# Personality Prediction From Text Based on the MBTI Model

Andrel Chew

# Overview

**1** Introduction

**2** Data

**3** Methodologies

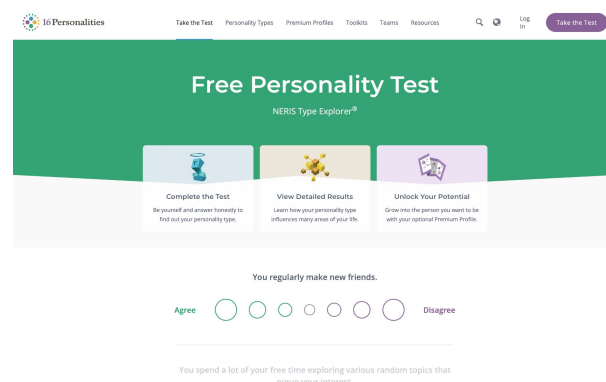**4** MBTI Prediction Tool

**5** Conclusion

# Introduction

- Myer Briggs Type Indicator - 16 personality types
- 4 Dimensions:
  - I/E Dimension: Introvert(I) or Extrovert(E)
  - N/S Dimension: Intuition(N) or Sensing(S)
  - T/F Dimension: Thinking(T) or Feeling(F)
  - J/P Dimension: Judging(J) or Perceiving(P)
- Applications:
  - Recommender Systems
  - Improve interpersonal relationships and job satisfaction
  - Criminal Profiling
- Other personality models: OCEAN, HEXACO, DiSC

# Problems

- Questionnaires are used to determine personality
- Biased Results:
  - Answering based on intended personality
  - Response bias in job interviews → fabricated answers
- Data quality issues in online surveys

# Objective & Scope

- Experiment with deep learning for automated personality prediction
- Mitigate shortcomings faced in online assessments.
  - Beyond Machine Learning - Comparing performance of neural networks, fine-tuned transformers models, multi-task learning etc.
- Implementation of an MBTI prediction tool
  - Purpose: to conveniently predict personality
  - Alternative to existing online personality assessments (eg. 16Personalities)
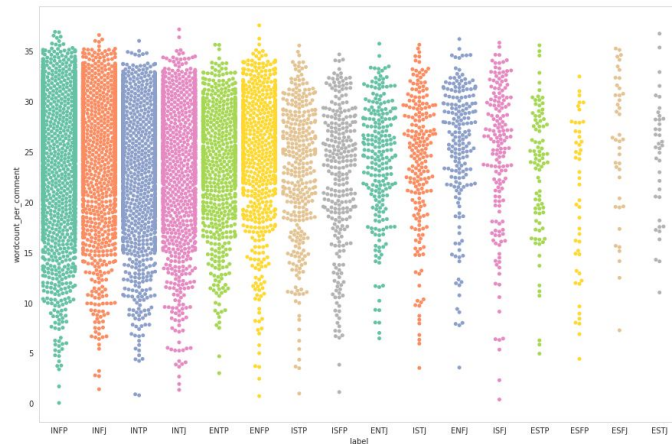
# Data

# Data Extraction



- Personality Cafe Forum dataset from Kaggle
- Reddit comments in 2019 from Google's Bigquery

# PersonalityCafe Data

- 8.6k rows of comments
- 50 comments per row, splits by "|||"

| | text | label |
|---|---|---|
| 0 | 'http://www.youtube.com/watch?v=qsXHcwe3krwIII... | INFJ |
| 1 | 'I'm finding the lack of me in these posts ver... | ENTP |
| 2 | 'Good one _____ https://www.youtube.com/wat... | INTP |
| 3 | 'Dear INTP, I enjoyed our conversation the o... | INTJ |
| 4 | 'You're fired.IIIThat's another silly misconce... | ENTJ |

# Reddit Data

- 2 million rows of comments
- Requires extraction of MBTI types from Reddit flairs

| | flair | text | subreddit | user |
|---|---|---|---|---|
| 0 | eStJ gAnG | I didn't reach puberty until I lead my Empire's armies on a series of conquests. Vanquishing all who stood in my path. | shittyMBTI | gggggggggee |
| 1 | ESTJ: The Supervisor | I didn't reach puberty until I lead my Empire's armies on a series of conquests. Vanquishing all who stood in my path. | shittyMBTI | gggggggggee |
| 2 | INFP: The Dreamer | This map is fucking swell compared to Euphrates Bridge. \n\nThat map is an abomination through and through. | modernwarfare | DankMatter3000 |
| 3 | INFP: The Dark Lady | &gt;I definitely had the thought that it would be nice if rustc knew that the way I was blocking should make my code race-free\n\nWell, there's Rice's Theorem. Informally, any program that decides whether a program has a particular property must have at least one flaw:\n\n* it can't analyze all programs\n* it sometimes makes mistakes\n* it sometimes fails to find an answer\n\nIf you *do* manage to write a flawless analyzer, then the property must be "trivial" - always true or always false. It's kinda like the Second Law of Compiledynamics.\n\nSafe Rust is intended to have two of those flaws:\n\n* it sometimes rejects programs that would be safe (a "soundness bug" is when it does the opposite - that's *not* intended but does happen)\n* type-checking isn't guaranteed to terminate and can even be tricked into performing arbitrarily complex computation, such as this [fractal type error](http://www.treblig.org/daveG/rust-mand.html). \n\nSo the "no data races" guarantee comes with the... | rust | claire_resurgent |
| 4 | INFP: The Dreamer | Or just has long hair... | memes | Sgt-Thunder-3 |

| | text | label |
|---|---|---|
| 0 | Hitman please! Favorite game is Persona 5 | INFJ |
| 1 | If there were a dedicated ARAM mode people wou... | INFJ |
| 2 | It's obscure game, but Starmade might be more ... | ENTP |
| 3 | I can't even count the reasons why it's comple... | INTP |
| 4 | Thanks for the tip. Please forgive me, I am st... | ENTJ |

# Data Cleaning

- Duplicated rows
- Lowercasing
- Punctuations
- URLs
- Numbers
- Special symbols (eg. emojis)
- Stopwords
- Lemmanization (NLTK)

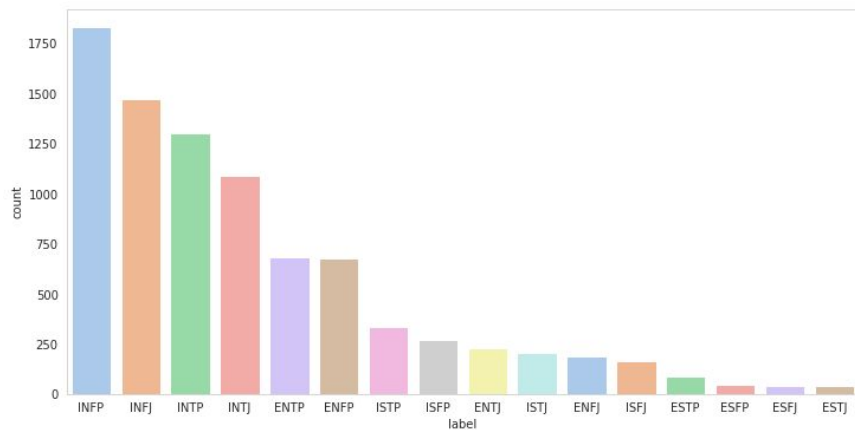|   | text | label |
|---|------|-------|
| 0 | intj moment play experience life repeat today ... | INFJ |
| 1 | find lack post alarm sex boring position often... | ENTP |
| 2 | good course say know absolutely positive good ... | INTP |
| 3 | dear rule arbitrary construct create dear entj... | INTJ |
| 4 | silly misconception approach logically go key ... | ENTJ |

PersonalityCafe Dataset Cleaned

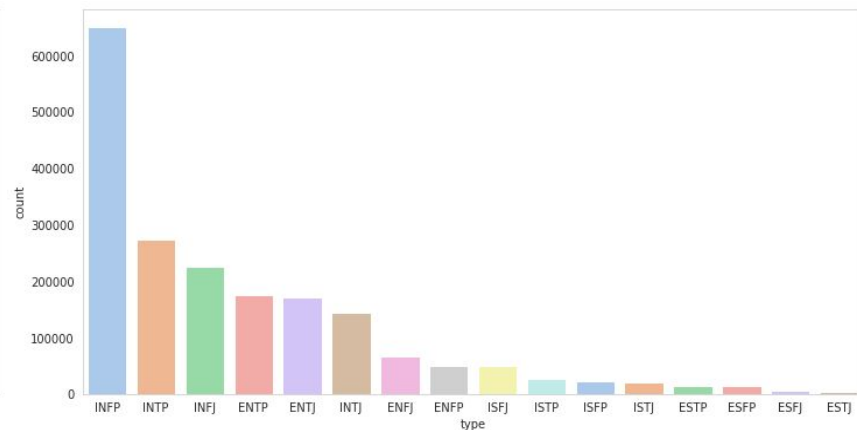|   | text | label |
|---|------|-------|
| 0 | favorite game persona | INFJ |
| 1 | dedicated mode people would stop play mode | INFJ |
| 2 | obscure may alley actually fly fight ship buil... | ENTP |
| 3 | even count reason completely insane | INTP |
| 4 | thank tip forgive stupid | ENTJ |

Reddit Dataset Cleaned
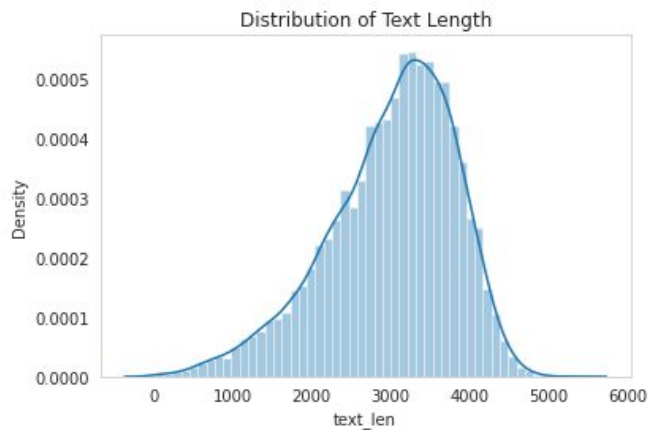
# Analysis

# PersonalityCafe vs Reddit



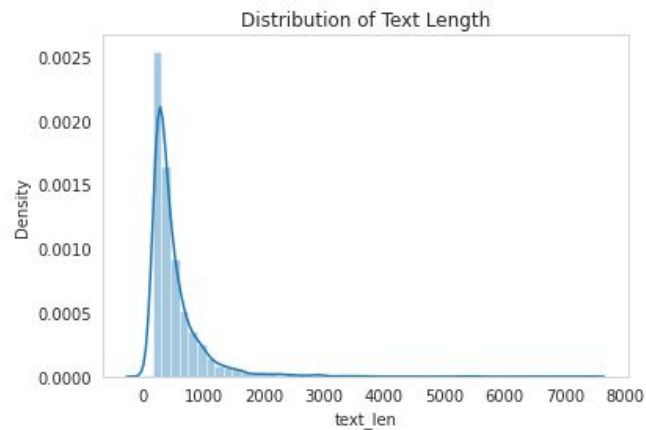Distribution of PersonalityCafe Dataset

Distribution of Reddit Dataset

# PersonalityCafe vs Reddit



Distribution of PersonalityCafe Dataset
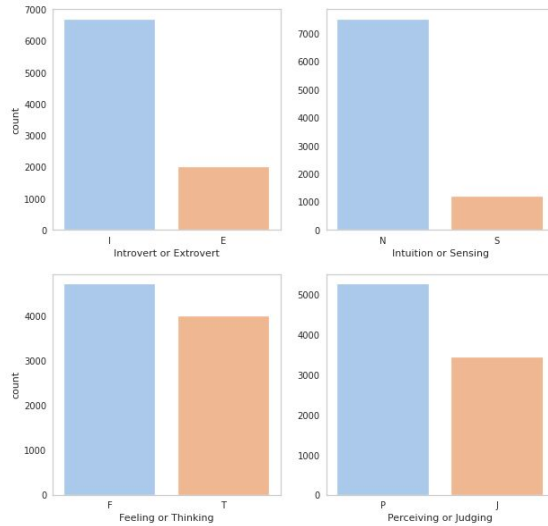
Distribution of Reddit Dataset

# Exploratory Classification Comparison
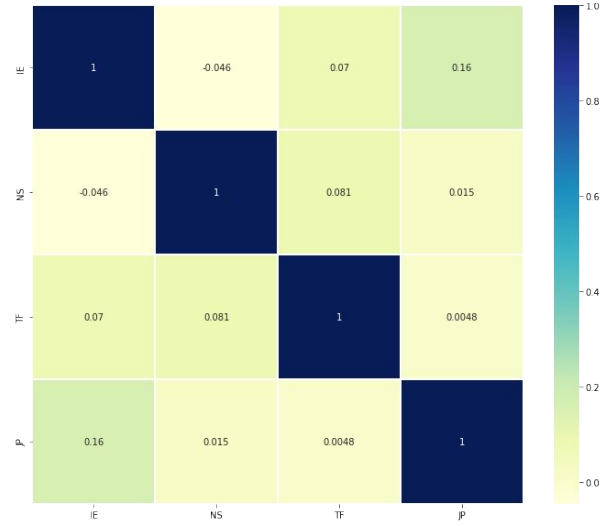
**Proceeding with P.Cafe dataset!**

| Model | 16 Types Multiclass (PersonalityCafe Dataset) | 16 Types Multiclass (Reddit Dataset) |
|---|---|---|
| Adaboost | 0.2324 | **0.0865** |
| CatBoost | 0.2422 | 0.0843 |
| Logistic Regression | **0.2715** | 0.0842 |
| Naive Bayes | 0.2635 | 0.0849 |
| Random Forest | 0.2091 | 0.0740 |
| XGBoost | 0.2406 | 0.0813 |

Exploratory Classification With
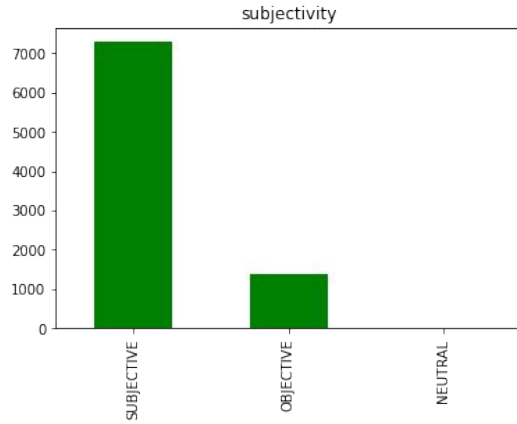Bag of Words

# Data Analysis on PersonalityCafe Data
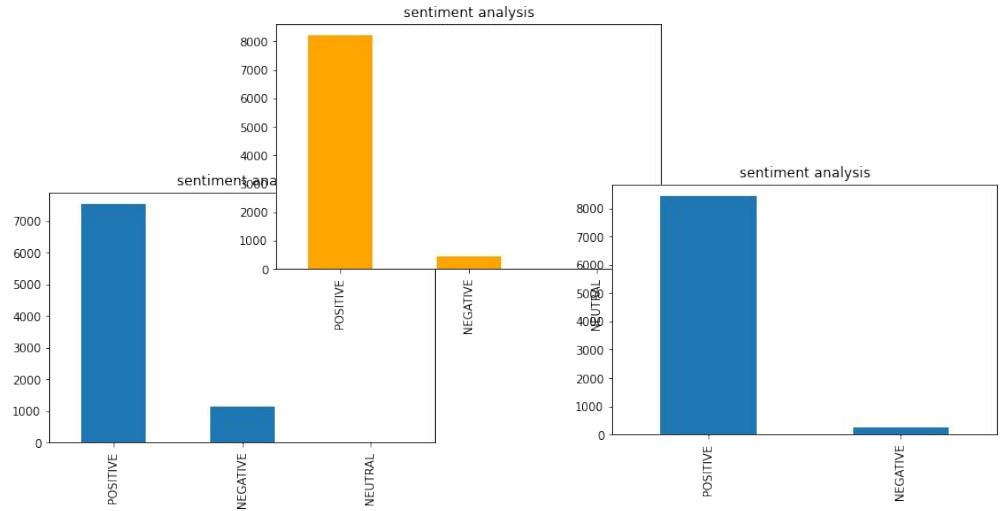


Distribution of 4 Dimensions

Correlation Matrix of the 4 Dimensions

Before Cleaning

After Cleaning

4 Dimensions

16 Personalities
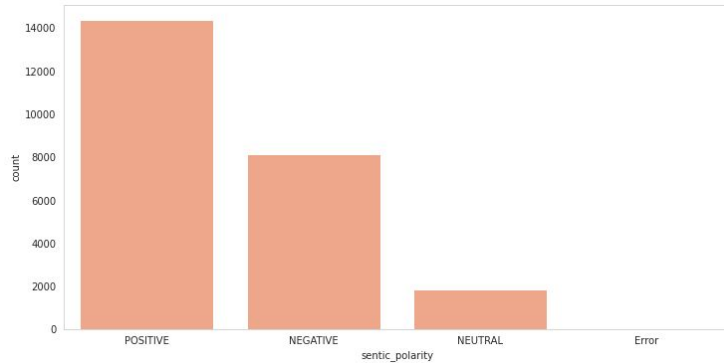
# Subjectivity Detection



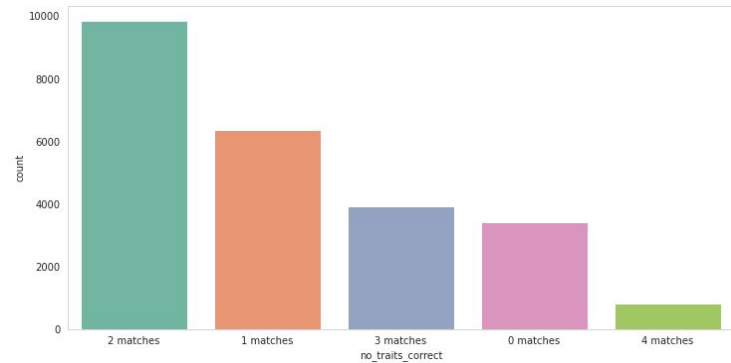# Sentiment Analysis



Processing using SenticNet, NLTK & TextBlob

# Analysis of SenticNet APIs
# (with 24k Reddit Data)



SenticNet's Polarity Classification



Matches of SenticNet's MBTI Prediction

```
print(sentic['sentic_polarity'][sentic['sentic_polarity'].str.contains('414 Request-URI Too Long')])

1938     <!DOCTYPE HTML PUBLIC "-//IETF//DTD HTML 2.0//EN">\n<html><head>\n<title>414 Request-URI Too Long</title>\n</head><body>\n<h1>Request-URI Too Long</h1>
11627    <!DOCTYPE HTML PUBLIC "-//IETF//DTD HTML 2.0//EN">\n<html><head>\n<title>414 Request-URI Too Long</title>\n</head><body>\n<h1>Request-URI Too Long</h1>
23659    <!DOCTYPE HTML PUBLIC "-//IETF//DTD HTML 2.0//EN">\n<html><head>\n<title>414 Request-URI Too Long</title>\n</head><body>\n<h1>Request-URI Too Long</h1>
Name: sentic_polarity, dtype: object
```

URL Error for long text

# Methodologies

# Methods

**Machine Learning**

LDA, SMOTE Experiment,
16 Classification Models

**1**

**Transformers &
Transfer Learning**

BERT, DistilBERT,
ZeroShot Pipeline

**2**

**Neural Networks**

GloVe Embeddings with
CNN, LSTM, GRU etc.

**3**

**Ensemble Learning**

Binary & Multiclass
Neural Network Ensemble

**4**

**Multi-task Learning**

MBTI, Sentiment, Subjectivity

**5**

**Hybrid Model**

DistilBERT Embeddings,
Self-Attention Mechanism

**6**

# Types of Classification



**#1**

**Multiclass**
16 MBTI Types

INFP

ESTJ

**#2**

**Binary**
4 MBTI Axes

I/E

N/S

T/F

J/P

# Machine Learning
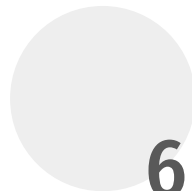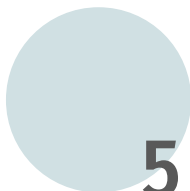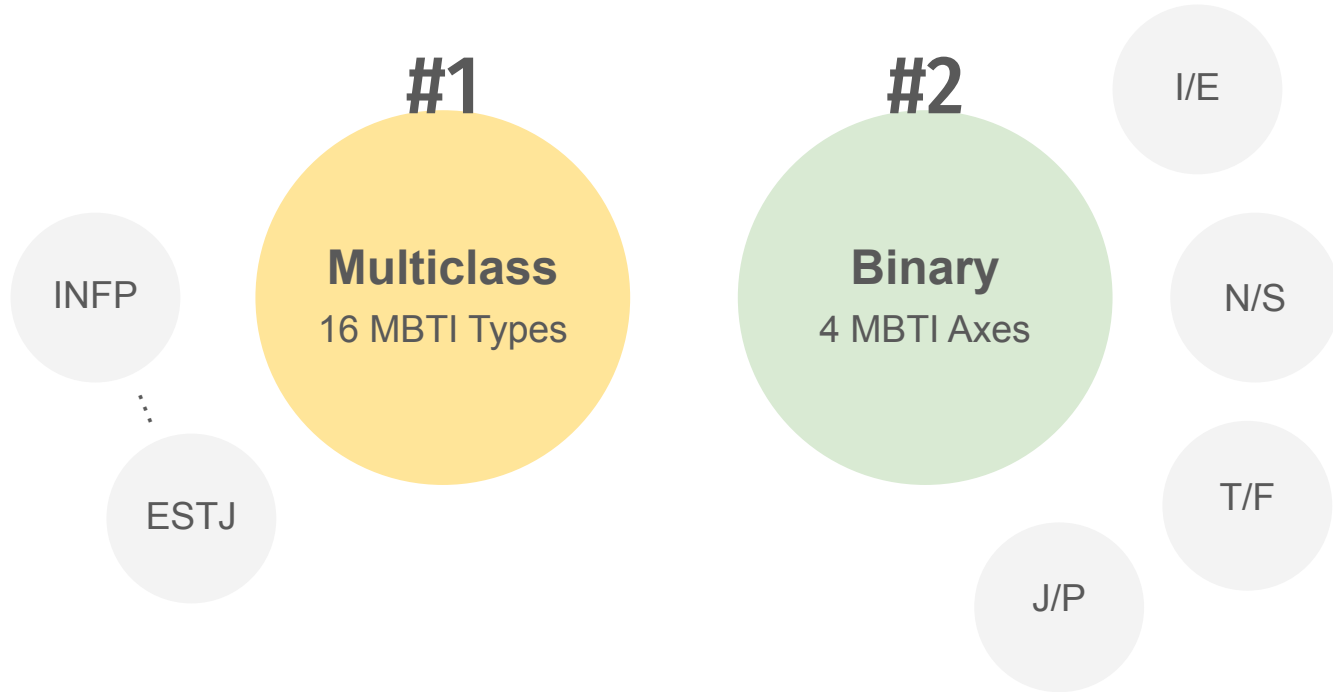
- LDA Topic Modelling on 4 topics
- Class Balancing techniques:
  - SMOTE, Random Oversampling, Random Undersampling & SMOTEENN
- 16 Machine Learning Models with Pycaret
  - Best models: Logistic Regression, Naive Bayes, SVM

| | text | label | IE | NS | TF | JP | Topic_0 | Topic_1 | Topic_2 | Topic_3 | Dominant_Topic |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | intj moment play experience life repeat today ... | INFJ | I | N | F | J | 0.001276 | 0.695496 | 0.159435 | 0.143793 | Topic 1 |
| 1 | find lack post alarm sex boring position often... | ENTP | E | N | T | P | 0.299757 | 0.366157 | 0.305049 | 0.029037 | Topic 1 |
| 2 | good course say know absolutely positive good ... | INTP | I | N | T | P | 0.246211 | 0.495362 | 0.123393 | 0.135034 | Topic 1 |
| 3 | dear rule arbitrary construct create dear entj... | INTJ | I | N | T | J | 0.254627 | 0.000650 | 0.728423 | 0.016300 | Topic 2 |
| 4 | silly misconception approach logically go key ... | ENTJ | E | N | T | J | 0.117287 | 0.440013 | 0.442000 | 0.000700 | Topic 2 |

Latent Dirichlet Allocation (LDA)

| | Topic_0 | Topic_1 | Topic_2 | Topic_3 | IE | Label |
|---|---|---|---|---|---|---|
| 0 | 0.551252 | 0.286153 | 0.162140 | 0.000454 | I | I |
| 1 | 0.269050 | 0.005653 | 0.724851 | 0.000446 | I | I |
| 2 | 0.520067 | 0.077223 | 0.168415 | 0.234296 | E | I |
| 3 | 0.288667 | 0.202691 | 0.508182 | 0.000461 | I | I |
| 4 | 0.393778 | 0.604691 | 0.000765 | 0.000765 | I | I |

Example of predicted results for
I/E Axis (Logic Regression)

# Transformers & Transfer Learning

- Fine-tuning of pre-trained models
  - BERT
  - DistilBERT -- best model

- Zero-Shot Learning
  - BART-large-mnli model pipeline

```
***** Running training *****
  Num examples = 6244
  Num Epochs = 3
  Instantaneous batch size per device = 16
  Total train batch size (w. parallel, distributed & accumulation) = 16
  Gradient Accumulation steps = 1
  Total optimization steps = 1173
```
[1173/1173 30:09, Epoch 3/3]

| Step | Training Loss |
|------|---------------|
| 500  | 1.967200      |
| 1000 | 1.393800      |

Fine-tuning DistilBERT

# Neural Networks



- 6 Models
  - CNN -- best model
  - GRU
  - Bi-GRU
  - LSTM
  - Bi-LSTM
  - MLP
- Training on 4 Binary & 1 Multiclass Classifiers
  - GloVe word embeddings
  - Binary model - sigmoid activation function, binary cross-entropy
  - Multiclass model - softmax activation function, categorical cross-entropy
  - Dropout layers to prevent overfitting
- Difficulty in predicting personality on dataset with class imbalance

# Model Enhancements

# Ensemble Learning

- Stacking Ensemble Approach
- Ensemble of top 3 models in both multiclass & binary classifiers
- Requires model retraining
  - Used non-independent data when training DNN
  - New data allocation with holdout dataset
  - Controlled random states of training data
- Observed better performance with standalone models
  - Due to weak base classifiers < 50% accuracy
  - Lesser training data with new allocation

# Data Allocation with Controlled Random States

# Stack Ensemble with Multiclass Classifiers

# Multi-Task Learning

- Hard-parameter sharing MTL
- Learn tasks simultaneously to increase efficiency
- 2 groups of MTL Tasks
  - 1: MBTI with 4 dimensions
  - 2: MBTI with 16 types, Sentiment & Subjectivity
- 2 Models
  - Keras Tokenizer + MLP
  - GloVe Embeddings + CNN

# Multi-Task Learning

# Data Labelling for MTL

- Sentiment
  - Fill undetected sentiment when processing with VADER (server was down)
  - Checking SenticNet with VADER → 82.76% similarity

- Subjectivity
  - Single Row of Data with Neutral Subjectivity will cause issues with training
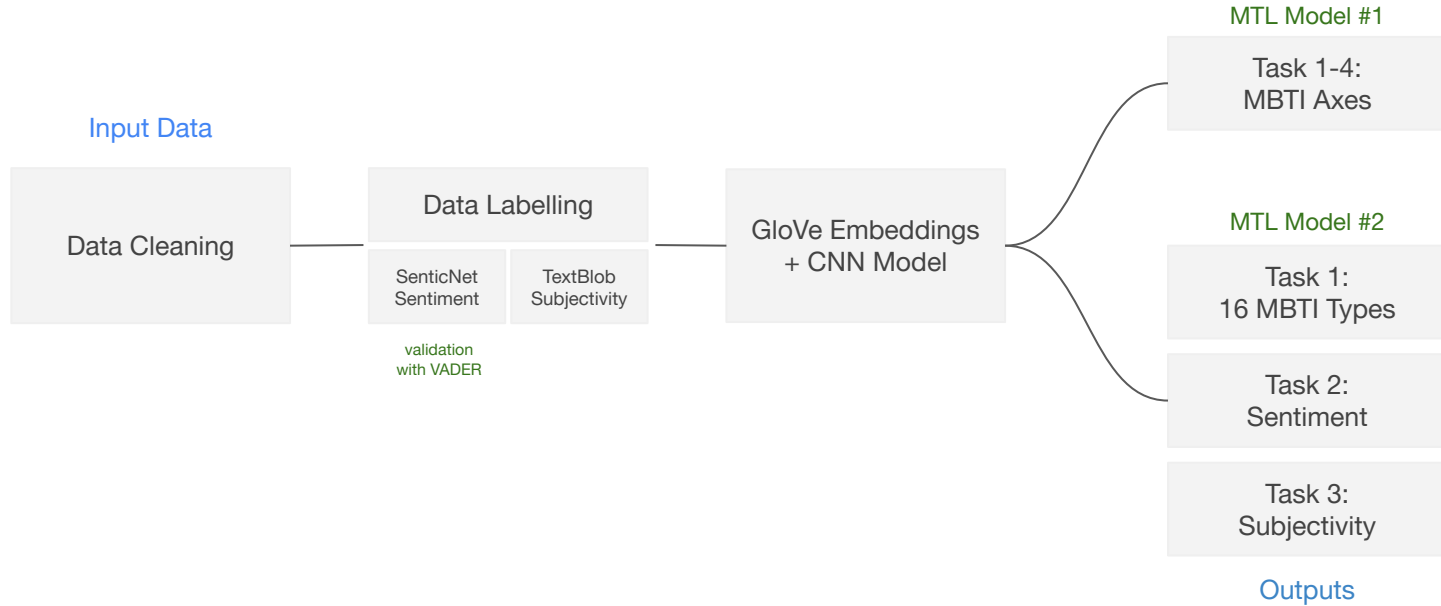  - 2 classes instead of 3 classes (either Subjective or Objective)

| | text | mbti | sentiment |
|---|---|---|---|
| 0 | intj moment play experience life repeat today ... | INFJ | POSITIVE |
| 1 | find lack post alarm sex boring position often... | ENTP | POSITIVE |
| 2 | good course say know absolutely positive good ... | INTP | POSITIVE |
| 3 | dear rule arbitrary construct create dear entj... | INTJ | POSITIVE |
| 4 | silly misconception approach logically go key ... | ENTJ | POSITIVE |

```
text            suppose must live election night ad
mbti                                           INTP
sentiment                                  POSITIVE
subjectivity                                NEUTRAL
Name: 6277, dtype: object
```

# Hybrid Model with Self-Attention Mechanism

- Using best models: DistilBERT fine-tuned, GloVe-CNN Neural Network
  - With a Self-Attention Mechanism - help to retain semantic information
- 2 Hybrid Models
  - DistilBERT embeddings + CNN
  - DistilBERT embeddings + CNN + Scaled Dot Product Self-Attention
- Better Accuracy + More classes predicted with GloVe-CNN
  - Therefore GloVe > DistilBERT embeddings

Input Data

DistilBERT Embeddings

Data Cleaning

CNN Model

Predicted MBTI Type

Scaled Dot-Product Self-Attention

Output

# MBTI Prediction Tool

- Build with Streamlit
    - Prediction Tool & Data Visualization pages
- 2 Models trained on 2 different datasets:
    - GloVe + CNN on PersonalityCafe dataset
    - GloVe + CNN on Reddit dataset
- Simple text processing: remove digits, punctuations, lowercase
- Returns MBTI + Subjectivity (TextBlob) + Sentiment (VADER)
    - Multitask Model did not perform well

Hmm..

The genie thinks you're an INFP + the words are Subjective and Neutral🙂

SCAN ME

# Deployment

# MBTI Prediction Page



**MBTI Prediction Tool**

**Let's guess your MBTI...**

Enter some text here

dancing in the moonlight

Choose a Model for Prediction
- ● CNN + PersonalityCafe Data
- ○ CNN + Reddit Data

**Hmm..**

The genie thinks you're an INFP + the words are Subjective and Neutral😊

Made with Streamlit

MBTI Prediction

**MBTI Prediction Tool**

**Let's guess your MBTI...**

Enter some text here

dancing

Choose a Model for Prediction
- ● CNN + PersonalityCafe Data
- ○ CNN + Reddit Data

Invalid text! Enter text with more than 10 letters

Made with Streamlit

Error Handling:
for text length < 10

# Data Visualization Page



Option 1: Text Analysis

Option 2: Personality Types

# Evaluation on Crawled Tweets

- 3K Crawled Tweets with mention of 'food'
- Predictions on both CNN models deployed
- Results:
  - Reddit Model had predictions from all 16 MBTI types
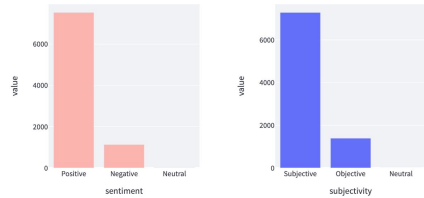  - PersonalityCafe model had predictions from 11 different MBTI types (all undetected are extroverts - probably not enough training data)
  - 4.897% of rows have identical matches

```
food_df.prediction_pcafe.value_counts()

INFP    1637
INTP     785
ENFP     520
INFJ     478
INTJ     451
ENTP      48
ISTP      32
ENTJ      12
ISTJ      11
ISFJ       5
ISFP       3
Name: prediction_pcafe, dtype: int64
```

```
food_df.prediction_reddit.value_counts()

ENTJ    1183
INTJ     657
ENFP     519
INFJ     370
INTP     218
ESFP     160
ESFJ     155
ENFJ     151
ISTJ     133
INFP     100
ESTJ      83
ISFJ      72
ISTP      70
ESTP      54
ISFP      45
ENTP      12
Name: prediction_reddit, dtype: int64
```

| | text | prediction_pcafe | prediction_reddit |
|---|---|---|---|
| 0 | recipe of the day creamy parmesan baked acorn ... | INFP | INTJ |
| 2 | this restaurants service felt very much racist... | INTJ | ENFJ |
| 3 | food pics and ambience is captured as well goo... | INTJ | ESFJ |
| 4 | ktla lapublichealth why do you guys keep raid... | INFP | ENFP |
| 5 | what a performance diggins says she fought ... | INFP | ENTP |

- Conclusion: Imbalanced dataset gives more accurate (eg. INFP with the most data), but fragmentary predictions (with 5 unpredicted classes)

# Conclusion

# Conclusion: Summary of Findings

- 6 Classification Approaches + MBTI Prediction Tool
  - Dataset with class imbalance issue - difficult to learn efficiently (INFP vs ESFJ)
  - Option of using 4 dimensions may offer better analysis for class imbalance datasets
- Findings
  - Best <u>Binary</u>: Fine-tuned DistilBERT; <u>Multiclass</u>: GloVe-CNN neural network
  - MTL can boost efficiency in saving time: training more tasks with similar accuracy
  - Hybrid Model Results: GloVe > DistilBERT embeddings
- Future Work: MTL with other NLP tasks (eg. NER, POS) & new datasets

| Binary Models | I/E | N/S | T/F | J/P | Acc. (%) | Multiclass Models | Acc. (%) |
|---|---|---|---|---|---|---|---|
| Fine-tuned DistilBERT | 83.6 | 89.9 | 83.2 | 77.4 | <u>**48.4**</u> | Glove-CNN (NN) | <u>**43.6**</u> |
| Glove-CNN (NN) | 80.0 | 87.3 | 74.9 | 69.1 | 36.1 | Glove-LSTM (NN) | 40.2 |
| Glove-CNN (MTL) | 76.3 | 84.7 | 53.8 | 54.8 | 19.1 | DistilBERT-CNN (Hybrid) | 33.6 |

Better than Random Guessing of 1/16 = <u>6.25%</u>

# Project Gantt Chart

# Limitations

- Class Imbalance of PersonalityCafe Data
  - Consideration of Reddit Dataset → lower accuracy in exploratory classification
  - No holdout data → model retraining for ensemble learning model
- Memory/GPU limitations
  - Used several google colab notebooks for classification
  - Unable to use entire Reddit dataset for processing (contains 2 mil rows)

# Future Work

- Multi-task Learning with other NLP tasks
  - Eg. Aspect-Based Sentiment Analysis (ABSA), Part of Speech (POS) Tagging, Named Entity Recognition (NER).
- Advanced state-of-the-arts transformers / word embeddings for improvements
  - Eg. GPT-NEO with text-generation may predict words choices of each personality type
  - Allows better understanding of human psychological aspects
- Other MBTI datasets
  - Labels with MBTI, Big Five, horoscopes, relationships, career paths etc.
  - Eg. INFP x ENFP relationship compatibility

# Thank You🙂

MBTI Characters Reference: 16Personalities. "Free Personality Test," 16personalities.com. [Online]. Available: https://www.16personalities.com. [Accessed: 20-Mar-2022].

ARCHIVED

# Gantt Chart

| Task | Aug 21 | | | | Sep 21 | | | | Oct 21 | | | | Nov 21 | | | | Dec 21 | | | | Jan 22 | | | | Feb 22 | | | | Mar 22 | | | | Apr 22 | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | 1 | 2 | 3 | 4 | 1 | 2 | 3 | 4 | 1 | 2 | 3 | 4 | 1 | 2 | 3 | 4 | 1 | 2 | 3 | 4 | 1 | 2 | 3 | 4 | 1 | 2 | 3 | 4 | 1 | 2 | 3 | 4 | 1 | 2 | 3 | 4 |
| Data Preprocessing | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| Data Analysis | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| Data Classification | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| Model Improvements | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| MBTI Prediction Tool | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| Documentation | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |