

# Mining Facebook Data via Text, Images, and Likes\*

Extended Abstract<sup>†</sup>

Wenfei Yin<sup>‡</sup>  
Institute of Technology  
Univ. of Washington: Tacoma  
Tacoma, WA 98402  
yinwf@uw.edu

Shreya Yembarwar<sup>§</sup>  
Institute of Technology  
Univ. of Washington: Tacoma  
Tacoma, WA 98402  
shreyay@uw.edu

Jonathan McFadden<sup>¶</sup>  
Institute of Technology  
Univ. of Washington: Tacoma  
Tacoma, WA 98402  
mcfaddja@uw.edu

## ABSTRACT

## KEYWORDS

ACM proceedings, L<sup>A</sup>T<sub>E</sub>X, text tagging

### ACM Reference format:

Wenfei Yin, Shreya Yembarwar, and Jonathan McFadden. 1997. Mining Facebook Data via Text, Images, and Likes. In *Proceedings of TCSS 555: Machine Learning Term Project, Tacoma, WA USA, June 2018 (TCSS 555: Spring 2018)*, 3 pages.  
[https://doi.org/10.475/123\\_4](https://doi.org/10.475/123_4)

## 1 INTRODUCTION

mnbmb

## 2 METHODOLOGY

hvljh

### 2.1 Text

For the Text part, we separate it into three parts. The first one is using the text to predict the accuracy of age and gender, then we use the 82- LIWC features to predict the big five personalities. For each part, we also use different methods to build the basic model to test accuracy.

We tried four different methods in text part, Linear regression, Random forest, Naïve Bayes and Logistical regression. Finally, we find the best result from the method which uses logistical regression and Naïve regression.

The first step of the text part is data preprocessing. In the files, we find that they are separated, thus we merge the profile file and text file according to the common primary key in each take. We add each text part into profile.csv. At the meanwhile, we also transfer the age to age groups, and we use 1 and 0 to replace the female and male in order to more effective.

After data preprocessing, we use the Naïve Bayes to predict the age and gender. As is known to us who focus on machine learning. Naïve Bayes classifiers are a family of simple probabilistic classifiers? based on applying Bayes theorem with strong independence assumptions between the features. However, we find the model only working well on gender part. The accuracy of age in Naïve Bayes model is 0.49 which is lower than baseline(0.59). Therefore, we have to change our method. After researching, we guess the Logistical Regression will work well on age part. At the beginning, we split the training data into two part, the one is training data which are 8000 rows, the other is test data which is 1500 rows. To our surprise, it does not work well. After trying different methods, like linear regression, random forest. We find the model cannot work well when we use text to predict the age, thus we try to use 82-LIWC features to build model.

The LIWC, The linguistic Inquiry and Word Count tool, is known text analysis software which is widely used in psychology studies, In the file, for each user, it has 82 features. We have to merge the LIWC and profile file. When we create the model, we delete the big five personalities which need to be predicted in the table. After testing in local machine, we find it work well in age predict. The accuracy of age group reach to 0.62.

For Five Big Personality, we have to used 82-LIWC features. Because of big data, we decide to use Linear regression to set the model. In statistics, Linear regression is a linear approach to modeling the relationship between a scalar response and one or more explanatory variables. Fortunately, Linear regression work well in the model, the results almost similar with baseline. Although we try to decrease the RMSE of five big personality, the methods we tried are almost same.

### 2.2 Images

jhljhv

### 2.3 Likes

When considering the Likes data, we first create a **User/LikeID Matrix** from the training data we are given. This process is described below in the Dataset and Metrics section. Based on the training data we were given, this resulted in a matrix with 536204 columns. Due to the large number of columns in this matrix, dimensionality reduction using Singular Value Decomposition or Principal Components Analysis was considered. Unfortunately, when the singular values of this matrix were computed (*up to the 9499th singular value*), there was no characteristic "cliff" drop-off in their magnitudes which would indicate a good point for truncation.

\*Produces the permission block, and copyright information

<sup>†</sup> The full version of the author's guide is available as `acmart.pdf` document

<sup>‡</sup>

<sup>§</sup>

<sup>¶</sup>

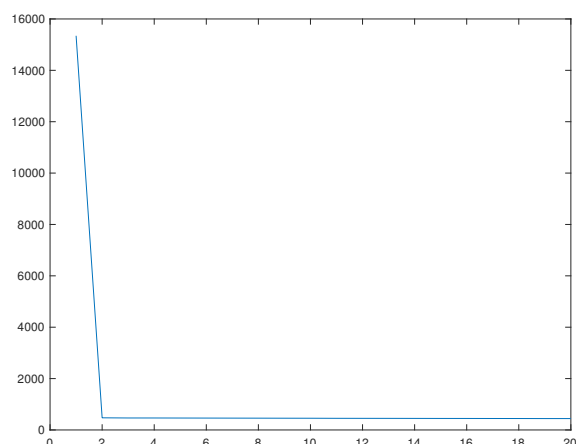
Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).

TCSS 555: Spring 2018, June 2018, Tacoma, WA USA

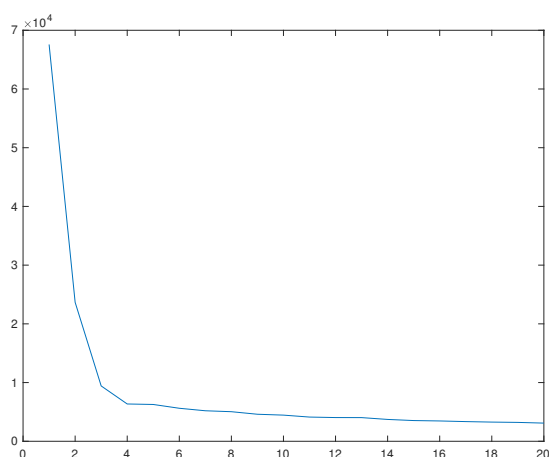
© 2016 Copyright held by the owner/author(s).

ACM ISBN 123-4567-24-567/08/06...\$15.00

[https://doi.org/10.475/123\\_4](https://doi.org/10.475/123_4)



**Figure 1: The type of 'cliff' drop-off in singular value magnitudes we wanted to see.**



**Figure 2: The singular value magnitudes we got instead of what we wanted to see.**

To illustrate this, the type of "cliff" in the singular value magnitudes we were seeking can be seen in **Figure ??**.

Unfortunately, the result was got did not have the sudden drop over just a few singular values. What we got instead can be seen in **Figure ??**.

### 3 DATASET AND METRICS

jkklkj

### 4 RESULTS

kjgkjhg

**Table 1: Text RESULTS for Age & Gender Accuracy**

Baseline	0.59	0.59
Linear Regression	0.52	0.60
Random Forest	0.51	0.60
Naive Bayes	0.49	0.71
Logistical Regression	0.61 <sup>a</sup>	0.59
	0.59 <sup>b</sup>	

<sup>a</sup>Using LIWC data

<sup>b</sup>Using text data

**Table 2: Text RESULTS for Personality Score RMSE**

	OPE	NEU	EXT	AGR	CON
Baseline	0.65	0.80	0.79	0.66	0.73
Linear Regression	0.65	0.79	0.79	0.65	0.72

#### 4.1 Text

For age prediction model, using Logistical regression based on 82-LIWC features is best choice in the method we tried. We also use other methods, line Linear regression, random forest, Naïve Bayes. We can see the results in the chart given in **Table 1**.

From the chart in **Table 1**, we can see that the best result in age prediction is the model when it uses logistical regression based on LIWC, and the best accuracy in gender predict model is Naïve Bayes. Therefore, our choice are Naïve Bayes based on text in gender and Logistical Regression based on LIWC in age prediction.

For Big five personality, we only use linear regression on the text and LIWC part, because we tried other methods, and the results are same. The results are in the chart given in **Table tab:txt-person**

From the chart in **Table tab:txt-person**, we can clearly see that the results under the predict model are similar with Baseline when we use Linear regression.

#### 4.2 Images

jhlglhg

#### 4.3 Likes

ljlkj

### 5 CONCLUSION

ljhgljh

#### 5.1 Text

For text part, we have tried many methods when we set up the model, but we do not have do more specific work on data preprocessing, such as feature selection, word choice. When we predict the age and gender, we can find more significant words, like she, he which will affect our results. Although we get a good result in the age prediction when we use Logistical regression based on 82 LIWC features, it is not perfect. In our last assignment, we find the neural

network may be good at decrease RMSE in big five personality, thus we can try it in the future work.

## **5.2 Images**

jbjlj

## **5.3 Likes**

bljk