Personality Profiling: How informative are social media profiles in predicting personal information?

Joshua Watt, Jonathan Tuke and Lewis Mitchell
School of Computer & Mathematical Sciences, The University of Adelaide, Adelaide SA 5005, Australia
{joshua.watt,simon.tuke,lewis.mitchell}@adelaide.edu.au

Abstract—Personality profiling has been utilised by companies for targeted advertising, political campaigns and vaccine campaigns. However, the accuracy and versatility of such models still remains relatively unknown. Consequently, we aim to explore the extent to which peoples' online digital footprints can be used to profile their Myers-Briggs personality type. We analyse and compare the results of four models: logistic regression, naive Bayes, support vector machines (SVMs) and random forests. We discover that a SVM model achieves the best accuracy of 20.95% for predicting someones complete personality type. However, logistic regression models perform only marginally worse and are significantly faster to train and perform predictions. We discover that many labelled datasets present substantial class imbalances of personal characteristics on social media, including our own. As a result, we highlight the need for attentive consideration when reporting model performance on these datasets and compare a number of methods for fixing the class-imbalance problems. Moreover, we develop a statistical framework for assessing the importance of different sets of features in our models. We discover some features to be more informative than others in the Intuitive/Sensory (p = 0.032) and Thinking/Feeling (p = 0.019) models. While we apply these methods to Myers-Briggs personality profiling, they could be more generally used for any labelling of individuals on social media.

Index Terms—personality profiling, targeted advertising, machine learning, social media, statistical framework, digital footprints, data science, natural language processing

I. INTRODUCTION

In 2023 there are thousands of social media applications and over 4.59 billion social media users worldwide, constituting approximately 60% of the world's population [1]. While this enables most of the world to be connected, it also creates an environment of mass data, defining what we refer to as the information environment. The huge amounts of individual-level data provided by each user is an important aspect of social media which is unique to this type of information environment. Consequently, it is crucial for scholars to understand how this aspect of social media may impact society. There exists a need to quantify the extent to which social media can be weaponized by governments and other organisations for influence.

Every time a user enters a social media application, they leave a unique data trace – information they have posted, liked, shared, commented, even how long they have spent viewing different material on the application. We refer to this unique trace of data as a user's online digital footprint. It has been suggested that someone's online digital footprint can expose actionable information about them; including their personality profile, relationship status, political opinions and even their

propensity to adopt a particular opinion or behavior [2, 3, 4, 5, 6, 7]. Cambridge Analytica was suggested to use online digital footprints to impact the result of the 2016 US election and the 2016 Brexit referendum [2]. However, the extent to which companies like Cambridge Analytica can determine this information from social media data is still questioned [3, 4, 5]. As a result, it is of interest for individuals to understand the extent of information that is attainable from their online digital footprint. This is also of key concern for governments, who seek to maintain democracies and the ethical use of such data, which can be abused by understanding personal data.

We seek to determine how informative online digital footprints are in predicting Myers-Briggs personality types. This is a theoretical model comprised of four traits/dichotomies, based on Jungian theory [8, 9]. Modelling personal information about individuals using their online information has previously enabled researchers to understand the accuracy of such models. We extend this work by creating a new labelled dataset of Myers-Briggs personality types on Twitter and a statistical modelling framework which can be generally applied to any labelled characteristic of online accounts. We aim to reconsider the personality profiling and political microtargetting performed by companies like Cambridge Analytica.

First we collect a labelled dataset of accounts with self-reported Myers-Briggs personality types. We then collect a number of different features for these accounts including social metadata features and linguistic features: LIWC [10]; VADER [11]; BERT [12]; and Botometer [13]. We then create independent logistic regression (LR), naive Bayes (NB), support vector machines (SVMs) and random forests (RF) models on each dichotomy to model the Myers-Briggs personality type of the accounts. As part of this, we consider four different weighting/sampling techniques to adjust for class imbalances. Lastly, we provide a statistical framework for analysing the importance of different features in these models. We consider the importance of features at an individual level and across groups of features for each dichotomy. Our main contributions are:

 A labelled dataset¹ of 68,958 Twitter users along with their Myers-Briggs personality types, the largest available dataset (to our knowledge) of labelled Myers-Briggs personality types on Twitter [14].

¹Dataset available at https://figshare.com/articles/dataset/Self-Reported_Myers-Briggs_Personality_Types_on_Twitter/23620554.

- A statistical framework to combine NLP tools and mathematical models to predict online users' personality types, which can be more broadly used to model any labelled characteristics about online accounts.
- A comparison of machine learning models on NLP features, and a comparison of various weighting/sampling techniques to address problems with class imbalance.
- Statistical methods which compare the importance of different features in NLP-based models at an individual level and across groups of features.

II. BACKGROUND

Myers-Briggs [8] is the most well-known personality model, being applied in hiring processes, social dynamics, education and relationships [15, 16, 17]. The Myers-Briggs Type Indicator (MBTI) handbook illustrates a four factor model of personality where people form their 'personality type' by attaining one attribute from each of four dichotomies; Extrovert/Introvert, Intuitive/Sensory, Thinking/Feeling and Judging/Perceiving. This gives 16 different personality types where a letter from each dichotomy is taken to produce a four letter acronym, e.g., 'ENTJ' or 'ISFP'.

The model has received substantial scrutiny, particularly from psychologists who question its validity and reliability [18, 19]. Nonetheless, we utilise the Myers-Briggs model in our analysis for the following reasons:

- Thousands of Twitter users self-report their MBTI on Twitter. This enables us to obtain a labelled dataset through appropriately querying for each of the 4 letter personality type acronyms that are unique to MBTI.
- The Myers-Briggs model has the largest number of selfreports on Twitter, enabling us to achieve the largest labelled personality dataset on Twitter.
- We aim to develop a framework for modelling personality profiles from social media data using statistical machine learning (ML) approaches. MTBI is a test case for our framework, which can be applied to other personality models (or other labelings/characteristics of individuals on social media) more generally.

Open-source labelled training data with Myers-Briggs personality types has not existed until recently. Plank and Hovy [20] modeled the MBTI of Twitter users through attaining a small dataset of 1,500 users and Gjurković and Šnajder [21] modeled the MBTI on a larger corpus of Reddit users. In 2017, Jolly [22] posted a labelled MBTI dataset on Kaggle, constituting the only known publicly available labelled dataset used for modelling the MBTI of social media users. The dataset was comprised of 8,675 users, their personality types and a section of their last 50 posts on an online forum called personalitycafe.com. This small online forum contains 153,000 members dedicated to discussing health, behavior, personality types and personality testing. The discussions are therefore quite different to those on other social media platforms, and likely a different demographic. Hence, this dataset is likely not generalisable to other platforms like Twitter and Facebook. It is also relatively small and imbalanced, limiting which models can be utilised on various feature sets. Class imbalance is considerable in all cases, and in one particular dataset some classes up to 28 times larger than their counterpart. Nevertheless, many papers apply machine learning models to such datasets without accounting for these class imbalances [4, 23, 24, 25, 3]. Consequently, the metrics reported often misrepresent model performance, and instead highlight the severity of class imbalances in the datasets.

III. DATA COLLECTION & PREPROCESSING

We discovered a number of Twitter accounts to self-report their labelled MBTI acronym on Twitter in the form of a regular expression. We therefore formulated two methods for querying and labelling the Myers-Briggs personality type of accounts. Let Ω define the set of 16 acronyms for Myers-Briggs personality types, then:

- M1 Query: $\{x : x \in \Omega\}$. We obtained the set of users who currently self-report their personality type in their username or biography.
- M2 Query: $\{(I \text{ am } x) \lor (I \text{ am a } x) \lor (I \text{ am an } x) : x \in \Omega\}$. We obtained the set of users who have self-reported their personality type in a Tweet since Twitter's creation (March 26, 2006). Note that we only searched for self-reports in Tweets, not Retweets, Quotes and Replies due to a number of users often not self-reporting their own MBTI when referencing MBTI acronyms in these forms of communication.

Note that in both cases, the queries were not case-sensitive.

The resulting dataset comprised of 68,958 users with their labelled MBTI; the dataset and more details on its collection are provided in [14]. In total, we collected 15,986 accounts by querying usernames and biographies, and 52,972 accounts from querying tweets, with misclassification rates 1.9% and 3.4% based on random samples of 1,000 accounts from each.

Next we obtained account characteristics for each user, including their: biography, most recent 100 tweets/quotes, as well as a set of Social Metadata (SM) features. The user's biography and the 100 tweets/quotes were used to generate a set of linguistic features, whereas SM features (Table I) are directly used as numeric features in the models.

We removed duplicate users, then combined the biography and tweets into a combined text for every account. We then:

- 1) Normalised the text and calculated each account's dominant language.
- 2) Removed non-English language using the Compact Language Detect 2 (PyCLD2) library.
- 3) Calculated (language-dependent) Botometer scores².
- 4) Converted text to lowercase, removed URLs, email addresses, punctuation and numbers.
- 5) Tokenized using the Tweet Tokenizer from the Natural Language Toolkit (NLTK) [26].
- 6) Removed empty tokens and any instances of the 16 MBTI acronyms.

²Further discussion: https://rapidapi.com/OSoMe/api/botometer-pro/details

Next, we formulated an inclusion-exclusion criteria to determine whether a personality could be profiled from a Twitter account: we kept accounts with over 100 tweets/quotes, over 50% English language, Botometer CAP score less than 0.8, and strictly one MBTI type referenced.

We use the Botometer CAP score because we are interested in the overall bot likelihood and not the sub-category bot likelihoods. Unfortunately, there is no consistency in the literature on thresholds for binary bot classification. Rather, authors define their threshold based on a false positive rate in the context of their problem. For instance, Wojcik et al. [27] use a threshold of 0.43 for their political analysis of the twittersphere, whereas Keller and Klinger [28] use a larger threshold of 0.76 for their analysis of social bots in election campaigns. To avoid large numbers of false positive bot classifications, we chose a high threshold of 0.8.

Finally, we extracted the LIWC, BERT and VADER features from the text. The data cleaning techniques above were performed only for LIWC feature extraction, whereas the BERT and VADER features can be extracted directly from the raw text output. Thus, we calculated the LIWC features on the combined text by micro-averaging the tokens present in each LIWC category for every user. Next, we calculated the BERT features on the raw Twitter output using BERTweet [29], a pre-trained language model for English Tweets. First, we averaged the embeddings for the tokens to form a single embedding vector for each tweet/quote, then averaged the embedding vectors for the tweets/quotes to create a single 768dimensional embedding vector for each user. We calculated the VADER features (sentiment, proportion of positive words and proportion of negative words) on the raw Twitter output for each user and include scores for both a user's biography and their tweets. We distinguish these because of contextual differences in the language; biographies often discuss oneself and tweets often discuss one's environment. We then have a total of 866 features; these are provided in Table I.

Category	Features
SM	followers_count, friends_count, listed_count,
	favourites_count, geo_enabled, verified, statuses_count,
	default_profile, default_profile_image,
	profile_use_background_image, has_extended_profile
Botometer	cap_english, english_astroturf, english_fake_follower,
	english_financial, english_other, english_self_declared,
	english_spammer
LIWC	function, pronoun, ppron, i, we, you, shehe, they, ipron,
	article, prep, auxverb, adverb, conj, negate, verb, adj,
	compare, interrog, number, quant, affect, posemo, negemo,
	anx, anger, sad, social, family, friend, female, male,
	cogproc, insight, cause, discrep, tentat, certain, differ,
	percept, see, hear, feel, bio, body, health, sexual, ingest,
	drives, affiliation, achiev, power, reward, risk, focuspast,
	focuspresent, focusfuture, relativ, motion, space, time,
	work, leisure, home, money, relig, death, informal, swear,
	netspeak, assent, nonflu, filler, total_word_count
BERT	$\{e_i \; ; \; i=1,\ldots,768\}$
VADER	tweets_sentiment, bio_sentiment, tweets_pos_words,
	bio_pos_words, tweets_neg_words, bio_neg_words

TABLE I: Features in our models, separated by category.

IV. EXPLORATORY DATA ANALYSIS

We performed an exploratory data analysis (EDA) on the dataset to determine important information about our dataset, prior to any modelling. We acknowledge and discuss two forms of potential bias in our dataset: (i) only considering MBTI types on Twitter; (ii) only selecting accounts which satisfy our inclusion-exclusion criteria as well as self-report their MBTI types on Twitter. Figure 1 demonstrates these biases through bar plots showcasing the proportions of the MBTI dichotomies in our dataset. We compare with a study reporting MBTI proportions on Twitter [30], and with the proportion of personality types in the general population [31].

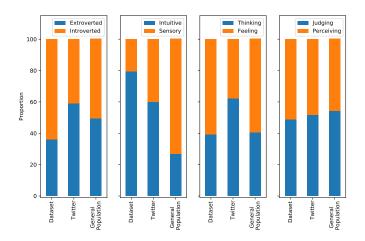


Fig. 1: Proportion of accounts displaying each dichotomous trait in our dataset, on Twitter and in the general population.

A noticeable imbalance in the Intuitive/Sensory dichotomy exists across all datasets in Figure 1. There are also observable imbalances in the Extrovert/Introvert and Thinking/Feeling dichotomies. Whereas, the Judging/Perceiving dichotomy is more balanced across each dataset than the other dichotomies. The imbalances in our dataset are mostly consistent with those from www.personalitycafe.com. The higher proportion of introverts in our dataset is consistent with [32] who find that introverts tend to use social media as a primary form of communication, whereas extroverts tend to prefer communicating in-person. The larger proportion of intuitives in our dataset is consistent with Schaubhut et al. [30] who discovered that more Intuitive individuals (13%) reported being active users of Twitter than individuals with a preference for Sensing (8%). The imbalance in the Thinking/Feeling dichotomy in our dataset is opposite to what we observe in the Twitter dataset. However, Schaubhut et al. [30] found that people displaying the Feeling trait are more likely to spend their personal time browsing, interacting and sharing information on Facebook. Provided the same is true for Twitter users, our inclusionexclusion condition requiring users to be active on Twitter (i.e. tweet/quote at least 100 times) may bias our dataset leading to more users exerting the Feelings trait.

Some authors don't assume independence between the dichotomies when modelling [23, 3], whereas most choose to model the dichotomies independently [33, 34, 35, 24, 25]. We take a data-driven approach, determining the dependency structure of the four MBTI dichotomies in our dataset using the bias-corrected version of the Cramér's V Statistic [36] (Table II). The Cramér's V statistic is small in every case, implying that the four Myers-Briggs dichotomies are independent in our dataset, and so we model them independently.

	E/I	N/S	T/F	J/P
E/I	1.00	0.03	0.00	0.10
N/S	0.03	1.00	0.02	0.08
T/F	0.00	0.02	1.00	0.11
J/P	0.10	0.08	0.11	1.00

TABLE II: Pairwise results of the bias-corrected Cramér's V Statistic between the MBTI dichotomies for our dataset.

We performed a Principal Component Analysis (PCA) on the features to discover if we could significantly reduce the dimension of the feature space, and multicollinearity between the features. The first principal component explains 25.1% of the variance in the data and the first 200 principal components explain 95.4% of the variance in the data. As a result, we utilise the first 200 PCA components in our machine learning models, significantly reducing both the dimension of the feature space and the multicollinearity of the features.

V. MODEL COMPARISON

We train LR, NB, SVM and RF classifiers on each of the four dichotomies in our dataset, using 10-fold cross validation. The class imbalances we observe for some dichotomies (particularly Intuitive/Sensory and Extrovert/Introvert), leads us to perform four different weighting/sampling techniques prior to model fitting:

- Weight the importance of classifying dichotomies,
- Upsample the minority class (with replacement),
- Perform the Synthetic Minority Oversampling Technique (SMOTE) on the minority class,
- Downsample the majority class.

Each model uses the first 200 principal components of the features in Table I as predictors. As an example, Figure 2 shows confusion matrices for the Intuitive/Sensory dichotomy under the standard LR model and the upsampled LR model.



(a) Standard logistic regression (b) Upsampled logistic regression

Fig. 2: Confusion matrices for modelling the N/S dichotomy.

This shows that the standard LR model primarily predicts the majority class, indicating that it exploits the class imbalance to make predictions on the test sets. In comparison, the upsampled model predicts significantly more of the minority class on the test sets, resulting in more accurate predictions for the minority class. We observe similar behavior for all other models, highlighting the importance of weighting/sampling techniques to ameliorate the effect of class imbalance for prediction. However, we observe a clear trade-off between accurately predicting the majority and minority classes, with an overall reduction in accuracy due to weighting/sampling techniques. We therefore report both accuracy and Area Under the Curve (AUC) metrics for each of our models in Table III. We report four types of accuracy depending on the number of accurately predicted dichotomies in each model. Of course, accuracy can be a misleading metric when assessing a model's performance on unbalanced data, so for comparison we report the accuracies for a random classifier and a majority class classifier. Moreover, we use an approach similar to other authors to report two types of AUC for each model [37, 38]: we macro-average and micro-average the true positive rate and false positive rate at each threshold of the ROC curve for the independent models of each dichotomy. This provides us with two ROC curves (and AUC metrics) for each model. The micro-averaged AUC aggregates the contributions of all samples in each model and weights individual predictions equally, so it is generally less sensitive to class imbalances. Table III compares the accuracies and AUCs of the best performing models from each method. In each case, we include the 'Standard' model and the weighted/sampling model which achieves the highest sum of micro- and macro-averaged AUC.

	Accura	Accurately Predicted Dichotomies			AUCs	
Model	4	≥ 3	≥ 2	≥ 1	Macro	Micro
Standard LR	20.82	60.43	89.35	98.82	0.6688	0.6547
SMOTE LR	13.89	48.63	82.51	97.65	0.6642	0.6620
Standard NB	14.20	49.17	81.91	97.40	0.5784	0.5867
Upsampled NB	13.75	48.06	80.82	97.18	0.5861	0.5917
Standard SVM	20.95	60.25	89.64	98.90	0.6693	0.6518
SMOTE SVM	13.56	48.61	82.54	97.61	0.6660	0.6554
Standard RF	19.69	57.96	88.69	98.67	0.6223	0.6273
Upsampled RF	19.70	58.16	88.48	98.76	0.6305	0.6264
Random Classifier	6.250	31.25	68.75	93.75	0.5000	0.5000
Majority Class	15.31	54.54	87.20	98.28	0.5000	0.5000

TABLE III: Accuracies and AUCs for the best performing models for each ML method. We include results from the 'Standard' model (with no weighting/sampling) and the best performing weighted/sampling model. Note that we determine the 'best performing weighted/sampling model' based on the sum of macro- and micro-averaged AUC.

Table III highlights the relatively small improvement in accuracy achieved by each model in comparison to the majority class classifier. It is clear that our standard SVM model is the best performing model on average. However, this model is only 5.64% more accurate at predicting a user's complete personality type compared to the majority class classifier. This is a reasonable and statistically significant improvement, but we remark based on the above discussion that the standard models are simply exploiting the class imbalances in our

dataset. Moreover, we achieve accuracies very similar to those obtained by Plank and Hovy [20], who produced the only other Twitter dataset of labelled MBTI's (to our knowledge). In particular, we achieve better accuracies for the T/F and J/P dichotomies, and only marginally worse accuracies for the E/I and N/S dichotomies – further evidencing that our models perform similarly to others in the literature.

Interestingly, the standard LR model most accurately predicts at least three out of four user dichotomies and is only marginally worse than the standard SVM model for all other metrics. The LR model is also significantly faster to train than the SVMs – making it the model of choice on larger datasets.

The AUC is important in discussions of model performance, especially for unbalanced datasets. This is because it equally weights the true positive rate and false positive rate, making it more robust for unbalanced datasets compared to accuracy. Most of our AUCs lie around 0.65 (65%), apart from the NB Classifiers. In particular, the best performance for the macroaveraged and micro-averaged AUCs is the standard SVM model and the SMOTE LR model, respectively. These AUCs are significantly larger than what we observe for both the random classifier and the majority class classifier, indicating there is certainly some 'signal' in our features. We therefore perform an in-depth analysis of feature importance next.

VI. FEATURE IMPORTANCE

We perform independent upsampled LR models on each of the four MBTI dichotomies because they performed well on our dataset (macro- and micro-averaged AUCs: 0.6676 and 0.6536). We choose an LR model because it is fast to train and, straightforward to interpret and perform feature selection on. Moreover, we use an upsampled model because it does not involve creating 'synthetic' data in the same way that SMOTE does – this is important for determining feature importance.

We consider the variable importance of the descriptive features in our models; these include all features except from BERT. For each dichotomy we fit the upsampled LR model and perform a stepwise feature selection to obtain a model with only significant features. In each case, we start with a null model and perform the stepwise selection algorithm on the p-values with a threshold in of 0.05 and a threshold out of 0.1. We determine the variable importance of features using the t-statistic for the parameter coefficients associated with each feature. For each dichotomy, we calculate the variable importance of each remaining feature after the stepwise selection algorithm is complete and display the absolute value of the variable importance. Figure 3 displays the 12 most important features for each model in a bar chart. We colour the bars based on the variable's preference for each class in the dichotomy.

Pennebaker and Francis [39] suggested function words such as pronouns (pronoun), personal pronouns (ppron), 1st person singular (i), 1st person plural (we), prepositions (prep), auxiliary verbs (auxverb) and negations (negate), can describe people. Figure 3 shows the function words that are significant predictors in our models, e.g., 1st person plurals are significant in the E/I model and prepositions are significant in the N/S

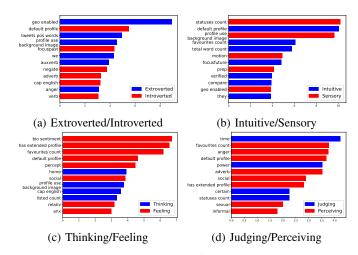


Fig. 3: Variable Importance Plots for an upsampled LR model for each dichotomy. Variables are sorted by the absolute value of the variable importance (from top to bottom). We colour bars by the feature preference for each class.

model. This reinforces the importance of function words, and that techniques such as stop-word removal may remove useful information, particularly for tasks like personality prediction.

Extroverts tend to be associated with more positive language, and introverts tend to have more of a focus on the past. Similarly, Chen et al. [40] suggested that extroverts display more positive emotion because they have a "dispositional tendency to experience positive emotions". Accounts which have a larger favourites count (i.e. the account likes more tweets) tend to be more intuitive, whereas accounts which write more statuses tend to be more sensory. Interpreting favourites as a proxy for the amount of information an account consumes, our results suggest that intuitives consume more information on Twitter, whereas sensory individuals write more. This proxy is of course not perfect, because people may consume information without liking it. Nonetheless, it is consistent with Myers-Briggs Foundation definitions, which state that intuitives pay "most attention to impressions or the meaning and patterns of the information", whereas sensors pay "attention to physical reality, what I see, hear, touch, taste, and smell" [41]. The strongest predictor for the J/P dichotomy (Figure 3d) is time; judgers are more likely to use words related to time and certainty compared to perceivers. 'End', 'until' and 'season' are examples of time-related words and 'always', 'never' are words related to certainty. This is also consistent with the Myers-Briggs Foundation, which states judgers "prefer a planned or orderly way of life, like to have things settled and organized" [41].

Next we explore how emoji usage relates to a Twitter user's MBTI. On Twitter, emojis often have multiple meanings. For instance, the rainbow flag can indicate support for LGBTQ+ social movements, the wave can symbolise a "Resister" crowd of anti-Trump Twitter, and the okay symbol can be used by white supremacists, some of which covertly use the symbol to indicate their support for white nationalism [42]. Hence,

emojis can indicate how these groups/movements interact with different personality types. We determine each emoji's frequency in a user's tweets and include these frequencies as predictors in upsampled LR models. Performing the same stepwise feature selection algorithm as above, we display the 12 most important predictors from the remaining models in Figure 4.

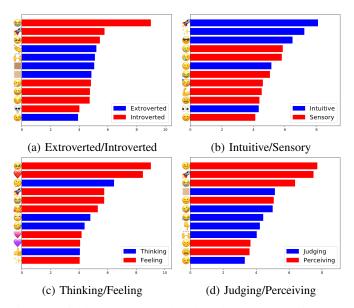


Fig. 4: Variable Importance Plots based on only emoji counts in the upsampled LR models. Variables are sorted by absolute value of the variable importance (from top to bottom). We colour bars by the feature preference for each class.

The rocket ship emoji is one the top 12 most important predictors across all models. An increase in this emoji's usage implies a higher likelihood of an account being introverted, intuitive, feelings-orientated and perceiving. The rocket ship emoji has been used by finance enthusiasts who use the emoji to denote a fast increase in a particular stock or cryptocurrency. Hence, it is possible that we are observing crypto enthusiasts to be more introverted, intuitive, feelings-orientated and perceiving. However, this emoji has other meanings like as an actual rocket ship, so we explore created word clouds of tweets containing the rocket emoji (Figure 5a), as well as the red heart emoji (Figure 5b). The rocket ship generally appears in crypto-related tweets discussing 'projects', 'great opportunities', 'developments' and 'cryptos'. However, it also appears in tweets discussing the 'moon' and 'space'. The red heart emoji mainly appears in emotive tweets discussing 'love' and 'happiness'. A number of the emojis making an account more introverted are sad/upset emojis, whereas no sad/upset emojis make an account more extroverted. This further confirms Figure 3a which suggested that extroverts prefer to display positive emotion online.

Next we wanted to consider the importance of different feature groups (including the BERT features) and discuss whether different groups of features are more informative in our models. Again, we do this by fitting an upsampled logistic





(a) Rocket Ship Emoji

(b) Red Heart Emoji

Fig. 5: Word clouds of tweets/quotes containing specific emojis in our dataset. Larger words appeared more frequently in the tweets - we present results for the rocket ship emoji (left) and the red heart emoji (right). Note that we remove stopwords as they do not provide much context for the tweets.

regression model to all of the features and performing the stepwise feature selection on each of the models. We use the same thresholds for accepting and removing features. We then use the number of remaining features in each feature group after stepwise selection to measure the importance of the different feature groups. For each model, Table IV displays both the number of predictors (in each feature group) and their proportion that remain after the stepwise feature selection algorithm. This proportion can be considered a measure of the importance of each feature group which is not biased by the number of features in each group. We introduce a more robust statistical framework to determine whether different groups of features are actually more informative of our data. We do this by performing a Chi-Squared Test on the number of features retained and excluded from each model. We test the null hypothesis that each feature group is equally as informative (per feature) and include the p-values from the Chi-Square Test in the sub-table captions displayed in Table IV.

The number of features selected depends on the type of model. For instance, 243 features are selected in the N/S model, whereas only 124 features are selected in the J/P model. Interestingly, the N/S model is also the most accurate and the J/P model the least accurate, implying a positive relationship between accuracy and number of features retained. This is consistent with the remark that more features are retained in a model when they are more informative about the data. Moreover, the SM features are on average the most-retained across models. Conversely, the Botometer features have worst payoff across the four models, having the smallest proportion retained on average. The most interesting comparison is between the LIWC and BERT features, which both aim to describe linguistic properties about users. In each model, the BERT features are more highly retained. However, only the results from the N/S model and the T/F model are significant at the 5% level. We therefore reject the null hypothesis that each feature group is equally as informative (per feature) for the N/S and T/F models. However, the Chi-Squared Test does not alone tell us what feature groups perform significantly better,

Feature	I	Prop.	Feature	1	Prop.
Type	#	Retained	Type	#	Retained
SM	4	36.4%	SM	7	63.6%
LIWC	15	20.3%	LIWC	18	24.3%
BERT	176	22.9%	BERT	217	28.3%
Botometer	1	14.3%	Botometer	0	0.00%
VADER	2	33.3%	VADER	1	16.7%
Total	198	22.9%	Total	243	28.1%
(a) E/I Feature	l .	0.720) Prop.	(b) N/S Feature	l (I	0.032) Prop.
Туре	#	Retained	Type	#	Retained
SM	5	45.5%	SM	4	36.4%
LIWC	11	14.9%	LIWC	8	10.8%
BERT	124	16.1%	BERT	112	14.6%
Botometer	1	14.3%	Botometer	0	0.00%
VADER	3	50.0%	VADER	0	0.00%
Total	144	16.6%	Total	124	14.3%

TABLE IV: Number of features and proportion of features retained in each group after stepwise feature selection. For each dichotomy, we perform Chi-Squared Tests on the null hypothesis that each feature group is equally informative, per feature (see *p*-values).

so we perform individual confidence intervals (CIs) for the binomial proportions of accepting/rejecting features in each group using the Wilson Score interval [43]. The CIs for each feature group and model are displayed in Figure 6.

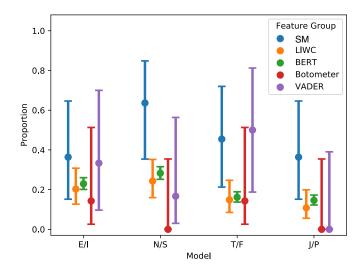


Fig. 6: 95% Wilson Score Binomial confidence intervals (CIs) for the proportion of retained features in each feature group. We display the CIs for each model and use the Wilson Score version to correct for having zero successes in some cases.

For the Intuitive/Sensory model, the 95% CI for the SM features lies completely above the 95% CIs for the LIWC and BERT features. This indicates that the SM features are more informative (per feature) than the LIWC and BERT features at the 5% level for this dichotomy. This highlights that attributes about a user's account are sometimes more important than the

language they use when modelling personality. This statement is also validated by the results for the Thinking/Feeling model, where the 95% CI for the SM features and VADER features lie completely above the 95% CI for the BERT features. We deduce that we are likely observing these results because the textual features are all fairly correlated with each other. Moreover, there is no evidence to suggest that the BERT features are any more informative than the LIWC features in determining someone's Myers-Briggs personality type.

VII. CONCLUSION

This paper contributes a labelled dataset of personality types from Twitter and a framework to model the personality types of these users. To our knowledge, this is the largest available Twitter dataset of labelled Myers-Briggs Personality Types. The only comparable dataset [20] contains only 1,500 labelled accounts. Moreover, the data collection techniques we used to collect this data are also novel, as they avoid the long, cumbersome questionnaires used in other research. Moreover, we develop a statistical framework which combines NLP tools and mathematical models to model/predict the personality type of users online. While we applied this framework to personality types, it can model any labelled characteristics of online accounts - political opinions, psychological properties or even someone's propensity to adopt an opinion or viewpoint. As part of this framework, we analyse and compare a number of different machine learning models. Since personality types in our dataset are unbalanced, we compare different weighting/sampling techniques to deal with issues arising from class imbalance. We discover that class imbalances are common in these types of datasets and are something which is often overlooked by many scholars. Because of this, we demonstrate why models on these datasets appear more accurate than they really are and we demonstrate why your digital footprint may be less informative of your personality type than you may think. Finally, we compare the importance of different features in our models on an individual and group level.

As we use a large number of features from a large dataset, a deep learning model would be applicable. However, this would give less interpretability than the models used here. It would also be interesting to consider different data collection methods. One limitation of our dataset is that we only have access to the classification of the four personality dimensions, when in reality these dimensions are represented on a numerical scale. For instance, two users may be extroverted but one user may be considerably more extroverted than the other. While performing questionnaires are long and expensive, it would enable us to obtain these personality dimensions on a numerical scale. We would expect this to have a significant improvement on the performance of our models.

ACKNOWLEDGMENT

LM acknowledges support from the Australian Government through the Australian Research Council's Discovery Projects funding scheme (project DP210103700).

REFERENCES

- S. Dixon, "Number of Worldwide Social Network Users 2027," https://www.statista.com/statistics/278414/number-of-worldwide-social-network-users/, Sep. 2022.
- [2] C. Wylie, Mindf*ck: Inside Cambridge Analytica's Plot to Break the World, main edition ed. London: Profile Trade, Aug. 2020.
- [3] S. M. Patil, R. Singh, P. Patil, and N. Pathare, "Personality Prediction Using Digital Footprints," in 2021 5th International Conference on Intelligent Computing and Control Systems (ICICCS), May 2021, pp. 1736–1742.
- [4] M. M. Tadesse, H. Lin, B. Xu, and L. Yang, "Personality Predictions Based on User Behavior on the Facebook Social Media Platform," *IEEE Access*, vol. 6, pp. 61959–61969, 2018.
- [5] T. Tandera, Hendro, D. Suhartono, R. Wongso, and Y. L. Prasetio, "Personality Prediction System from Facebook Users," *Procedia Computer Science*, vol. 116, pp. 604–611, Jan. 2017.
- [6] D. Weber, M. Nasim, L. Falzon, and L. Mitchell, "# arsonemergency and australia's "black summer": Polarisation and misinformation on social media," in *Disinformation in Open Online Media: Second Mul*tidisciplinary International Symposium, MISDOOM 2020, Leiden, The Netherlands, October 26–27, 2020, Proceedings 2. Springer, 2020, pp. 159–173.
- [7] J. Tuke, A. Nguyen, M. Nasim, D. Mellor, A. Wickramasinghe, N. Bean, and L. Mitchell, "Pachinko prediction: A bayesian method for event prediction from social media data," *Information Processing & Management*, vol. 57, no. 2, p. 102147, 2020.
- [8] M. Block, "How the Myers-Briggs Personality Test Began in a Mother's Living Room Lab," NPR, Sep. 2018.
- [9] C. G. Jung, Collected Works of C.G. Jung, Volume 6: Psychological Types, 1st ed. Princeton: Princeton University Press, Oct. 1976.
- [10] J. Pennebaker, R. Boyd, K. Jordan, and K. Blackburn, The Development and Psychometric Properties of LIWC2015, Sep. 2015.
- [11] C. Hutto and E. Gilbert, VADER: A Parsimonious Rule-based Model for Sentiment Analysis of Social Media Text. The AAAI Press, Jan. 2015.
- [12] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "BERT: Pre-Training of Deep Bidirectional Transformers for Language Understanding," May 2019.
- [13] M. Sayyadiharikandeh, O. Varol, K.-C. Yang, A. Flammini, and F. Menczer, "Detection of Novel Social Bots by Ensembles of Specialized Classifiers," in *Proceedings of the 29th ACM International Conference on Information & Knowledge Management*, Oct. 2020, pp. 2725–2732.
- [14] J. Watt, "Self-Reported Myers-Briggs Personality Types on Twitter,"
 7 2023. [Online]. Available: https://figshare.com/articles/dataset/ Self-Reported_Myers-Briggs_Personality_Types_on_Twitter/23620554
- [15] R. E. De Vries, "The Main Dimensions of Sport Personality Traits: A Lexical Approach," Frontiers in Psychology, vol. 11, 2020.
- [16] B. W. Walsh and J. L. Holland, "A Theory of Personality Types and Work Environments," in *Person–Environment Psychology: Models and Perspectives*. Hillsdale, NJ, US: Lawrence Erlbaum Associates, Inc, 1992, pp. 35–69.
- [17] H. W. Lane, M. L. Maznevski, M. E. Mendenhall, and J. McNett, The Blackwell Handbook of Global Management: A Guide to Managing Complexity. John Wiley & Sons, Feb. 2009.
- [18] D. Pittenger, "Measuring the MBTI ... and Coming up Short," Journal of Career Planning and Employment, vol. 54, Jan. 1993.
- [19] A. Grant, "Goodbye to MBTI, the Fad That Won't Die Psychology Today," https://www.psychologytoday.com/intl/blog/give-and-take/201309/goodbye-mbti-the-fad-won-t-die.
- [20] B. Plank and D. Hovy, "Personality Traits on Twitter—or—How to Get 1,500 Personality Tests in a Week," in *Proceedings of the 6th Workshop on Computational Approaches to Subjectivity, Sentiment and Social Media Analysis*. Lisboa, Portugal: Association for Computational Linguistics, Sep. 2015, pp. 92–98.
- [21] M. Gjurković and J. Šnajder, "Reddit: A Gold Mine for Personality Prediction," in *Proceedings of the Second Workshop on Computational Modeling of People's Opinions, Personality, and Emotions in Social Media*. New Orleans, Louisiana, USA: Association for Computational Linguistics, Jun. 2018, pp. 87–97.
- [22] M. Jolly, "(MBTI) Myers-Briggs Personality Type Dataset," https://www.kaggle.com/datasets/datasnaek/mbti-type.
- [23] S. Başaran and O. H. Ejimogu, "A Neural Network Approach for Predicting Personality from Facebook Data," SAGE Open, vol. 11, no. 3, p. 21582440211032156, Jul. 2021.

- [24] S. S. Keh and I.-T. Cheng, "Myers-Briggs Personality Classification and Personality-Specific Language Generation Using Pre-Trained Language Models," Jul. 2019.
- [25] M. H. Amirhosseini and H. Kazemian, "Machine Learning Approach to Personality Type Prediction Based on the Myers–Briggs Type Indicator®," *Multimodal Technologies and Interaction*, vol. 4, no. 1, p. 9, Mar. 2020.
- [26] S. Bird, E. Klein, and E. Loper, Natural Language Processing with Python: Analyzing Text with the Natural Language Toolkit, 1st ed. Beijing; Cambridge Mass.: O'Reilly Media, Aug. 2009.
- [27] S. Wojcik, S. Messing, A. Smith, L. Rainie, and P. Hitlin, "Bots in the Twittersphere," Apr. 2018.
- [28] T. R. Keller and U. Klinger, "Social Bots in Election Campaigns: Theoretical, Empirical, and Methodological Implications," *Political Communication*, vol. 36, no. 1, pp. 171–189, Jan. 2019.
- [29] D. Q. Nguyen, T. Vu, and A. Tuan Nguyen, "BERTweet: A Pre-Trained Language Model for English Tweets," in *Proceedings of the 2020 Con*ference on Empirical Methods in Natural Language Processing: System Demonstrations. Online: Association for Computational Linguistics, Oct. 2020, pp. 9–14.
- [30] N. Schaubhut, A. Weber, and R. Thompson, "Myers-Briggs Type and Social Media Report," themyersbriggs.com/contents/MBTI_and_Social_Media_Report.aspx, 2012.
- [31] M. Robinson, "How Rare Is Your Personality Type?" https://www.careerplanner.com/MB2/TypeInPopulation.cfm, 1998.
- [32] Knowledge Leader, "How Technology and Social Media Empower the Introvert," https://knowledge-leader.colliers.com/editor/how-technologyand-social-media-empower-the-introvert/, Sep. 2015.
- [33] F. Alam, E. A. Stepanov, and G. Riccardi, "Personality Traits Recognition on Social Network Facebook," Proceedings of the International AAAI Conference on Web and Social Media, vol. 7, no. 2, pp. 6–9, 2013.
- [34] D. Sumpter, Outnumbered: From Facebook and Google to Fake News and Filter-Bubbles – the Algorithms That Control Our Lives, illustrated edition ed. London: Bloomsbury Sigma, Jun. 2018.
- [35] S. Bharadwaj, S. Sridhar, R. Choudhary, and R. Srinath, "Persona Traits Identification Based on Myers-Briggs Type Indicator (MBTI) - a Text Classification Approach," in 2018 International Conference on Advances in Computing, Communications and Informatics (ICACCI), Sep. 2018, pp. 1076–1082.
- [36] H. Cramér, "Mathematical Methods of Statistics," in *Mathematical Methods of Statistics*, ser. Princeton Mathematical Series; 9. Princeton: Princeton University Press, 1946.
- [37] G. Gunčar, M. Kukar, M. Notar, M. Brvar, P. Černelč, M. Notar, and M. Notar, "An Application of Machine Learning to Haematological Diagnosis," *Scientific Reports*, vol. 8, no. 1, p. 411, Dec. 2018.
- [38] N. C. De, L. Cindolo, L. Sarchi, A. Iseppi, M. Rizzo, B. Riccardo, A. Minervini, F. Sessa, G. Muto, P. Bove, M. Vittori, G. Bozzini, P. Castellan, F. Mugavero, D. Panfilo, S. Saccani, M. Falsaperla, L. Schips, A. Celia, M. Bada, A. Porreca, A. Pastore, A. S. Yazan, G. Marco, G. Novella, R. Rizzetto, N. Trabacchin, M. Guglielmo, G. Pini, R. Lombardo, B. Rocco, A. Antonelli, and A. Tubaro, "Using a Machine Learning Algorithm to Predict Prostate Cancer Grade," *Journal of Urology*, vol. 203, no. Supplement 4, pp. e1236–e1236, Apr. 2020.
- [39] J. W. Pennebaker and M. E. Francis, "Cognitive, Emotional, and Language Processes in Disclosure," *Cognition and Emotion*, vol. 10, no. 6, pp. 601–626, Nov. 1996.
- [40] J. Chen, L. Qiu, and M.-H. R. Ho, "A Meta-Analysis of Linguistic Markers of Extraversion: Positive Emotion and Social Process Words," *Journal of Research in Personality*, vol. 89, p. 104035, Dec. 2020.
- [41] "The Myers & Briggs Foundation Take the MBTI® Instrument," https://www.myersbriggs.org/my-mbti-personality-type/take-the-mbti-instrument/, 2022.
- [42] C. Bronsdon, "What Do Different Twitter Emojis Mean?" https://conorbronsdon.com/blog/what-do-different-twitter-emojis-mean.
- [43] J. Reed, "Better Binomial Confidence Intervals," Journal of Modern Applied Statistical Methods, vol. 6, no. 1, May 2007.