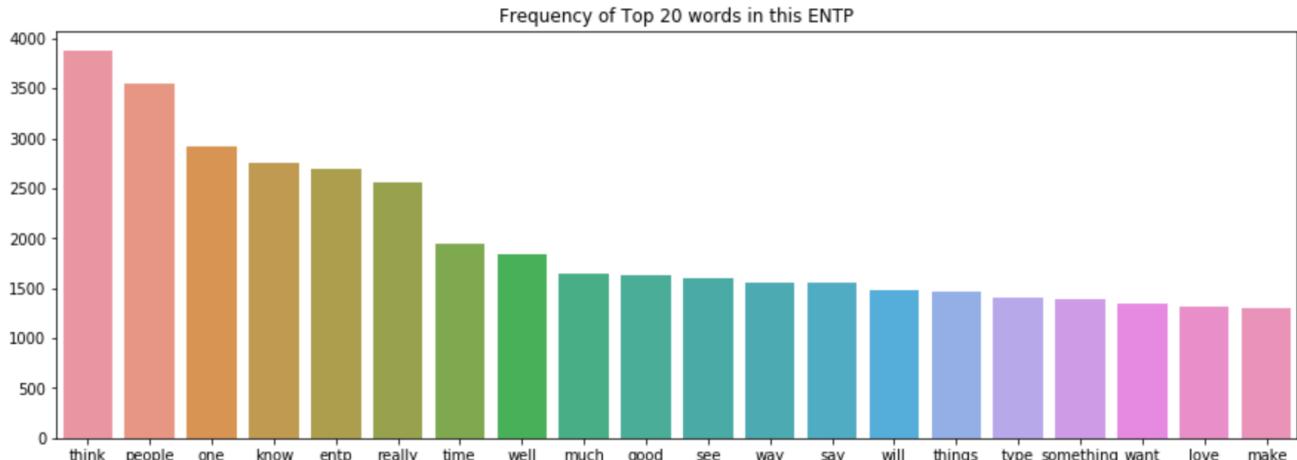
- Most Used Words and Word Cloud in all posts.



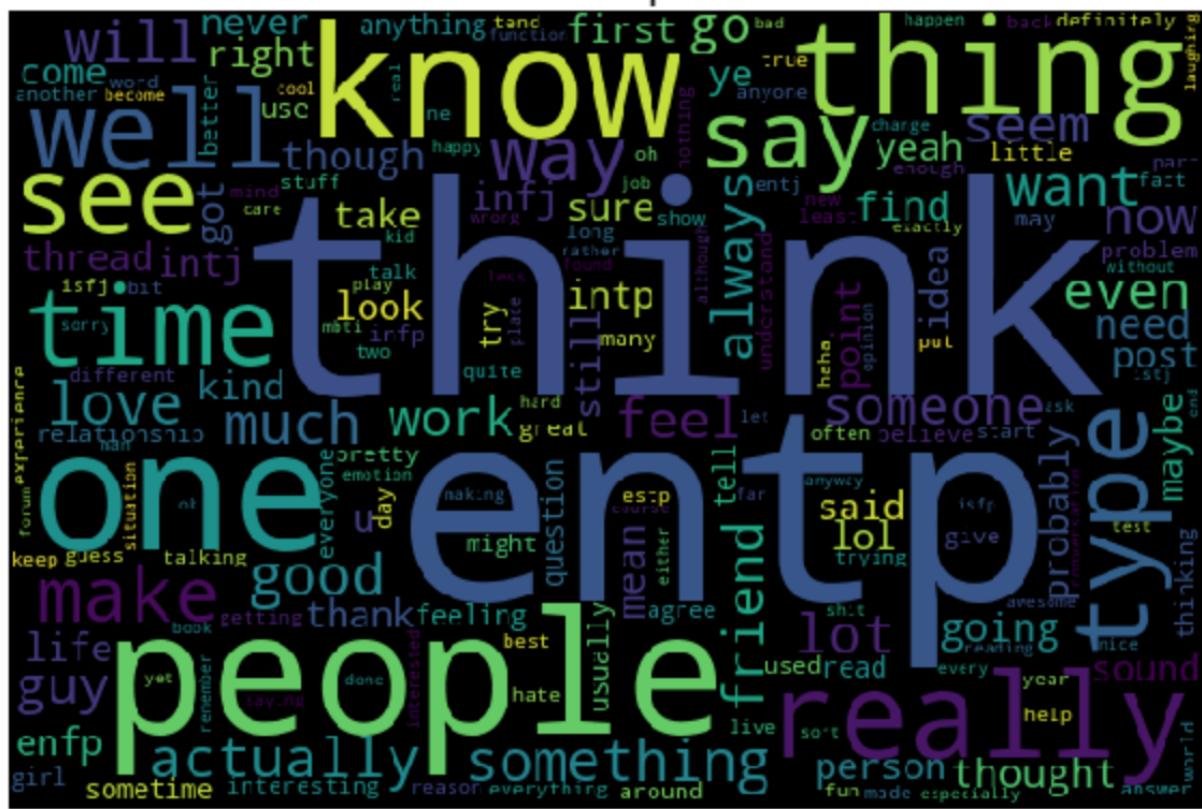
- Most Used Words and Word Cloud in each personality group
- INFJ



- ENTP

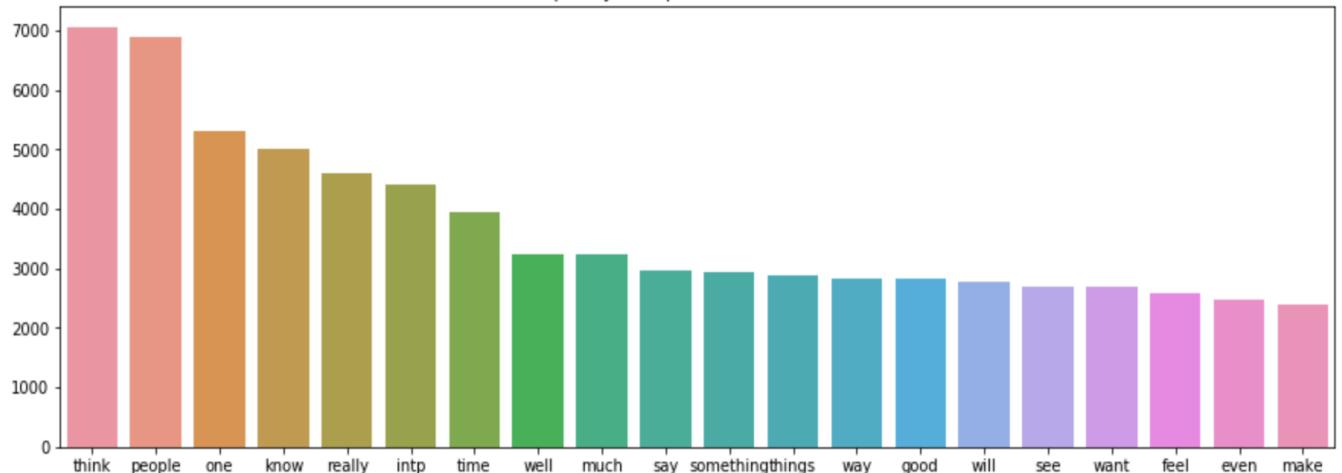


Wordcloud for posts ENTP



- INTP

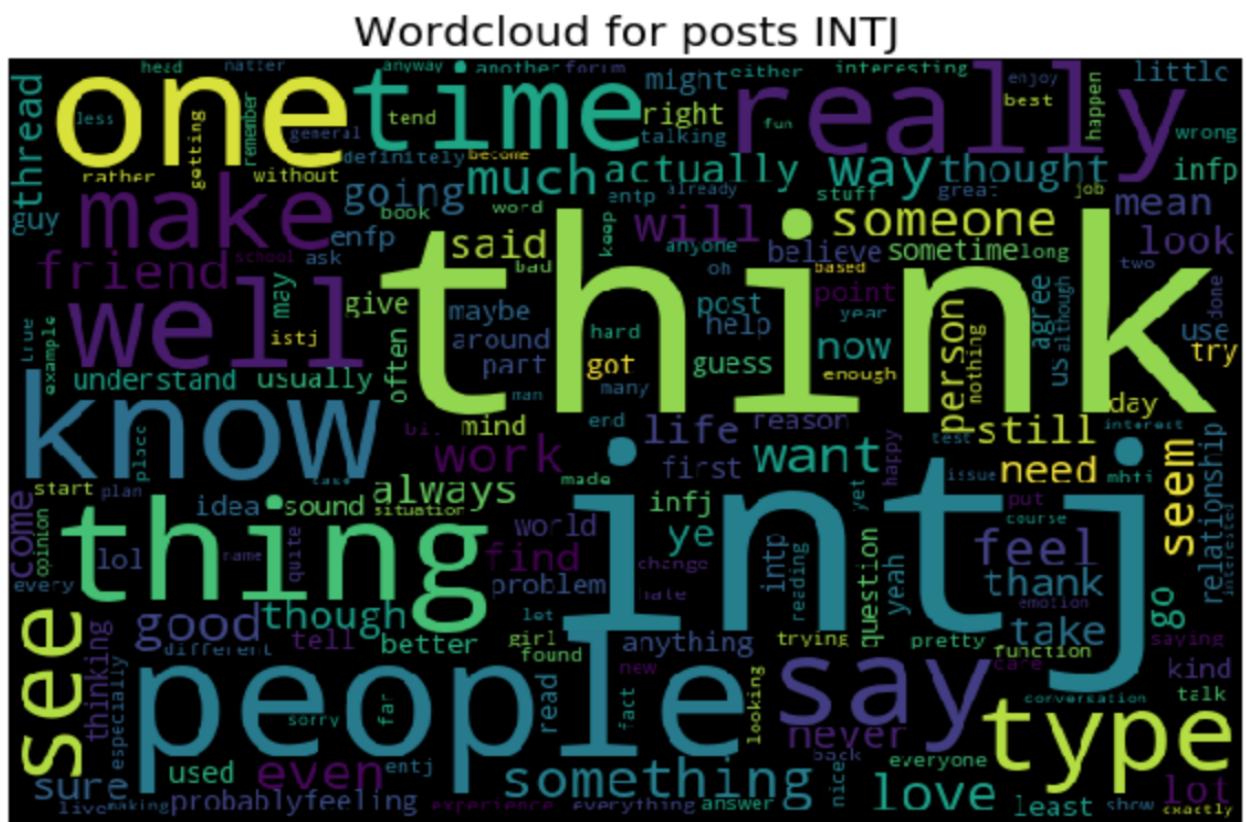
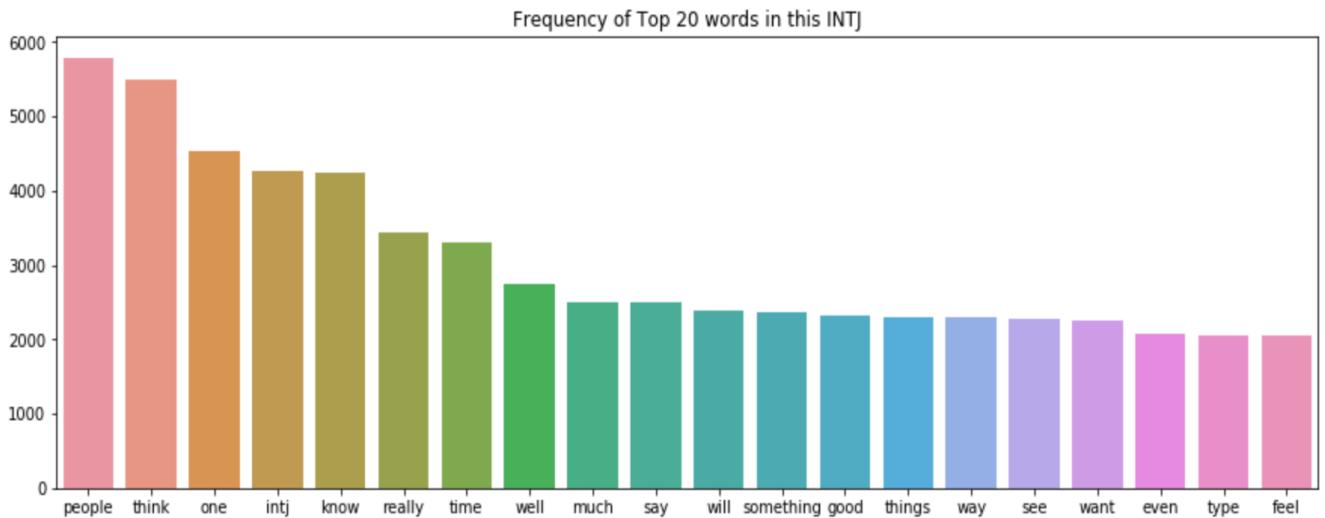
Frequency of Top 20 words in this INTP



Wordcloud for posts INTP

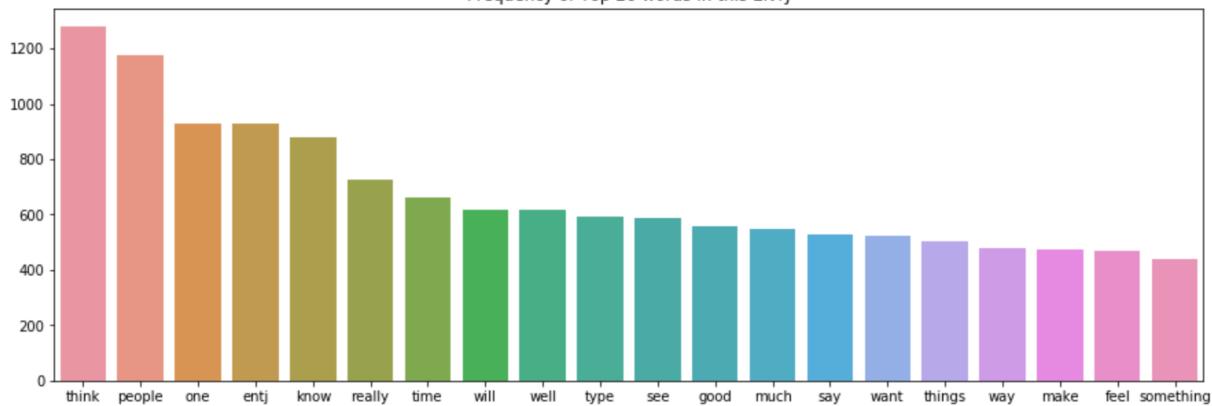


- INTJ



- INTJ

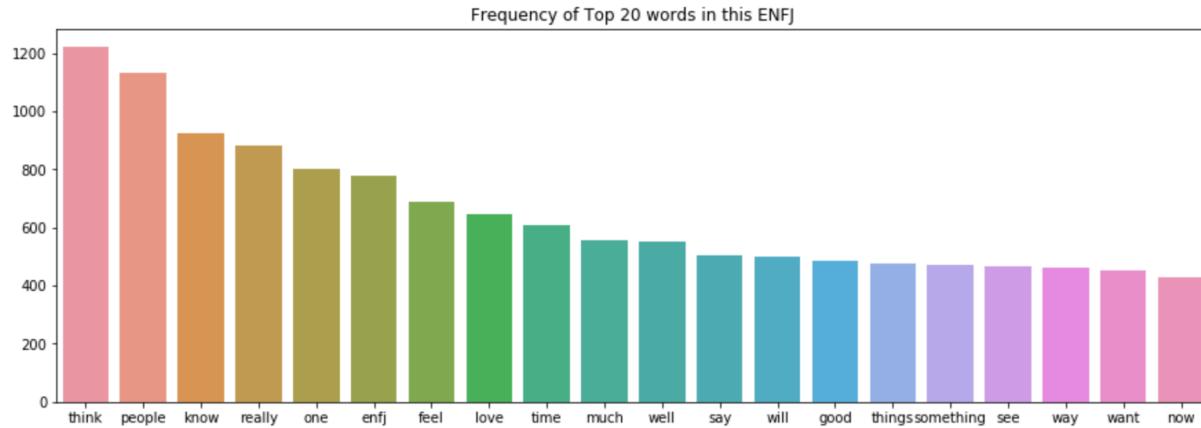
Frequency of Top 20 words in this ENTJ



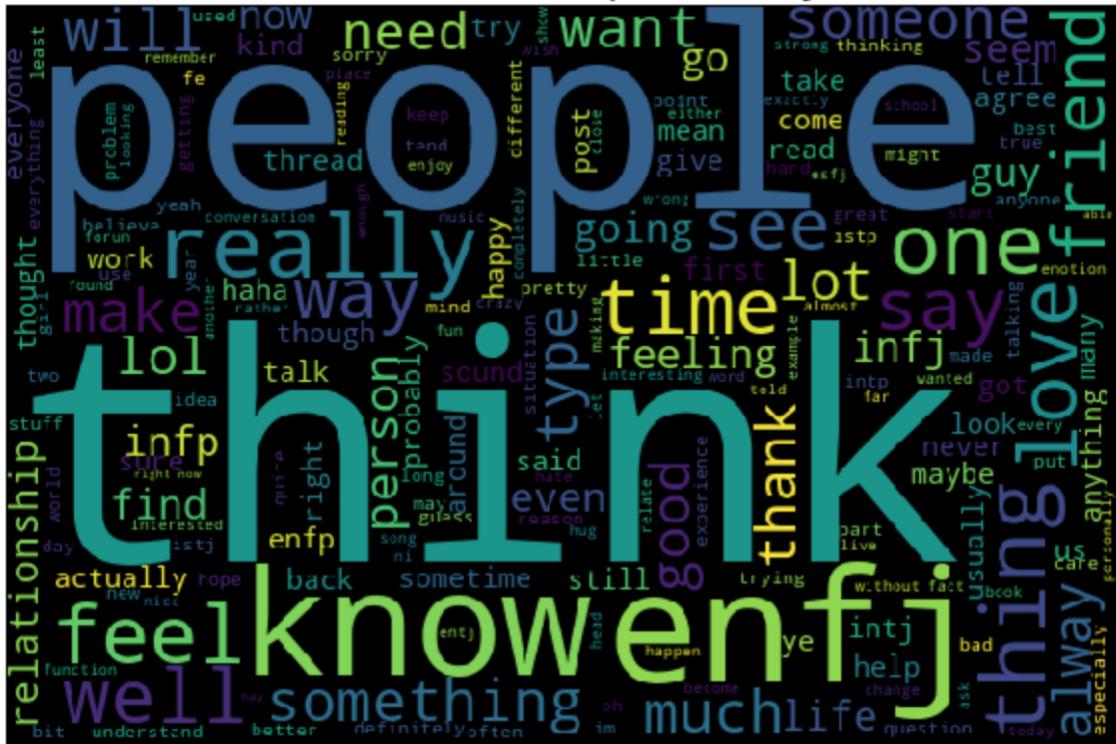
Wordcloud for posts ENTJ



● ENFJ

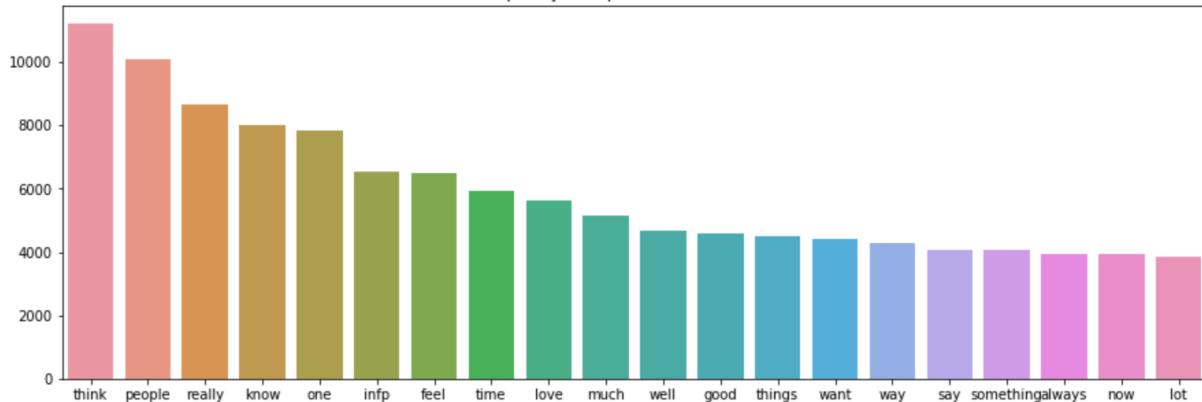


Wordcloud for posts ENFJ

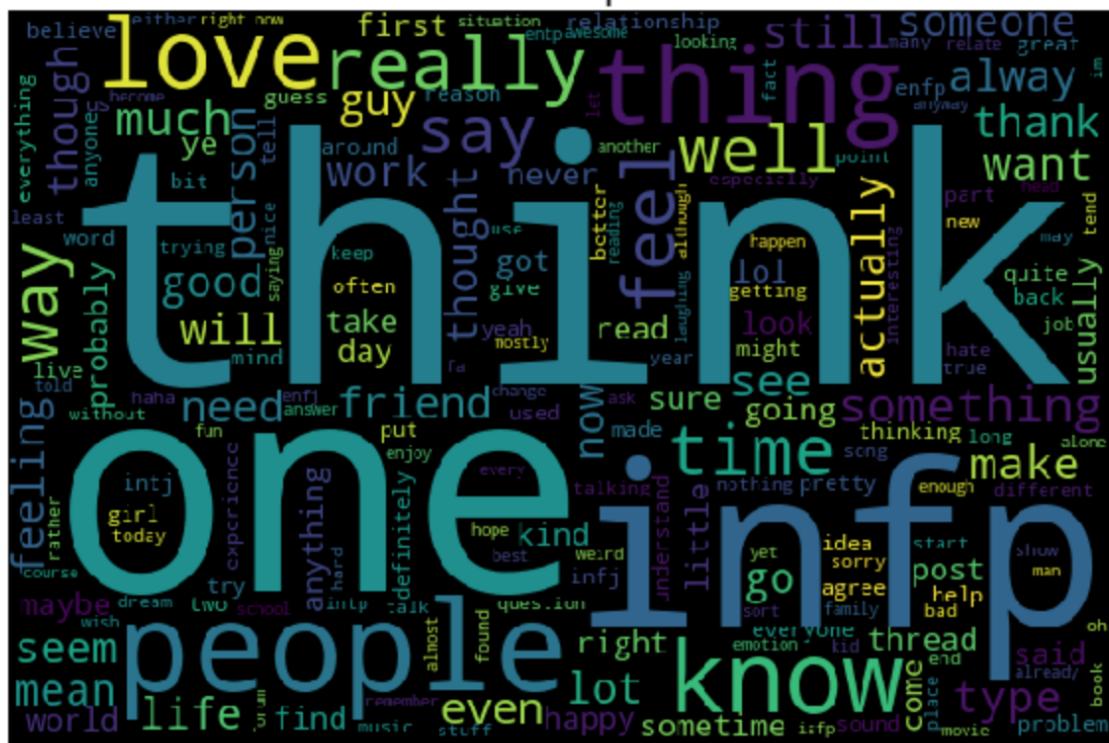


- INFP

Frequency of Top 20 words in this INFP

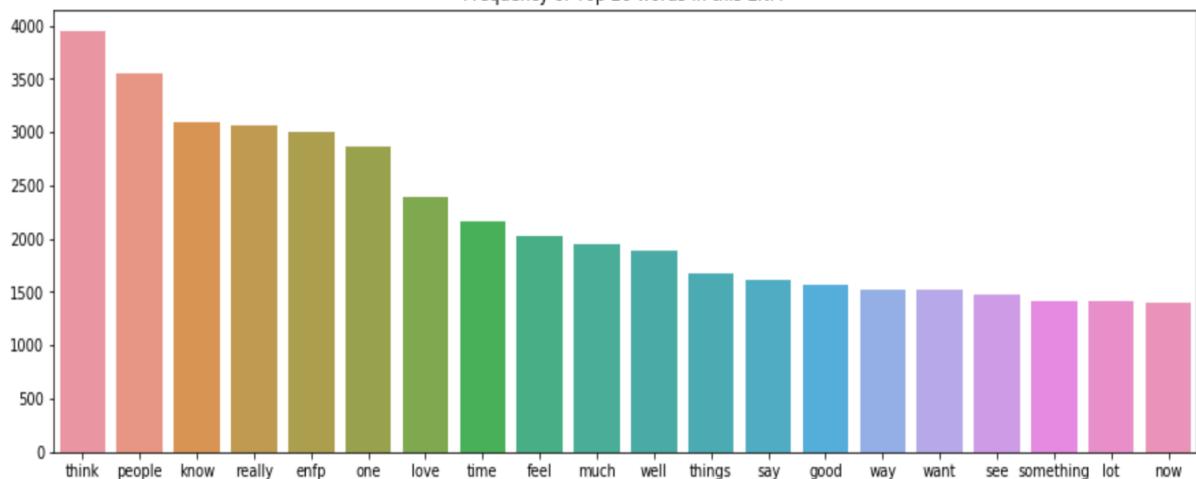


Wordcloud for posts INFP



- ENFP

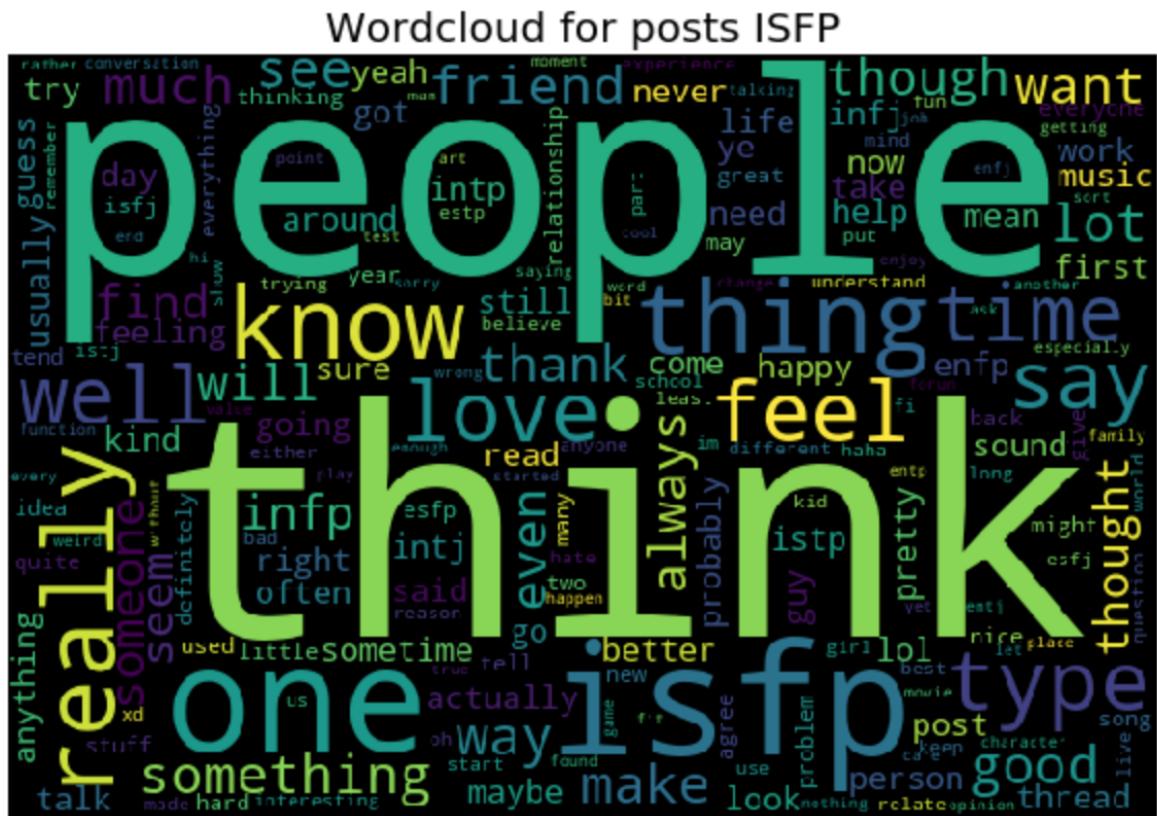
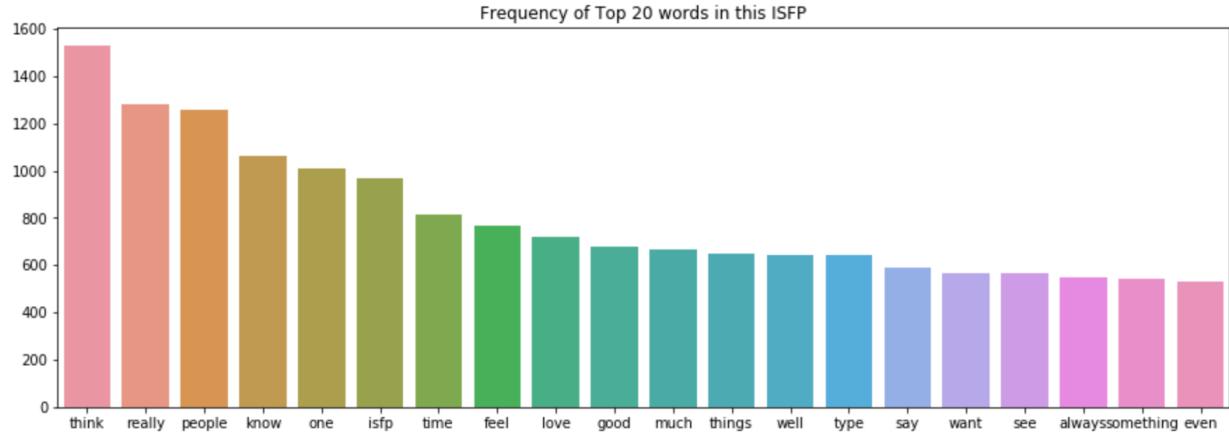
Frequency of Top 20 words in this ENFP



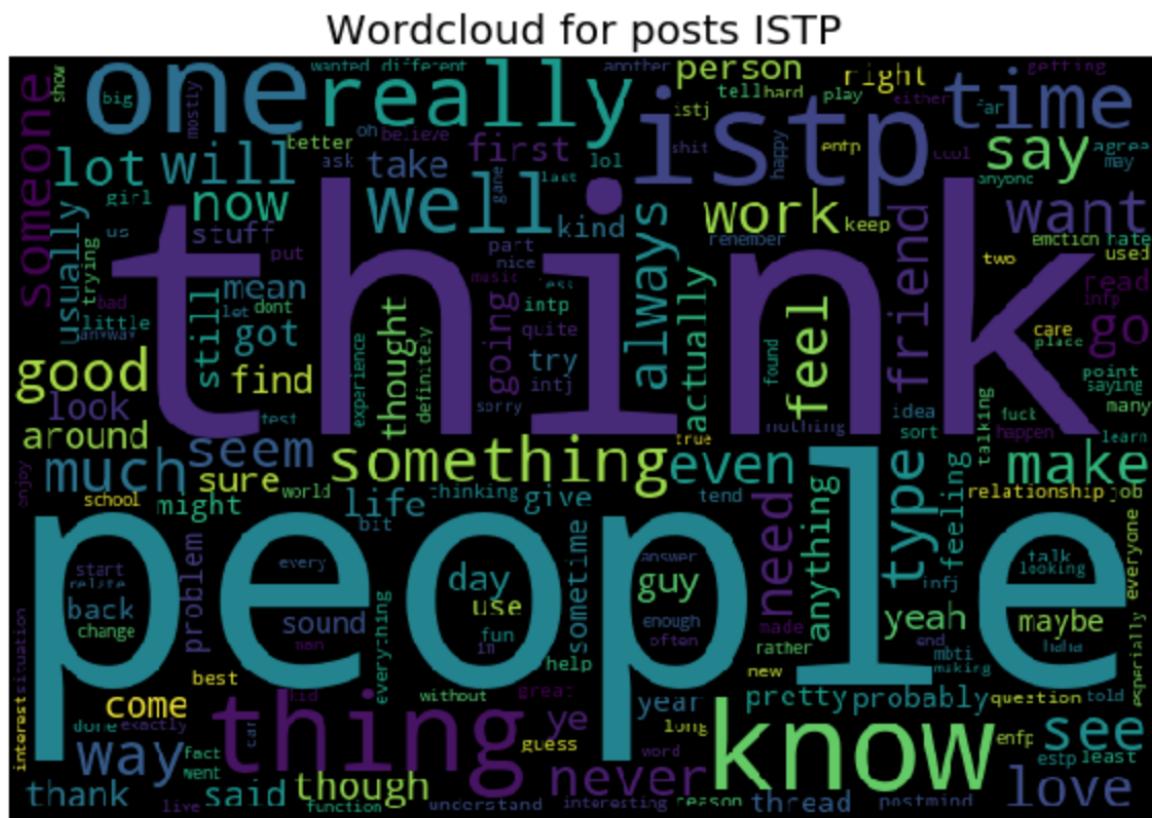
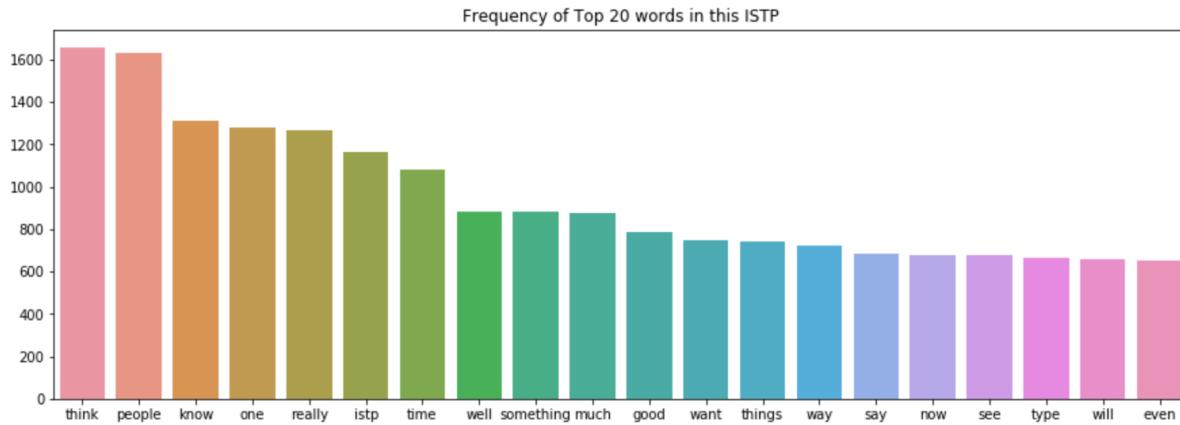
Wordcloud for posts ENFP



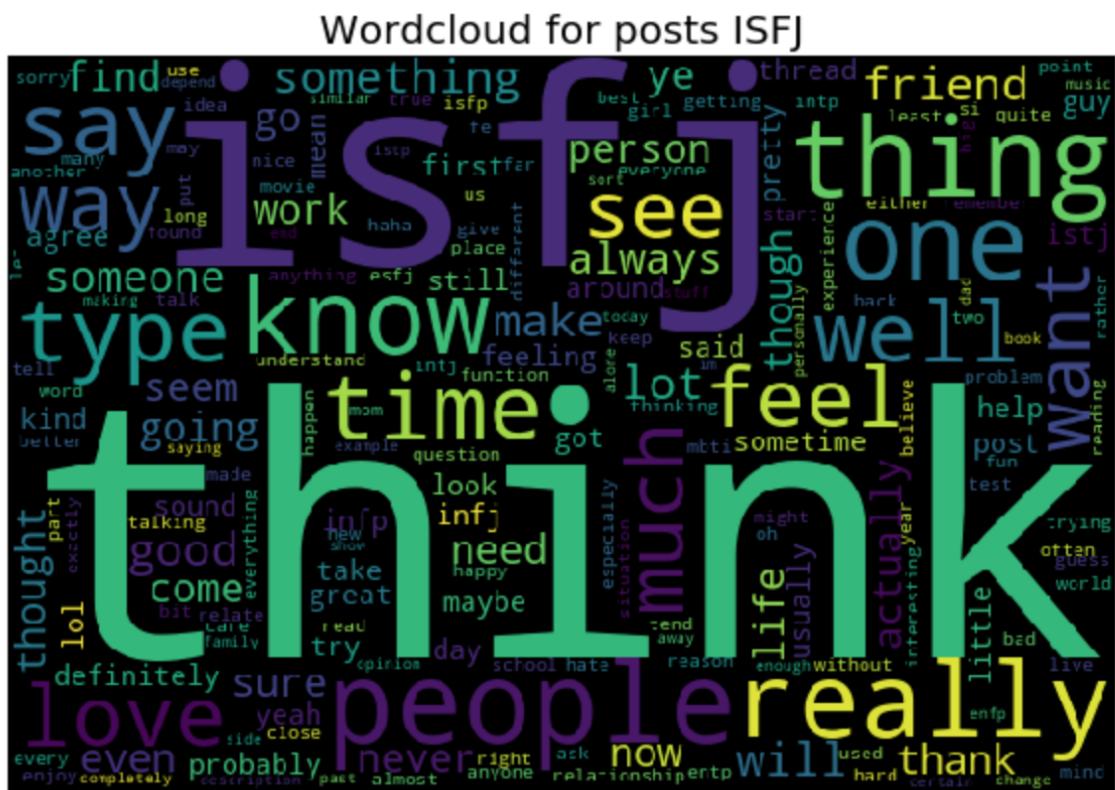
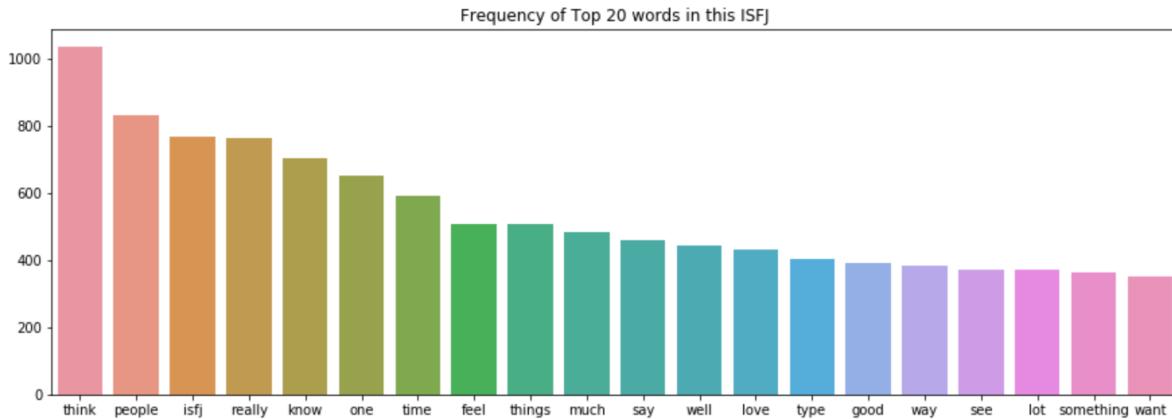
- ISFP



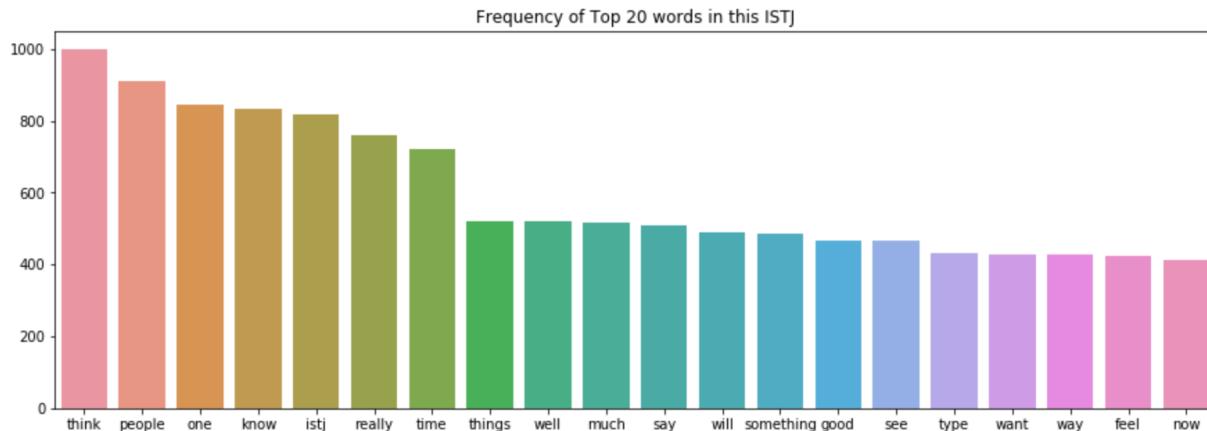
- ISTP



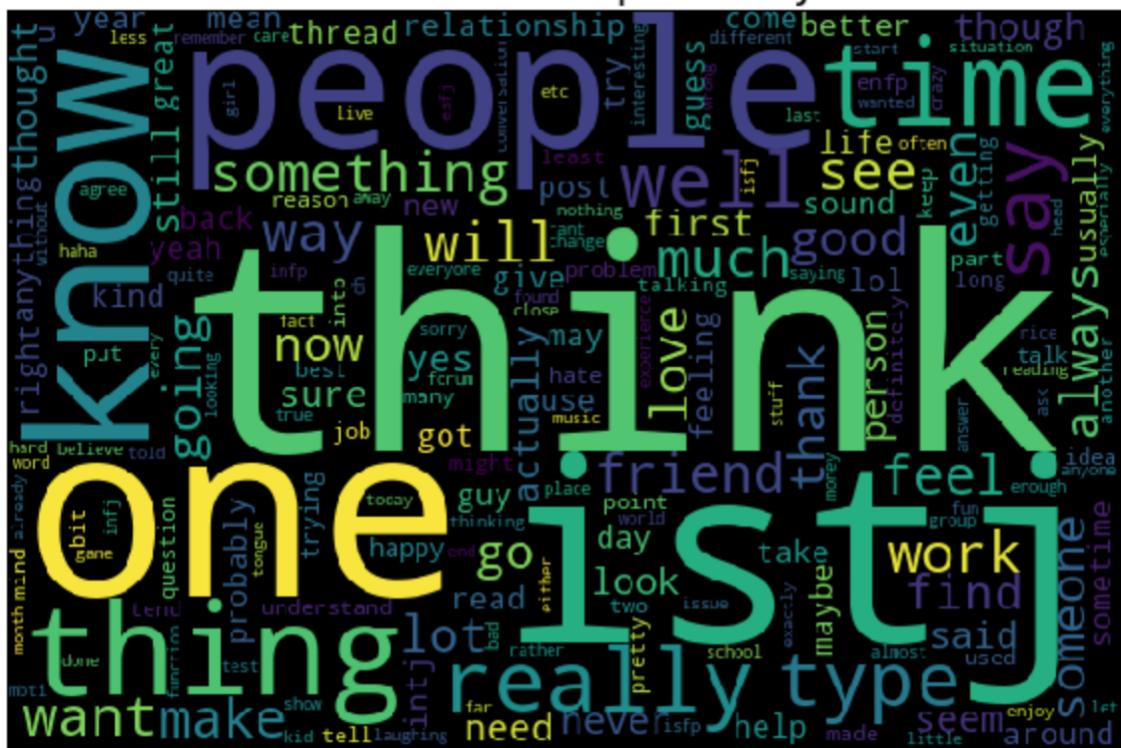
- ISFJ



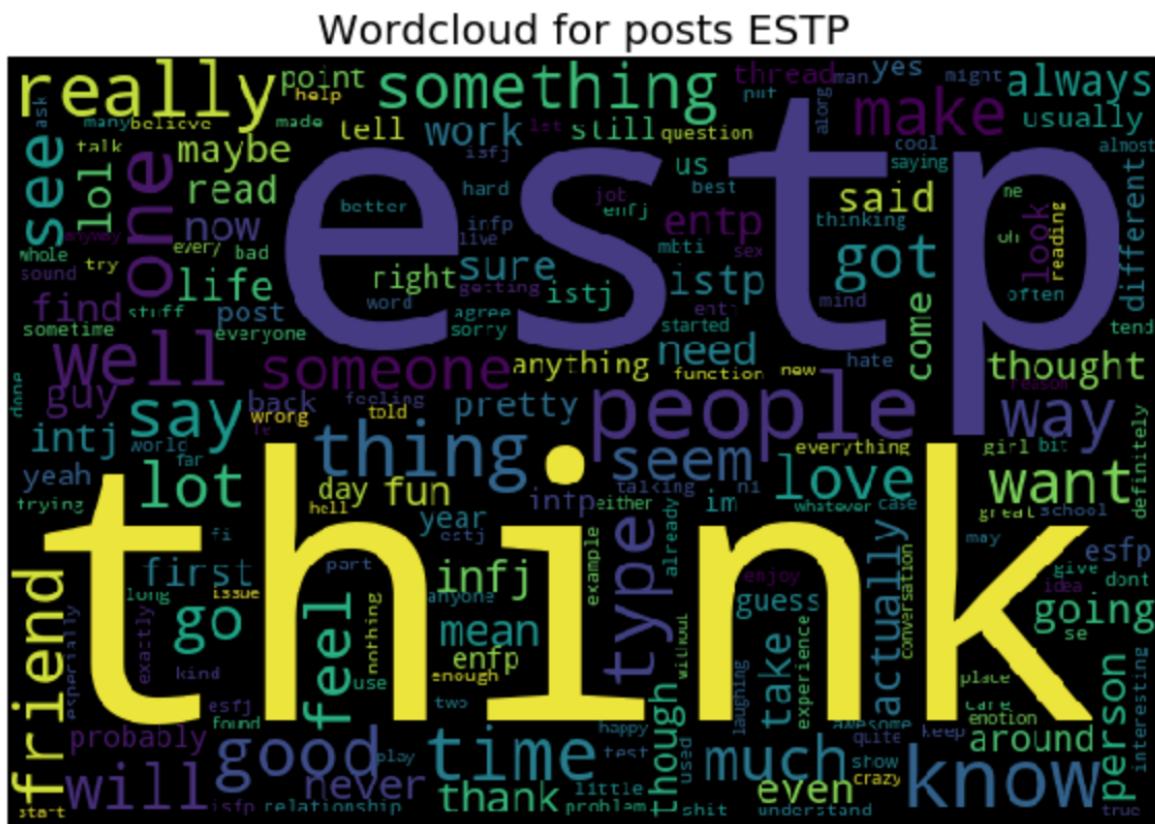
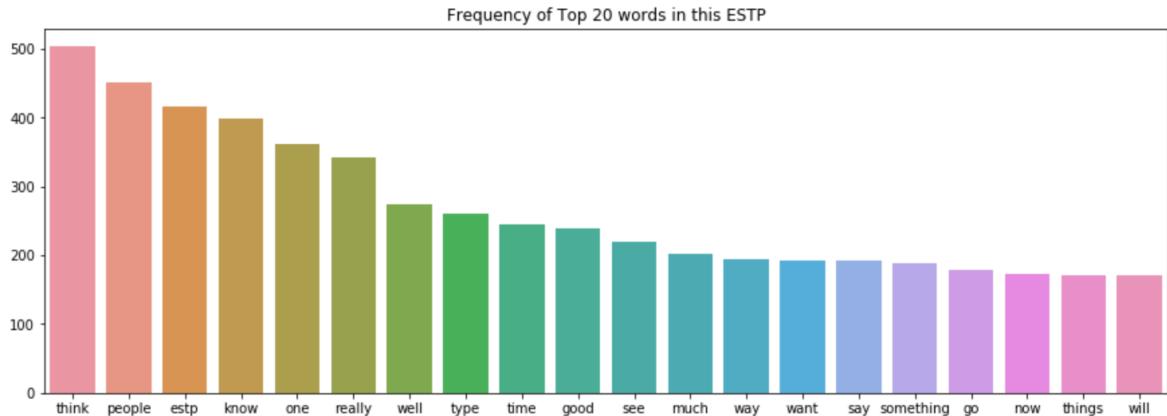
- ISTJ



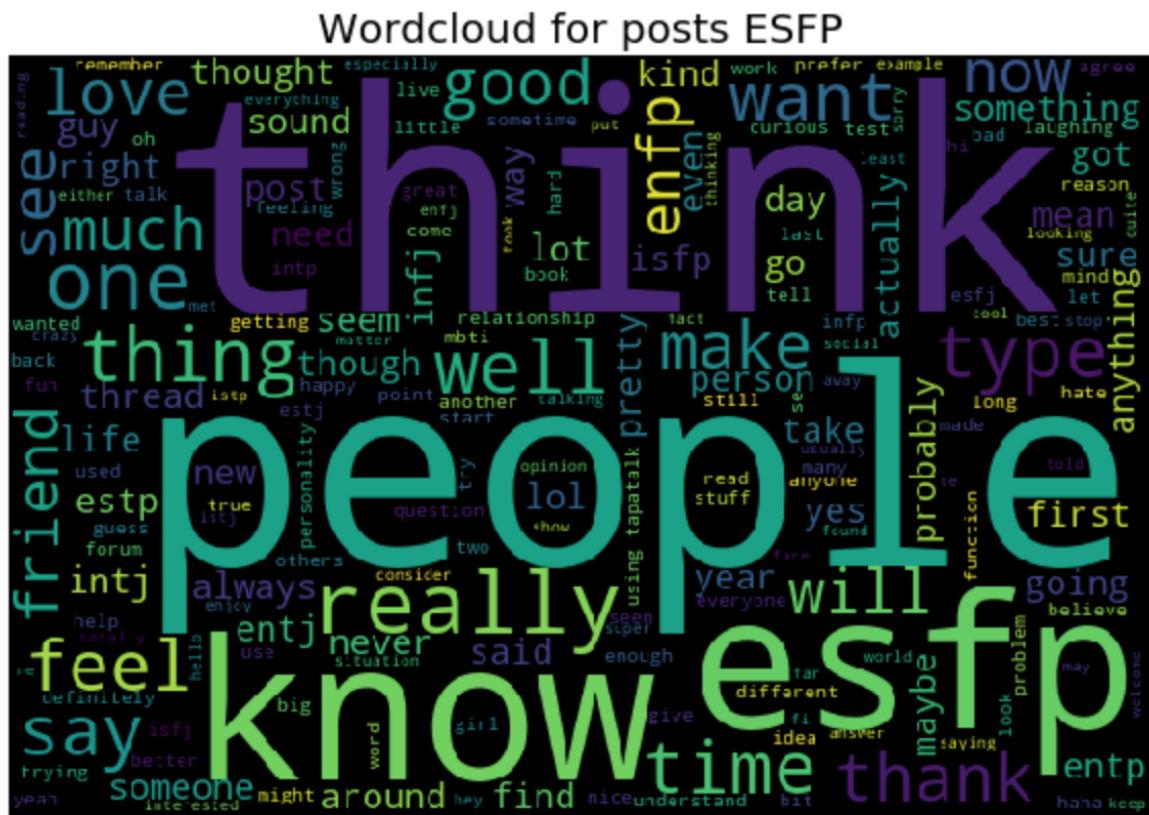
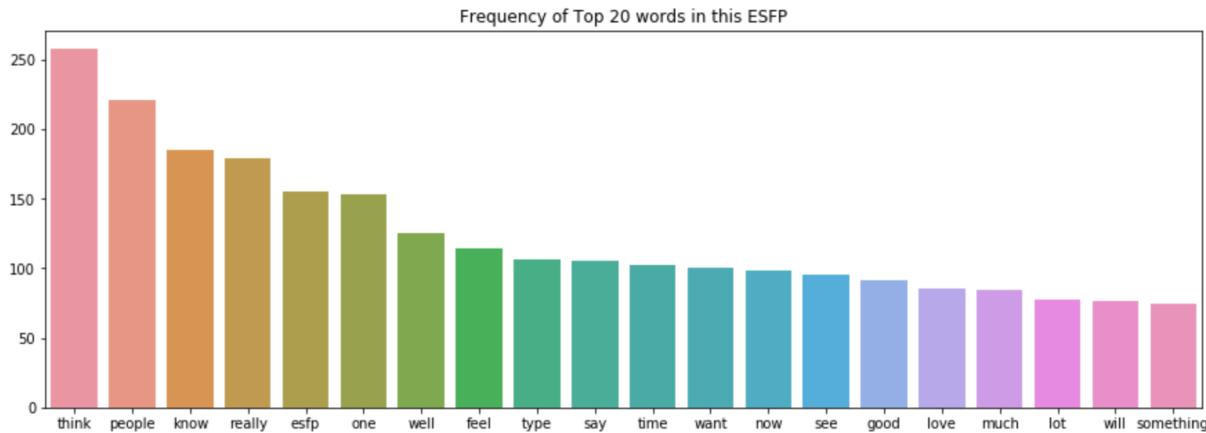
Wordcloud for posts ISTJ



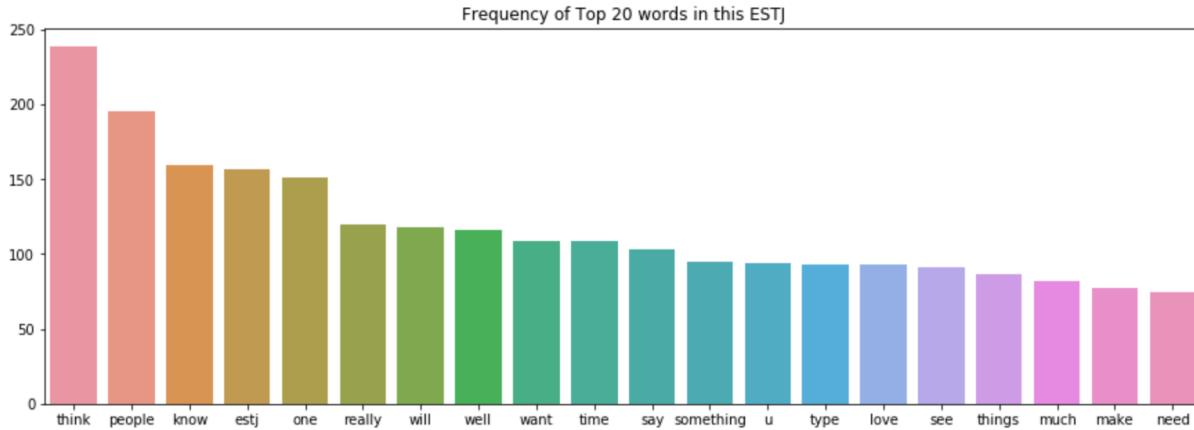
- ESTP



- ESFP



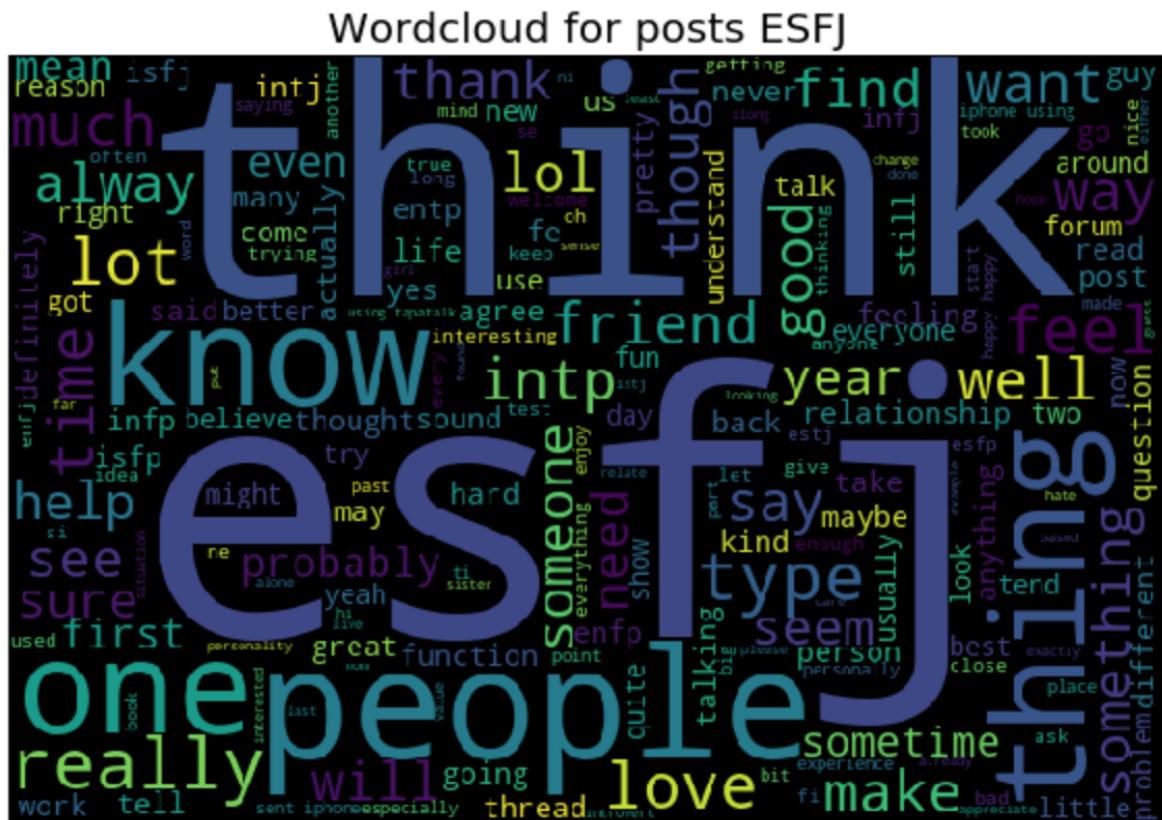
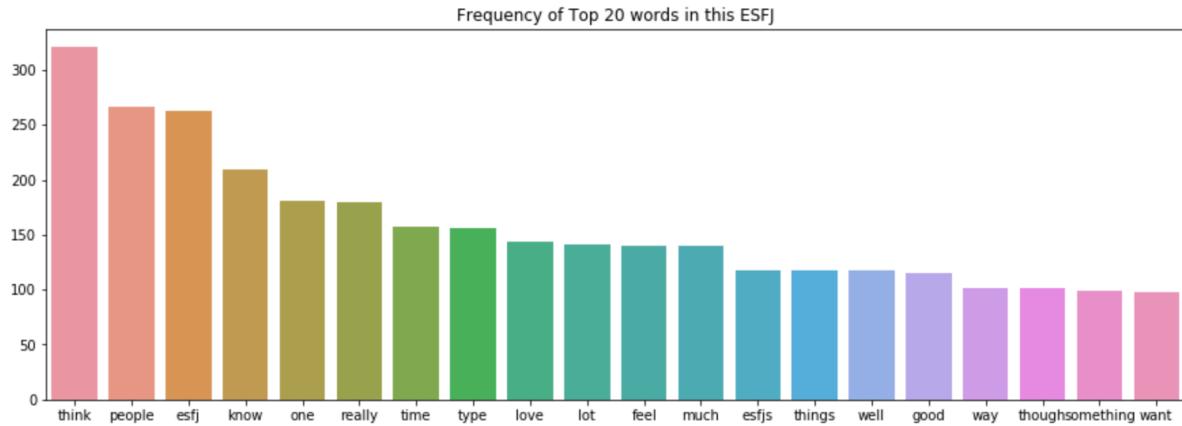
- ESTJ



Wordcloud for posts ESTJ



- ESFJ



Machine Learning Algorithms.

We are going to use **clean_posts** to predict personality type. First, we split data into train set and test set. The train set contains 75% of data points and the test set contains 25% of data points.

The **X_train**, **X_test** are sets of posts which are not the form that we can feed Machine Learning models. In order to apply Machine Learning models we need to have features expressed in numbers for **X**. We have to use **CountVectorizer** to convert the text documents to a vector of term/token counts.

We set **min_df=0.1** and **max_df=0.9**. It means we only consider words that appear in more than 10% of the posts and less than 90% of the posts. It makes sense here, since we do not want to take into account words that only appear in a few posts neither words that appear in almost every post.

count_train and **count_test** are sparse matrices whose columns are selected words and rows are posts. Each entry in the matrix is the frequency that the word (column) appears in the post(row). With **min_df=0.1** and **max_df=0.9**, we are going to apply the models(Logistic Regression, K Nearest Neighbors, Support Vector Machine, Decision Tree, Random Forest, Stochastic Gradient Descent and Naive Bayes) to the **count_train** and evaluate the accuracy of these models on **count_test**.

Models

Now we are ready to train a model and predict the required solution. Here,

y is value of target feature '**type_enc**'

and

X is a sparse matrix whose columns are selected words and rows are posts. Each entry in the matrix is the frequency that the word (column) appears in the post(row) which we will use to fit a machine learning model and predict the value of **y**.

There are so many predictive modelling algorithms out there to choose from. We must understand the type of problem and solution requirement to narrow down to a select few models which we can evaluate. Our problem is a classification problem. We are also performing a category of machine learning which is called supervised learning as we are training our model with a labeled dataset. With these two criteria - Supervised Learning plus Classification, we can narrow down our choice of models to a few. These include:

- K_Nearest_Neighbor
- Logistic Regression
- Support Vector Machine
- Stochastic Gradient Descent
- Random Forest
- Decision Tree
- Naive Bayes

Metrics

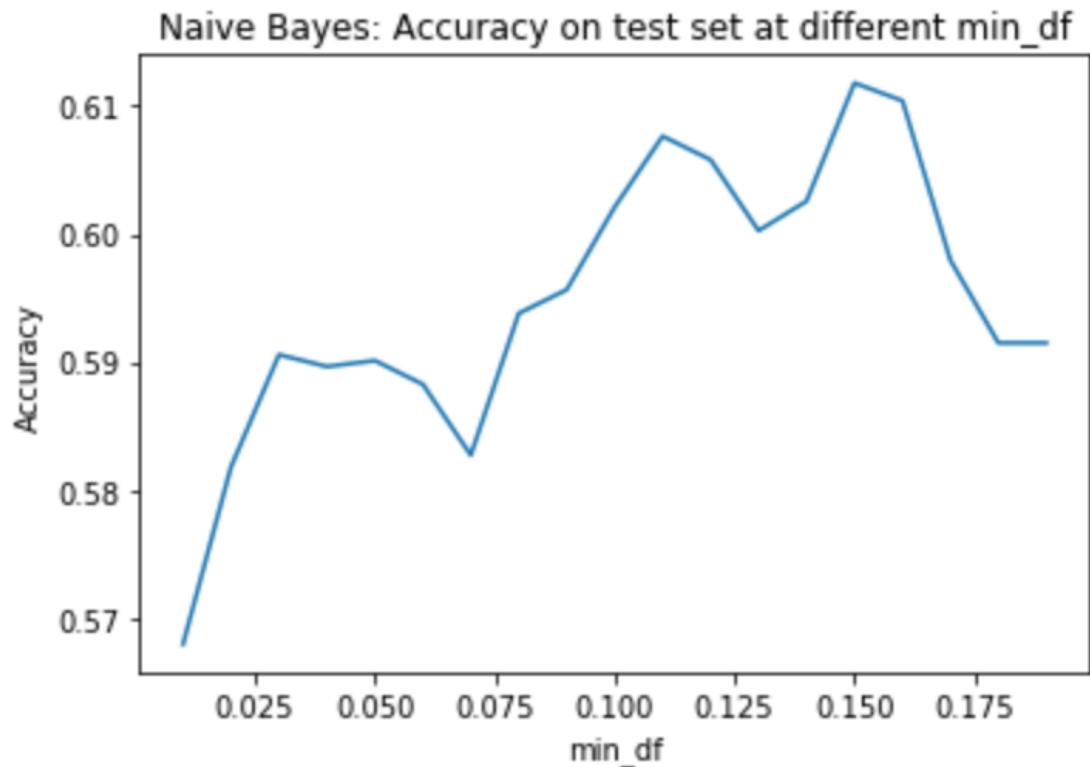
Beside the accuracy, we are going to use precision, recall, f1-score, AUC-ROC to evaluate our models.

| | Models | Accuracy | Precision | Recall | F1-score | AUC-ROC |
|--|-----------------------------|-------------|-------------|-------------|-------------|-------------|
| | Logistic Regression | 0.54 | 0.54 | 0.54 | 0.54 | 0.91 |
| | K-Nearest Neighbor | 0.36 | 0.36 | 0.36 | 0.35 | 0.77 |
| | Support Vector Machine | 0.56 | 0.56 | 0.56 | 0.56 | 0.93 |
| | Decision Tree | 0.47 | 0.47 | 0.47 | 0.47 | 0.72 |
| | Random Forest | 0.59 | 0.60 | 0.59 | 0.54 | 0.91 |
| | Stochastic Gradient Descent | 0.53 | 0.56 | 0.53 | 0.53 | 0.80 |
| | Naive Bayes | 0.60 | 0.62 | 0.60 | 0.61 | 0.92 |

Based on the summary table, Naive Bayes is the best model.

More about Naive Bayes

Out of all models we applied, Naive Bayes turns out to be the best model. We would like to look deeper into Naive Bayes model. We vary the values of `min_df` and `max_df` and apply the best model we have(Naive Bayes).



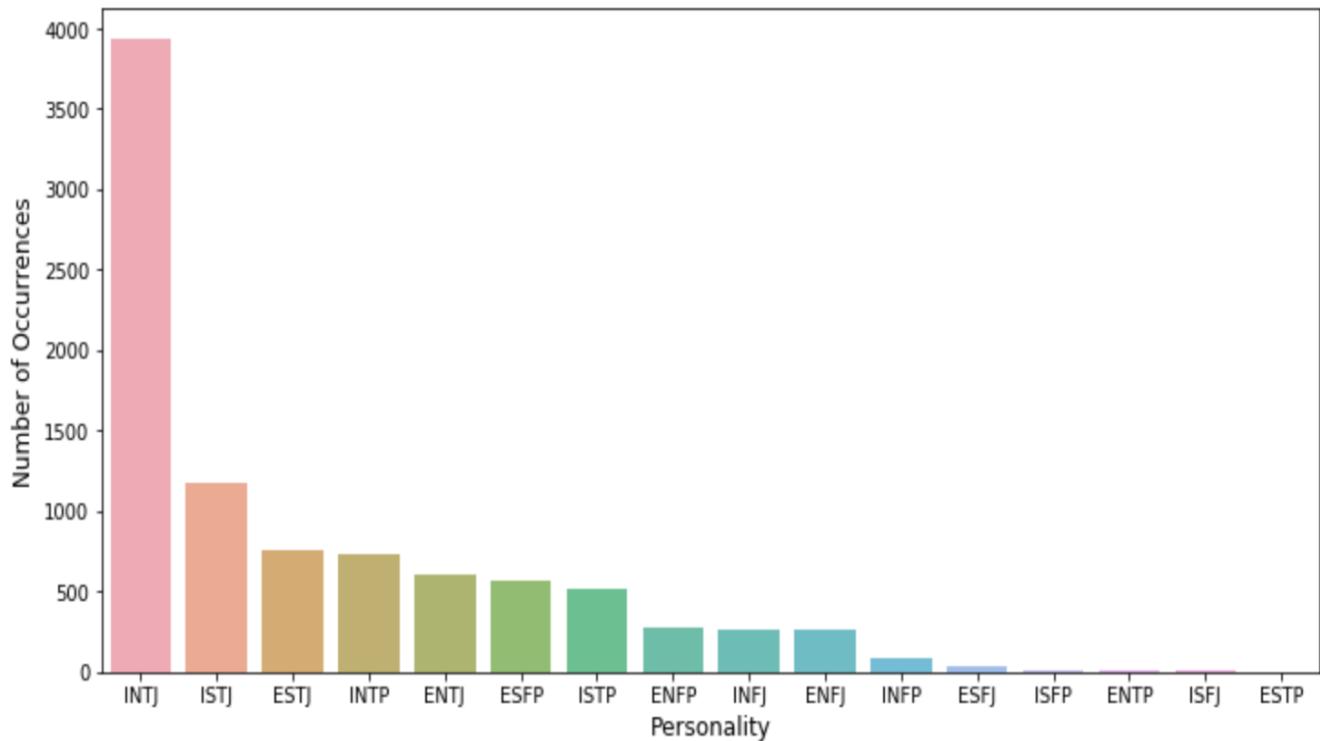
The highest accuracy is 0.612 when we only consider words that appear in more than 15% and less than 85% of the posts.

Predict Personality of Kaggle Users

We are going to apply our best model to find out Kaggle users personality based on their comments on Kaggle. The `ForumMessages.csv` is available on Kaggle at

<https://www.kaggle.com/kaggle/meta-kaggle?select=ForumMessages.csv>

| | PostUserId | Message | Message_length | clean_messages |
|---|------------|---|----------------|---|
| 1 | 368 | Here are some papers that analyze Eurovision v... | 32612 | papers analyze eurovision voting patterns migh... |
| 2 | 381 | <p>Hi Bdol.</p>\r\n<p>Please comment out the l... | 28821 | hi bdol please comment line dbclear error main... |
| 3 | 387 | <p>From an economic perspective let's look at ... | 526 | economic perspective look demand supply declar... |
| 6 | 393 | <p>[quote=Josette_BoozAllen;155818]</p>\n\n<p>... | 789 | thomasleleck aws credits opportunity will anno... |
| 7 | 412 | <p>Are there any criteria around model validit... | 789 | criteria around model validity nbsp e model ne... |



We observe that the majority of kaggle users have INTJ personality type.

Predict Personality of Public Figures

There are two US political figures that seem to stand for the opposite point of views: President Donald Trump and previous president Barack Obama. We will collect their posts and use our best model to predict their personality.

| | Name | Posts | posts_length | clean_posts |
|---|--------------|---|--------------|---|
| 0 | Barack Obama | Michelle and I have been spending a lot of tim... | 766 | michelle spending lot time together past month... |
| 1 | Donald Trump | No doubt many people told him his vision wasn... | 1146 | doubt many people told vision wasn t possible ... |

Based on our best model, they are both predicted to have INFP(Introversion-Intuition-Feeling-Perception) personality type.

Conclusion

People use social media to share their feelings and opinions. It is true that social media posts can reveal many things about the owner including his personality. Social media posts become golden data if they are studied correctly. My project focuses on using machine learning algorithms in supervised machine learning to predict the personality of a person from the type of posts they put on social media.

Psychologists can use this project to study personality. This project can also help companies who wish to learn more about their customers from their social media posts to provide appropriate services.