

Machine Prediction of Personality from Facebook Profiles

Randall Wald and Taghi Khoshgoftaar
Florida Atlantic University
{rwald1, khoshgof}@fau.edu

Chris Sumner
Online Privacy Foundation
chris@onlineprivacyfoundation.org

Abstract—An increasing number of Americans use social networking sites such as Facebook, but few fully appreciate the amount of information they share with the world as a result. Although studies exist on the sharing of specific types of information (photos, posts, etc.), one area that has been less explored is how Facebook profiles can share personality information in a broad, machine-readable fashion. In this study, we apply data-mining and machine learning techniques to predict users' personality traits (specifically, the traits of the Big Five personality model) using only demographic and text-based attributes extracted from their profiles. We then use these predictions to rank individuals in terms of the five traits, predicting which users will appear in the top or bottom 5% or 10% of these traits. Our results show that when using certain models, we can find the top 10% most Open individuals with nearly 75% accuracy, and across all traits and directions, we can predict the top 10% with at least 34.5% accuracy (exceeding 21.8%, which is the best accuracy when using just the best-performing profile attribute). These results have privacy implications in terms of allowing advertisers and other groups to focus on a specific subset of individuals based on their personality traits.

Index Terms—Facebook, Big Five, privacy, data mining, personality prediction

I. INTRODUCTION

In today's interconnected society, many people employ social networking sites such as Facebook and Twitter. These mediums of communication permit users to post information relevant to their lives and stay updated on their friends and favorite celebrities. The motivations for social networking differ from user to user: some focus on broadcasting information about themselves, while others are more interested in passively consuming information produced by others. Nonetheless, all but the most networking-averse do post information which can be used to gain insight into their personality and behavior, information which is extremely relevant to marketers.

It is well known that social networking sites use information from their users to target advertisements; a popular idiom states that if the user does not pay for the product, they *are* the product. Companies have towed the line in how to share this information, sometimes drawing the ire of users. In one well-known case, Facebook drew criticism over its Beacon project over concerns that it would publish users' activities on non-Facebook pages to their profiles. Other companies are not immune to privacy concerns either: Twitter account hacks have led to embarrassing tweets from high-profile users, and Google Buzz was sued over allegations that it leaked users' contact lists to their friends. Users of social networking sites are becoming more aware of the information they post, and more concerned about how this information can be used to identify them.

Most concerns pertaining to social networking posts focus on raw demographic information or specific damaging posts. For example, religious or political affiliation may impact on how an employer views a potential employee, as might an inflammatory post or salacious photograph. However, relatively little attention has focused on aggregating information from posts and profile text to create a broad picture of the user. The Big 5 Experiment [20], an experiment performed by the Online Privacy Foundation, has endeavored to do

just this. Users were asked to take a short personality survey, to determine their archetype according to the Big 5 personality model (a five-factor model which gives individuals a score for five different aspects of personality, specifically Agreeableness, Conscientiousness, Extroversion, Openness, and Neuroticism), and were also asked to provide information from their Facebook profiles. Using raw statistical tools, some useful patterns (such as a connection between number of friends and extroversion) were discovered.

This paper expands upon the Big 5 Experiment by employing a number of techniques from the field of data mining. Although data mining has been used in many application domains, its use for extrapolating information from social networking data for research purposes is relatively limited. Through these techniques, we have found that when ordering individuals based on our models, we could predict the top or bottom 10% based on any attribute and reliably categorize 34.5% of those individuals correctly (i.e., they were in fact in the top or bottom 10% 35% of the time). Furthermore, certain choices of attribute and direction (top or bottom) led to even better results: we could correctly select the top 10% in terms of Openness nearly 75% of the time, the bottom 5% in terms of Extroversion (e.g., the 5% most introverted) 64% of the time, and the top 5% most Agreeable 57% of the time. These results are significant because they would allow marketers and other interested parties to focus on specific subsets of users based on their profile information and create advertising more closely tailored to those users. They also have implications for more nefarious uses of social networks, where attackers seeking to perform social engineering attacks could determine which subset of the population is most susceptible.

The remainder of this paper will be organized as follows: Section II will present pertinent work in the study of social networking and its privacy implications. Section III will discuss the methods employed in this paper. Section IV presents the details of the Big 5 Experiment, including the data collection and preprocessing techniques. Section V presents our results, discussing the models and their predictions. Finally, section VI concludes and presents directions for future research.

II. RELATED WORK

Social networking sites, and Facebook in particular, have become a major part of our everyday culture. According to one study in 2011 [18], 51% of Americans have a Facebook profile. Unsurprisingly, many studies have examined the motivations and goals behind the increasing popularity of social networking. Some have focused on the use of Facebook to improve members' social capital (resources gained through use of peer networks) [4], [21], while others analyze what specific uses individuals gain from the site [9], [15], [17]. Further studies consider the changes in use over time [14] or how individuals find new people to add on Facebook [13]. Overall, there is a great deal of interest in understanding how and why users interact with Facebook.

One particularly important subset of research focuses on the perceptions and realities of privacy on sites such as Facebook. These studies examine such questions as how user demographics and knowledge pertain to online behavior [1], [7], how users cope with privacy challenges [3], and what specific threats to privacy most concern users [10]. The general conclusions have been that although users are becoming increasingly aware of the importance of privacy and security measures, they still often fail to implement appropriate measures to ensure their own privacy, implicitly assuming that breaches only concern others. These behaviors cut across all demographic lines, with race and gender not significantly impacting the use of privacy features.

While it is well-understood that Facebook profiles contain certain types of overt private information (e.g., contact information, potentially embarrassing photographs or posts, etc.), less recognized is the overall personality picture which can be understood from a Facebook profile. One popular tool for analyzing personality traits is the Big-Five Personality Index [5]. This defines human personality along five axes: Agreeableness, Conscientiousness, Extroversion, Neuroticism, and Openness. The Big Five metric has gained popularity in recent years, and some studies have examined how it can be used to better understand behavior on Facebook. One study by Gosling et al. [6] compared how individuals' Big Five values compared when these were ascertained by self-assessment, peer assessment, and independent reading of their Facebook profiles (by human volunteers); they found that Facebook profiles gave consistent values for the Big Five traits, and that aside from Neuroticism, these values were similar to the peer assessments given by users' friends. Another study, Back et al. [2], further examined the question of whether online profiles represent idealized or actual personality traits. As with Gosling et al., this included self-assessment, peer assessment, and independent profile reviewers. The researchers found that online profiles most closely represent real, not ideal, personality traits, especially in terms of Extroversion and Openness. A third study, Ross et al. [19], looks at the Big Five traits in a different fashion: by considering how they predict common Facebook uses. Although Extroversion and Openness had some correlation with some factors (for example, Extroversion was associated with membership in more Facebook groups), overall the personality factors did not have as much influence as the authors expected.

It is important to note that in the above studies, although the Big Five personality model was used to explore Facebook profiles, there was no effort to automatically (through software-based methods) interpret personalities based on Facebook profiles. Instead, the Big Five traits were used as intermediaries to compare how individuals are viewed in different contexts, or to understand how different types of personalities led to different online behavior. The present study is the first work to use machine learning techniques to process Facebook profiles and build a picture of users' personalities without direct human input.

III. METHODS

In many application domains, too much data is generated to be easily processed by humans. Although patterns exist, they require connecting vast quantities of data together in ways which would take years to be discovered manually. The field of data mining encompasses a wide range of techniques for extracting information from large datasets [22]. Many of these techniques are predicated on the idea that a dataset is a collection of instances, each representing one example of the type of data being analyzed. For studying social networks, as in this paper, an instance is a person. Each instance may

be considered as a set of (attribute, value) pairs; all instances will have the same number of pairs representing the same attributes, but their values will differ (and may even be "null" for some instances in some datasets). The specific goals of data mining depend on the nature of these attributes. Oftentimes, one or more of these attributes will be the "class," or dependent attributes; these are the most important attributes, and are used to understand each instance. For this study, the five classes (dependent attributes) are five different aspects of the Big 5 personality model, which will be discussed in greater detail in Section IV. The remaining attributes are referred to as features or independent attributes. These are pieces of information collected from each instance, not necessarily useful on their own but valuable for what they tell about the class attributes. Again, the specifics of what the features (independent attributes) of our study signify will be discussed later in the paper.

Depending on the nature (or presence) of the class attributes, different forms of data mining may be applied to a given dataset. If the class attribute is nominal (taking one of a fixed set of possible non-numeric values, called "classes"), classification may be performed. This is the task of building a model to assign each instance to one of the possible class values. This task is particularly useful for providing class labels to future, unlabeled instances. On the other hand, if the class attribute is numeric (taking a value from some integer or decimal range), the goal is numeric prediction: using the independent attributes to determine the numeric value of the class attribute. Finally, some datasets lack a class attribute altogether; for this data, the appropriate data mining approach is clustering — putting the instances into logical groups based solely on their independent attributes. Although clustering cannot be directly used to impart information about the instances, domain knowledge can sometimes use clustering as a stepping-stone to begin assigning class attributes. In this particular study, as the dependent attributes are numeric values ranging from 1 to 5, we employ numeric prediction models.

Numeric prediction models, as the name implies, concern themselves with predicting numeric values for a given attribute. They do this using the features of each instance along with a model built using information from the entire dataset. A number of different numeric prediction models have been employed throughout the literature, including linear regression, REPTree, and decision tables. These three models were chosen for use in this study because preliminary analysis showed they had the greatest performance. All models built in this study were implemented using the WEKA machine learning toolkit [8], using default values for all parameters except as noted below.

Linear regression (LinR) is a simple and widely-used form of numeric prediction. In this approach, each instance has its independent attribute values multiplied by a chosen constant (typically different for each attribute), and the results are summed together (along with a final weighting constant) to produce a predicted value for the class attribute. Although this model is simple, it presents the important parts of numeric prediction: to use the model, first a training phase must be employed (here, to generate the appropriate constants for each attribute), and then the instances may be run through the model to predict the class value.

Another popular form of numeric prediction is the decision tree-based learner known as REPTree. As with many decision trees, REPTree starts with a root node and picks one attribute to divide the instances into sub-trees. Then for each subtree, a new attribute (possibly the same as a previous attribute, only with a different threshold to divide the sub-trees) is used to further subdivide the instances, creating additional branches. Finally, leaf nodes are created,

and at each node the tree will predict a chosen value for the class attribute. To use the tree for numeric prediction, an instance starts at the top of the tree (root node) and is tested at each branch point: it will continue down whichever path contains the values appropriate for that instance. When it terminates this process at the leaf node, the value there is the predicted value of the class attribute for that instance. For this study, a total of seven folds were used for building and pruning the tree.

Decision tables (DTable) [12], despite their name, do not share much in common with decision trees. The table in question contains a list of instances and features. Although most implementations will include all training instances in the table, only a small fraction of the features are selected to build the model. Since only a few of the features are included, many instances will form into groups with identical independent attributes; in the case where some of the features are nominal (rather than numeric), discretization is employed to create small clusters of identical instances. The mean value of the class attribute is found for each group and becomes the prediction for all instances which fall into that group. When attempting to predict the value for a test instance, the features used by the decision table are considered, allowing the instance to be placed into one of these groups; if its values for these features are different from all other instances in the decision table (meaning it falls into none of the groups), then the mean value of the entire set is used to predict the result.

Although the models used in this study are based on numeric prediction, an additional processing step was performed using the predicted values for the class attributes. Because the goal is to understand which individuals are at the extreme ends (top or bottom) of the population in terms of the class attributes, the instances are sorted according to the predicted values for the class attribute. This process is similar to Module Order Modeling [11], which also evaluates the quality of a numeric prediction based on how it orders the instances, but differs in how the quality of the ordering is measured numerically. In the present work, this quality is found by considering, as the value of k changes, how many of the individuals predicted as being in the top/bottom k of the population (based on the class values generated by the model) are actually in the given k (based on the actual class values found in the data). For example, if $k = 50$, this test means determining how many of the individuals predicted as being in the top 50 are actually found in the top 50. Because there are a total of 537 individuals overall, the values of $k = 54$ and $k = 27$ can be used to describe the top 10% and top 5%, respectively. Note also that the absolute value of the predicted values does not directly affect the fraction of individuals within the predicted top k who are also within the actual top k : only the relative order among the predicted values matters, because it is this order which is used to determine which individuals are predicted as being within the top k .

IV. CASE STUDY

This study uses data from the Big 5 Experiment [20] performed by the Online Privacy Foundation. In this experiment, a total of 537 Facebook users were given a 45-question survey to categorize their personality according to the Big 5 personality index. The Big 5 personality index breaks down human personalities into a number of axes (specifically, five of them) and assigns a number to each axis representing where that individual falls along the axis. Each of these axes is named after the trait represented by being in the extreme upper end of the scale; being in the lower end signifies having the opposite trait. The five axes are Agreeableness, Conscientiousness,

Extroversion, Neuroticism, and Openness. Briefly, these traits may be defined as follows: Agreeableness represents a tendency towards compassion and cooperation, while the opposite is a prevalence of suspicion and antagonism. Conscientiousness shows an individual's self-discipline and devotion to duty; a lack of conscientiousness would be an increase in spontaneous and unplanned behavior. Extroversion is an individual's preference for outgoing social experiences, while the opposite (introversion) is a desire for a lower level of social involvement. Openness represents an individual's curiosity and appreciation for new experiences; the converse would be a greater respect for traditional and well-traveled experiences. Neuroticism, which is related to emotional instability, represents the tendency to experience negative emotions and a lower tolerance to stress. For the purposes of the Big 5 Experiment, all five of these axes are measured on a scale from 1 to 5 (decimal values permitted), with 1 being the lowest possible expression of the axis and 5 the maximal possible expression.

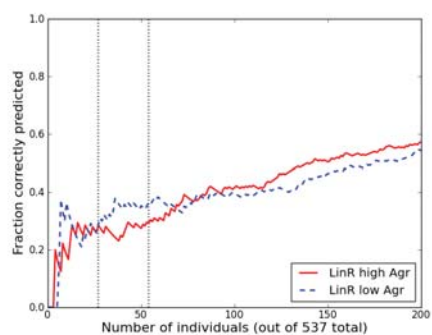
In addition to collecting the scores of the 537 users in the survey, information from their profile was used to generate a set of attributes defining each individual. These attributes were created solely from information the users already were sharing with their friends, not requiring any specific input or surveys. 32 of these attributes are so-called "demographic" data, hard numbers (or categories) acquired from different parts of the profile. Examples of these data include age, gender, locale, length of biography and quotes, relationship status, and the number of friends, photos, interests, and comments provided by the user. One of these properties, political affiliation, was removed from the data due to its free-form nature (users are free to input any text string into this field, and thus parsing it is challenging). Thus, a total of 31 demographic attributes were used for each individual.

Beyond these demographic attributes, 80 text-based attributes were acquired from each user. Rather than directly turning different aspects of the profiles into categorical or numeric data, these attributes arose from taking the users' profiles, photo descriptions, and public posts and processing them using the Linguistic Inquiry and Word Count software package [16]. This tool analyzes text to pull out "language dimensions" including positive and negative emotional content, self-reference, key topics, writing style, and more. It allows large blocks of text to be interpreted in terms of a fixed array of parameters, so that data mining can act on the "soft" content of language beyond its fixed numeric values.

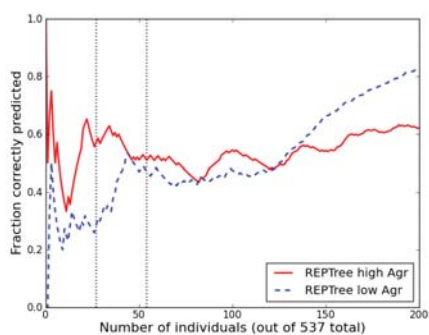
In this preliminary study, all 111 attributes (31 demographic + 80 text-based) were used to build models. Future work may consider the use of feature selection to reduce the size of the feature space.

V. RESULTS

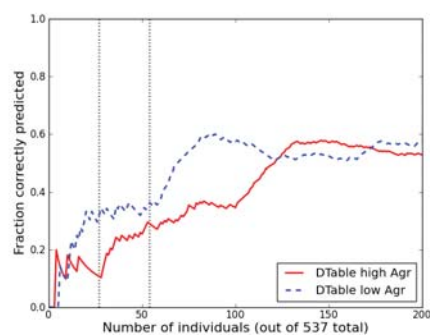
Using the dataset discussed in Section IV and the numeric prediction techniques in Section III, we predicted the five personality traits (Agreeableness, Conscientiousness, Extroversion, Neuroticism, and Openness) using the 111 (31 demographic + 80 text-based) independent attributes. Following this, the instances were ordered in terms of the predicted values for all five class attributes (creating five ranked lists of instances). We compared the top and bottom k instances from each of these lists to the instances found within the equivalent lists ordered using the real (rather than predicted) class attribute values. These percentages are shown in Figures 1 through 5. Each figure shows the results for one of the five class attributes, and the three subfigures within each hold the results for the three different numeric prediction approaches (LinR, REPTree, DTable). Throughout, the solid red line reflects the number of instances predicted as being within the top k which really were in the top k



(a) LinR

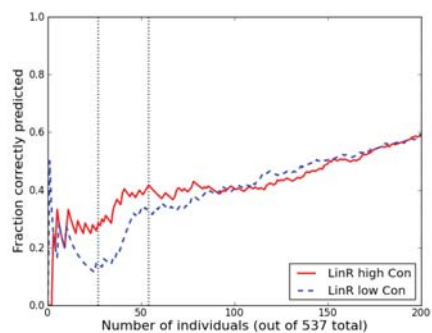


(b) REPTree

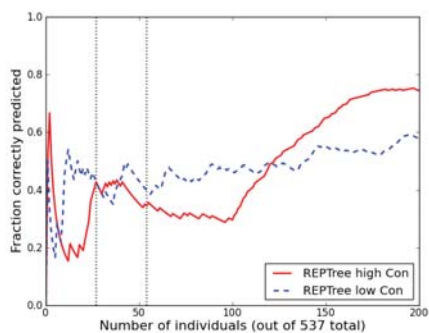


(c) DTable

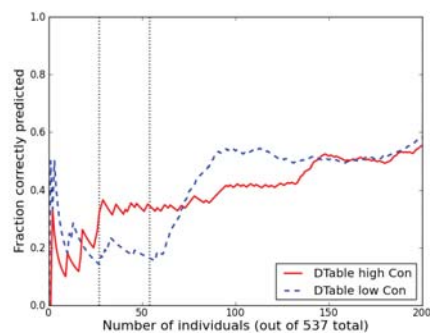
Fig. 1: Results for Agreeableness



(a) LinR

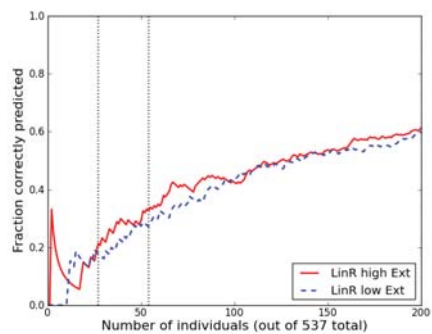


(b) REPTree

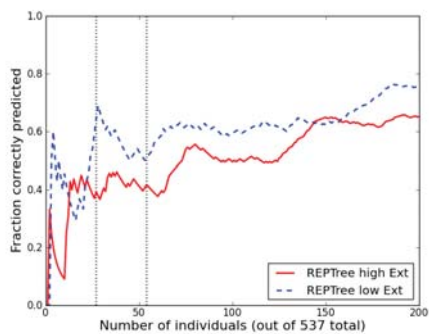


(c) DTable

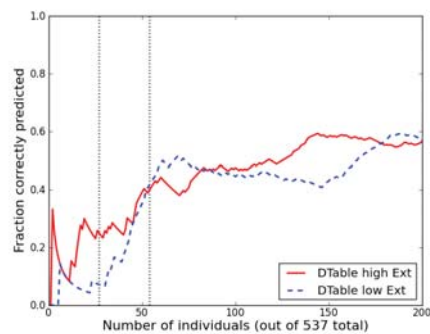
Fig. 2: Results for Conscientiousness



(a) LinR



(b) REPTree



(c) DTable

Fig. 3: Results for Extroversion

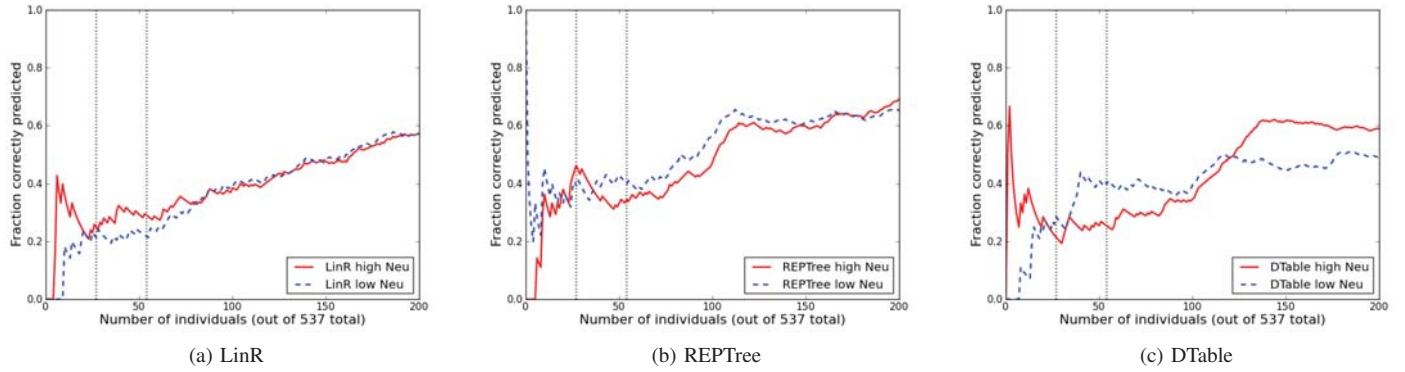


Fig. 4: Results for Neuroticism

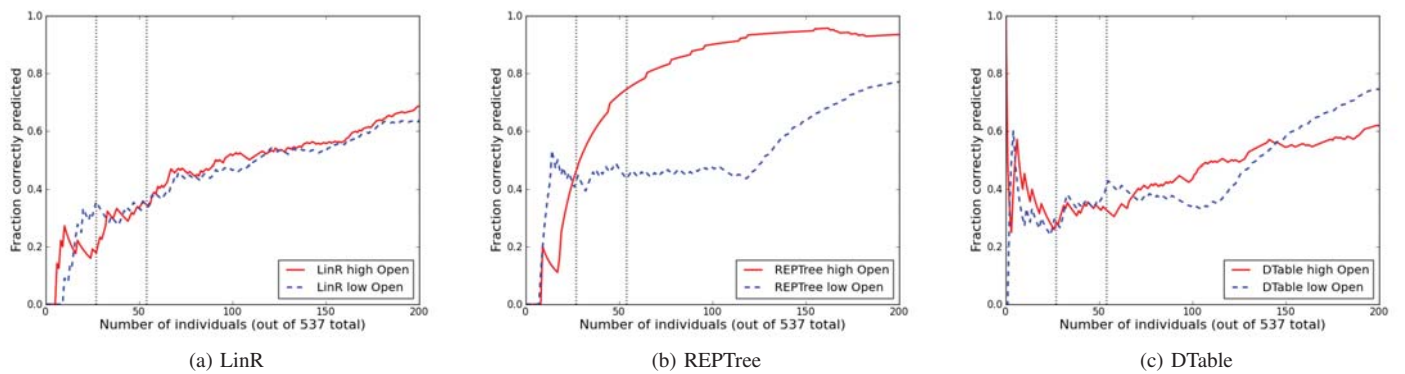


Fig. 5: Results for Openness

instances; the dotted blue line is for those predicted as being within the bottom k which really were in the bottom k . For all graphs, k is varied from 1 to 200. For convenience, the values of $k = 27$ (representing 5% of the instances) and $k = 54$ (representing 10% of the instances) have been marked on the graphs by including a dotted vertical line at each (the leftmost line is 27, the rightmost is 54).

Due to the importance of these values, Table I presents the exact percentage of instances correctly predicted for each combination of trait (class), direction, and numeric prediction approach (method) when the top/bottom 27 (5%) or 54 (10%) individuals are selected. Note that in the final two rows, rather than using one of the three numeric prediction methods to predict the rank, each individual's total number of friends (NoF) was employed. This is presented to compare the use of numeric prediction with simply using a single feature, and the values for Agreeableness and NoF are presented because these create the most accurate prediction models across all traits and features.

The first important observation from the figures and table is that when using REPTree to predict the Openness of individuals based on their Facebook profiles, the model was very good at predicting the top (high) results. Of those predicted to be in the top 5%, 46.2% of these were accurate predictions, and when going to the top 10%, 74.5% of the predictions were accurate. This means that for finding out which users are the most open to new ideas and experiences (a valuable demographic), it is possible to find nearly 3/4 of those in the top 10%. This result is extremely useful because although it does

not directly enable accurate prediction of individual values for the Openness class attribute, it allows broad-scale demographic analysis to specifically focus marketing efforts on a particular subset of users.

Two additional specific values are worth noting. For the top 5% in terms of Agreeableness and the bottom 5% in terms of Extroversion (e.g., the top 5% in terms of Introversion), the REPTree prediction model was able to produce a reasonably accurate ranking (57.1% for Agreeableness, 64.3% for Introversion). Both of these are not entirely consistent results: both drop off to lower performance levels as the values of k move on past the 5% line. Nonetheless, at least within the sweet spot for each combination of REPTree, class attribute, and direction, very high performance (well over 50%) is possible. As with the Openness results above, this information can be used to process a dataset and decide which set of individuals are most likely to have a given personality trait, with reasonable confidence in this result.

Beyond looking at particular choices of class attribute and direction with good results, it is interesting to observe the overall patterns from each of the three numeric prediction methods. LinR generally has consistent results: the fraction correctly predicted increases mostly monotonically as the number of individuals selected increases. In addition, by the time it has selected 10% of the individuals in the population, at least 30% of those were correctly predicted, except in the following cases: high Agreeableness (only 29.1%), low Extroversion (only 27.3%), high Neuroticism (only 29.1%), low Neuroticism (only 21.8%). Of these, only the low Neuroticism example is significantly below 30%; for all other combinations, LinR is a safe choice of

Trait	Direction	Model	5%	10%
Agreeableness	High	LinR	28.571%	29.091%
		REPTree	57.143%	50.909%
		DTable	10.714%	29.091%
	Low	LinR	28.571%	36.364%
		REPTree	28.571%	47.273%
		DTable	32.143%	36.364%
Conscientiousness	High	LinR	28.561%	41.818%
		REPTree	42.857%	34.545%
		DTable	32.143%	34.545%
	Low	LinR	14.286%	32.727%
		REPTree	42.857%	40.000%
		DTable	14.286%	16.364%
Extroversion	High	LinR	21.429%	32.727%
		REPTree	39.286%	41.818%
		DTable	25.000%	40.000%
	Low	LinR	17.857%	27.273%
		REPTree	64.286%	50.909%
		DTable	7.143%	41.818%
Neuroticism	High	LinR	25.000%	29.091%
		REPTree	46.429%	34.545%
		DTable	21.429%	25.455%
	Low	LinR	21.429%	21.818%
		REPTree	39.286%	40.000%
		DTable	28.571%	40.000%
Openness	High	LinR	17.857%	34.545%
		REPTree	46.429%	74.545%
		DTable	28.571%	32.727%
	Low	LinR	35.714%	34.545%
		REPTree	42.857%	43.636%
		DTable	28.571%	41.818%
Agreeableness	High	NoF	10.714%	18.181%
	Low	NoF	17.857%	21.818%

TABLE I: Fraction Correctly Predicted for Top/Bottom 5% and 10%

model, because regardless of the number of individuals chosen, a good (and ever-increasing) fraction will be correctly predicted.

DTable follows a different pattern: generally it has very good results for very few individuals chosen, drops back to some lower performance level, and then steadily rises thereafter. This suggests that DTable is very good at correctly predict the most extreme individuals in the groups, but is not as capable of selecting the group that is slightly farther away from the extreme. Nonetheless, by the time that 10% of the individuals are selected, DTable often has very good results, comparable with those of LinR. Only with low Conscientiousness does DTable have a particularly bad run, not selecting an appreciable fraction of correct instances until more than 1/6 of the individuals have been selected.

Unlike LinR and DTable, REPTree's results are difficult to classify. There is a tendency to have a wide peak of good performance well after the extreme values, with lower performance on either side and good performance eventually returning as more individuals are selected. Partially because this peak often corresponds to the range from 5% to 10% of individuals chosen, REPTree has particularly good results for these values. When selecting the top 10% of individuals, at least 34.5% are always correctly classified. In addition, the highest values overall are found using REPTree (as discussed previously). Overall, the only concern with REPTree is whether the peak of high performance corresponds with the desired fraction of individuals. Because performance will degrade as more users are chosen, it is important to ensure that the number of individuals chosen falls within the specified range.

One final question which may arise is how much the numeric prediction models are adding to the rankings: they are able to

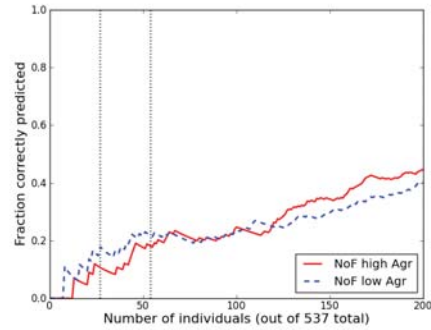


Fig. 6: Results for Agreeableness (Number of Friends only)

consistently find 34.5% of the instances in the top or bottom 10% of each class attribute, but how does this compare to simpler methods, such as sorting in terms of the independent attributes on their own (with no numeric prediction model)? To investigate this, we present the results for ordering the instances in terms of Number of Friends (NoF), and compare the top and bottom k from this ranking with the actual values for the Agreeableness class attribute. We chose these values because NoF and Agreeableness produce the best ranking of individuals when considering all combinations of feature and personality trait, with the highest performances at the 5% and 10% lines. The results are presented in Figure 6, and in the final two lines of Table I. As can be seen, the predicted ranking only rises to a value of 21.8%, and stagnates at that level until more than 100 individuals have been selected (contrast this slope with that seen in Figure 4a, the worst-performing LinR model at the 10% level). Thus, using multiple attributes together to build a better predicted value can help produce much more accurate rankings, allowing for the top and bottom fractions of the instances to be extracted with some degree of confidence in the result.

VI. CONCLUSION

In this study we have used a dataset with 537 individuals, consisting of 31 demographic and 80 text-based independent attributes, to build models to predict the values of the five traits of the Big Five personality model (Agreeableness, Conscientiousness, Extroversion, Openness, Neuroticism). We found that when using the REPTree numeric prediction technique, we are able to rank individuals in terms of predicted Openness and have the top 10% of the predicted list contain 74.5% of the users in the real top 10% most Open. In addition, for both the top 5% Agreeable and Introverted, REPTree gave models with accuracy exceeding 50%. Overall, using the REPTree-based models, we found that for every combination of personality trait and direction (top or bottom), we were able to predict the top 10% with 34.5% accuracy, exceeding 21.8%, the accuracy of simply ordering the users based on the best-performing independent attribute.

The significance of these results is that when considering large groups of users, automatic analysis can be used to identify those with specific personality traits. This can be scaled up far more easily than manual human assessment, and can be used for targeting specific groups of users for advertising, social engineering attacks, or finding influential users. While many are aware of the risks associated with posting inappropriate photographs or posts on social networking sites, this research demonstrates that even innocuous-seeming bibliographic and status information can be used to discern information which users may prefer to remain hidden.

Future research will explore additional techniques to interpret users' personalities through automated data mining analysis of Facebook profiles (for example, by applying feature selection to reduce the number of independent attributes), as well as examine additional datasets from other social networks.

REFERENCES

- [1] A. Acquisti and R. Gross, "Imagined communities: Awareness, information sharing, and privacy on the Facebook," in *Privacy Enhancing Technologies*, ser. Lecture Notes in Computer Science, G. Danezis and P. Golle, Eds. Springer Berlin / Heidelberg, 2006, vol. 4258, pp. 36–58, 10.1007/11957454_3. [Online]. Available: http://dx.doi.org/10.1007/11957454_3
- [2] M. D. Back, J. M. Stopfer, S. Vazire, S. Gaddis, S. C. Schmuttle, B. Egloff, and S. D. Gosling, "Facebook profiles reflect actual personality, not self-idealization," *Psychological Science*, vol. 21, no. 3, pp. 372–374, 2010. [Online]. Available: <http://pss.sagepub.com/content/21/3/372.short>
- [3] B. Debatin, J. P. Lovejoy, A.-K. Horn, and B. N. Hughes, "Facebook and online privacy: Attitudes, behaviors, and unintended consequences," *Journal of Computer-Mediated Communication*, vol. 15, no. 1, pp. 83–108, 2009. [Online]. Available: <http://dx.doi.org/10.1111/j.1083-6101.2009.01494.x>
- [4] N. B. Ellison, C. Steinfield, and C. Lampe, "The benefits of Facebook "friends": social capital and college students' use of online social network sites," *Journal of Computer-Mediated Communication*, vol. 12, no. 4, pp. 1143–1168, 2007. [Online]. Available: <http://dx.doi.org/10.1111/j.1083-6101.2007.00367.x>
- [5] L. R. Goldberg, "An alternative "description of personality": The Big-Five factor structure," *Journal of Personality and Social Psychology*, vol. 59, no. 6, p. 1216, 1990.
- [6] S. D. Gosling, S. Gaddis, and S. Vazire, "Personality impressions based on Facebook profiles," in *International AAAI Conference on Weblogs and Social Media*. AAAI Press, March 2007, pp. 1–4. [Online]. Available: <http://www.icwsm.org/papers/3--Gosling-Gaddis-Vazire.pdf>
- [7] R. Gross and A. Acquisti, "Information revelation and privacy in online social networks," in *Proceedings of the 2005 ACM workshop on Privacy in the electronic society*, ser. WPES '05. New York, NY, USA: ACM, 2005, pp. 71–80. [Online]. Available: <http://doi.acm.org/10.1145/1102199.1102214>
- [8] M. Hall, E. Frank, G. Holmes, B. Pfahringer, P. Reutemann, and I. H. Witten, "The WEKA data mining software: An update," *SIGKDD Explor. Newsl.*, vol. 11, no. 1, pp. 10–18, 2009.
- [9] A. N. Joinson, "Looking at, looking up or keeping up with people?: motives and use of Facebook," in *Proceedings of the twenty-sixth annual SIGCHI conference on Human factors in computing systems*, ser. CHI '08. New York, NY, USA: ACM, 2008, pp. 1027–1036. [Online]. Available: <http://doi.acm.org/10.1145/1357054.1357213>
- [10] H. Jones and J. H. Soltren, "Understanding privacy settings in Facebook with an audience view," in *Proceedings of the 1st Conference on Usability, Psychology, and Security*. USENIX Association Berkeley, CA, USA, 2005, pp. 1–8.
- [11] T. M. Khoshgoftaar and E. B. Allen, "Ordering fault-prone software modules," *Software Quality Journal*, vol. 11, no. 1, pp. 19–37, 2003.
- [12] R. Kohavi, "The power of decision tables," in *Machine Learning: ECML-95*, ser. Lecture Notes in Computer Science, N. Lavrac and S. Wrobel, Eds. Springer Berlin / Heidelberg, 1995, vol. 912, pp. 174–189, 10.1007/3-540-59286-5_57. [Online]. Available: http://dx.doi.org/10.1007/3-540-59286-5_57
- [13] C. Lampe, N. Ellison, and C. Steinfield, "A face(book) in the crowd: social searching vs. social browsing," in *Proceedings of the 2006 20th anniversary conference on Computer supported cooperative work*, ser. CSCW '06. New York, NY, USA: ACM, 2006, pp. 167–170. [Online]. Available: <http://doi.acm.org/10.1145/1180875.1180901>
- [14] C. Lampe, N. B. Ellison, and C. Steinfield, "Changes in use and perception of Facebook," in *Proceedings of the 2008 ACM conference on Computer supported cooperative work*, ser. CSCW '08. New York, NY, USA: ACM, 2008, pp. 721–730. [Online]. Available: <http://doi.acm.org/10.1145/1460563.1460675>
- [15] T. A. Pempek, Y. A. Yermolayeva, and S. L. Calvert, "College students' social networking experiences on Facebook," *Journal of Applied Developmental Psychology*, vol. 30, no. 3, pp. 227–238, 2009. [Online]. Available: <http://www.sciencedirect.com/science/article/pii/S0193397308001408>
- [16] J. W. Pennebaker, M. E. Francis, and R. J. Booth, "Linguistic inquiry and word count: LIWC 2001," *WORD: Journal Of The International Linguistic Association*, pp. 1–21, 2001. [Online]. Available: http://homepage.psy.utexas.edu/homepage/faculty/pennebaker/reprints/LIWC2007_OperatorManual.pdf
- [17] J. Raacke and J. Bonds-Raacke, "MySpace and Facebook: Applying the uses and gratifications theory to exploring friend-networking sites," *CyberPsychology & Behavior*, vol. 11, no. 2, pp. 169 – 174, 2008. [Online]. Available: <http://search.ebscohost.com/login.aspx?direct=true&db=aph&AN=31696364&site=ehost-live>
- [18] B. Rose and T. Webster, "The infinite dial 2011: Navigating digital platforms," Arbitron Inc. and Edison Research, Tech. Rep., April 2011. [Online]. Available: http://www.arbitron.com/study/digital_radio_study.asp
- [19] C. Ross, E. S. Orr, M. Sisic, J. M. Arseneault, M. G. Simmering, and R. R. Orr, "Personality and motivations associated with Facebook use," *Computers in Human Behavior*, vol. 25, no. 2, pp. 578–586, 2009. [Online]. Available: <http://www.sciencedirect.com/science/article/pii/S0747563208002355>
- [20] C. Sumner, A. Byers, and M. Shearing, "Determining personality traits & privacy concerns from Facebook activity," in *Black Hat Briefings*, Dec. 2011. [Online]. Available: https://media.blackhat.com/bh-ad-11/Sumner/bh-ad-11-Sumner-Concerns_w_Facebook_WP.pdf
- [21] S. Valenzuela, N. Park, and K. F. Kee, "Is there social capital in a social network site?: Facebook use and college students' life satisfaction, trust, and participation1," *Journal of Computer-Mediated Communication*, vol. 14, no. 4, pp. 875–901, 2009. [Online]. Available: <http://dx.doi.org/10.1111/j.1083-6101.2009.01474.x>
- [22] I. H. Witten and E. Frank, *Data Mining: Practical machine learning tools and techniques*, 2nd ed. Morgan Kaufmann, 2005.