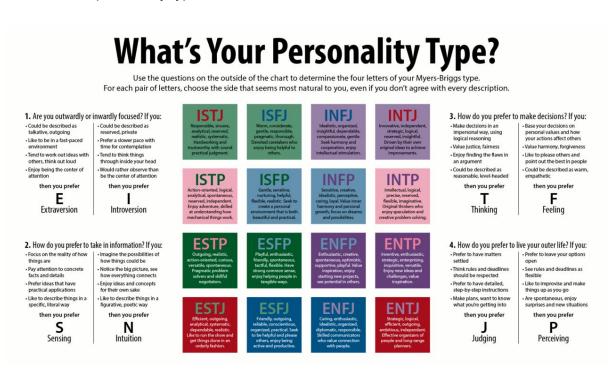**Project Outline**

**Determining MBTI's personality type based on user's posts.**

- **Explanation:** MBTI is a personality test indicator developed by Katharine Cook Briggs and Isabel Briggs Myers. It is based on the conceptual theory proposed by Swiss psychiatrist Carl Jung. It categorizes people into 16 personality types. My idea is to take the MBTI Kaggle dataset that classifies +8000 users from the website Personality Cafe into each of these 16 types and train a model to see if it is possible to determine personality types through written communication.

- **Problem Identification:** Uploaded on GitHub the problem identification slide which contains the following points:
    - Problem Statement
    - Context
    - Criteria for Success
    - Scope of solution space
    - Constraints within the solution space
    - Stakeholders to provide key insight
    - Key data sources

- **Framework:** The following image developed by Jake Beech provides a summary of the 16 personality types and the 4 axes.

- **Expected Process:** I will follow the Data Science Method proposed by Aiden V. Johnson although later in the process we might need to make adjustments with the newly found information. The process goes as follows:

  - **Problem Identification:** As aforementioned, this document provides the starting point where I explain my initial hypothesis.

  - **Data Collection, Organization, and Definitions:** This section is very simple within this project since the data is readily available and complete. I will take the opportunity to create new features from the available data.

  - **Exploratory Data Analysis:** I will check for correlations to determine if there are variables we can drop, also to see trends and patterns that might give us more ideas for feature creation.

  - **Pre-processing and Training Data Development:** There are many pre-processing options at our disposal, like separating the posts, doing lemmatization, removing symbols and stopwords, doing sentiments analysis, etc. At the moment we do not know what is the best process in terms of what pre-processing steps will yield the best results.

  - **Model Creation:** With the newly generated dataset I will try to classify each user into a personality type (target column) by training classification models. I am currently not versed in this field but I imagine I use things like Support Vector Machines, K-Nearest Neighbors, Random Forest Classifier, etc. The process will include:
    - Fit Models with Training Data Set.
    - Review Model Outcomes - Iterate over additional models as needed.
    - Identify the Final Model
    - Apply the Model to the Complete Data Set

  - **Documentation**
    - Review the Results — Share your findings
    - Finalize Code and Documentation

- **Further analysis**: Maybe it could be interesting to do some follow-up analysis in this project. Some ideas that come into mind:
  - **Expand dataset:** Since the dataset is skewed, it might be interesting to scrape some more posts from the profiles that have less labeled data.
  - **Enneagram of Big 5 Models**: do a similar analysis with the enneagram or the Big 5 models (other personality test types) and see if there is a relationship between the tests. This has been a longstanding topic of discussion and could be interesting to shed some light here.
  - **LIWC:** The Linguistic Inquiry and Word Count (LIWC) is a text analysis tool that differentiates documents according to broad themes expressed in writing: psychological, relativity, and contextual themes. LIWC was created by Dr. James Pennebaker in order to examine relationships between language and personality. (See: M. C. Komisin) Using this technique we could extract some interesting insights as well.