

Feature Selection and Parameter Tuning on Age and Gender Classification for English Tweets using SVMs A report for Automatic Classification and Kernel Methods

Roy Khristopher Bayot

Universidade de Évora, Department of Informatics,
Rua Romão Ramalho nº59, 7000-671 Évora, Portugal
d11668@alunos.uevora.pt

1 Introduction

Author profiling has been of importance in the recent years. From a forensic standpoint for example, it could be used to determine potential suspects by getting linguistic profiles and identifying characteristics. From a business intelligence perspective, companies could target specific people through online advertising. By knowing the profile of the authors, companies would easily find what a specific group of people talk about online and devise strategies to advertise to these people. They could also analyze product reviews and know what types of products are liked or disliked by certain people.

Part of the reason why the interest in author profiling grows is because the growth of the internet where text is one of the main forms of communication. Through this growth, various corpora could be extracted, curated, assembled from different sources such as blogs, websites, customer reviews, and even twitter posts. Of course, this presents some problems. For example, people from different countries who use the same online platform such as Twitter or Blogger could behave differently in terms of text usage. This presents a difficulty in profiling. This work tries to take this difficulty into account by studying which kind of features are useful for different languages.

The aim of this work is to investigate the parameters for support vector machines in terms of classification using the dataset given in PAN 2015 [11]. The dataset contains twitter data from 4 different languages which are used to profile an author based on age, gender, and 5 personality traits - agreeability, conscientiousness, extrovertedness, openness, and stability. The four languages are English, Dutch, Italian, and Spanish. However, the focus of this work is on English alone and only in age and gender classification. Furthermore, the investigation is more on using different kernels and different parameters for the classification.

2 State of the Art

One of the first few works on author profiling is that of Argamon et al. in [1] where texts are categorized based on gender, age, native language, and personality. For personality, only neuroticism was checked. The corpus comes from different sources. The age and gender have the same corpus taken from blog postings. The native language corpus was taken from International Corpus of Learner English. Personality was taken from essays of psychology students from University of Texas in Austin. Two types of features were obtained: content-based features and style-based features and Bayesian Multinomial Regression was used as a classifier. Argamon et al. had some interesting results where from the gender task, they were able to achieve 76.1% accuracy using style and content features. For age task with 3 classes, the accuracy was at 77.7% also using style and content features. For the native language task, the classifiers were able to achieve 82.3% using only content features. And finally, in checking for neuroticism, the highest obtained was 65.7% using only style features.

There has also been some research that uses datasets collected from social media. A particular example is that of Schler et al. in [13] where writing styles in blogs are related to age and gender. Stylistic and content features were extracted from 71,000 different blogs and a Multi-Class Real Winnow was used to learn the models to classify the blogs. Stylistic features included parts-of-speech tags, function words, hyperlinks, and non-dictionary words. Content features included word unigrams with high information gain. The accuracy achieved was around 80% for gender classification and 75% for age identification.

The work of Argamon et al. [1] became the basis for the work in PAN. It is an ongoing project from CLEF with author profiling as one of its tasks. It currently has three editions. In the first edition of PAN [10] in 2013, the task was age and gender profiling for English and Spanish blogs. There were a variety of methods used. One set includes content-based features such as bag of words, named entities, dictionary words, slang words, contractions, sentiment words, and emotion words. Another would be stylistic features such as frequencies, punctuations, POS, HTML use, readability measures, and other various statistics. There are also features that are n-grams based, IR-based, and collocations-based. Named entities, sentiment words, emotion words, and slang, contractions and words with character flooding were also considered. After extracting the features, the classifiers that were used were the following - decision trees, Support Vector Machines, logistic regression, Naïve Bayes, Maximum Entropy, Stochastic Gradient Descent, and random forests. The work of Lopez-Monroy in [4] was considered the winner for the task although they placed second for both English and Spanish with an accuracy of 38.13% and 41.58% respectively. They used second order representation based on relationships between documents and profiles. The work of Meina et al. [7] used collocations and placed first for English with a total accuracy of 38.94%. On the other hand, the work of Santosh et al. in [12] gave a total accuracy of 42.08% after using POS features for Spanish.

In PAN 2014 [9], the task was profiling authors with text from four different sources - social media, twitter, blogs, and hotel reviews. Most of the approaches

used in this edition are similar to the previous year. In [3], the method used to represent terms in a space of profiles and then represent the documents in the space of profiles and subprofiles were built using expectation maximization clustering. This is the same method as in the previous year in [4]. In [5], ngrams were used with stopwords, punctuations, and emoticons retained, and then idf count was also used before placed into a classifier. Liblinear logistic regression returned with the best result. In [16], different features were used that were related to length (number of characters, words, sentences), information retrieval (cosine similarity, okapi BM25), and readability (Flesch-Kincaid readability, correctness, style). Another approach is to use term vector model representation as in [15]. For the work of Marquardt et al. in [6], they used a combination of content-based features (MRC, LIWC, sentiments) and stylistic features (readability, html tags, spelling and grammatical error, emoticons, total number of posts, number of capitalized letters number of capitalized words). Classifiers also varied for this edition. There was the use of logistic regression, multinomial Naïve Bayes, liblinear, random forests, Support Vector Machines, and decision tables. The method of Lopez-Monroy in [3] gave the best result with an average accuracy of 28.95% on all corpus-types and languages.

3 Dataset and Tools

The dataset for the problem at hand is composed of a set of tweets for English. Different models were made for each classification task - age and gender. There were 4 categories for the age classification - 18-24, 25-34, 35-49, and 50 and above. Gender has two categories - male and female. There were 152 users for English. Each user has different number of tweets. The dataset is balanced based on gender. Processing the data was done through Python using the scikits-learn [8] library.

4 Methodology

The focus of this study is to determine optimal parameters for classification using support vector machines. The approach is to do preprocessing, feature extraction, feature selection, then use the optimal features into support vector machines with different parameters. Evaluation was made through 10 fold cross validation and Wilcoxon signed rank test was used to compare the statistical significance of the results.

4.1 Preprocessing and Feature Extraction

For each language, xml files from each user are read. Then the tweets taken from each user are extracted and concatenated into one line to form one training example. The examples are then transformed by putting them all in lower case. No stop words are removed. Hashtags, numbers, mentions, shares, and retweets were not processed or transformed to anything else. The resulting file is used for feature extraction. Features extracted were simply *tfidf* features.

4.2 Feature Selection

Classification was done using Support Vector Machines [2] with a linear kernel, choosing $C=1$ as the default. Different number of *tfidf* features were tested. The number of features were ranked by either information gain or gain ratio.

The results for feature selection are given in table 1. When using all the different features, the highest accuracy could be attained. However, to check if the other accuracies are not statistically different, Wilcoxon signed rank test was done. The p-value results are given in table 2. For age classification with words ranked with information gain and gain ratio, the results begin to be statistically different when the number of features were 8000. The number of features chosen for further age classification experiments was 9000. On the other hand the results become statistically different when the number of features was 5000 for gender classification. Therefore, The number of features chosen for further age classification experiments was 7000.

Num Features	Age		Gender	
	info gain	gain ratio	info gain	gain ratio
26263	0.6913	0.6913	0.6825	0.6825
10000	0.6321	0.6321	0.6167	0.6167
9000	0.6321	0.6321	0.6163	0.6163
8000	0.6188	0.6188	0.6233	0.6233
7000	0.6188	0.6188	0.6300	0.6300
5000	0.6188	0.6188	0.4863	0.4863
2000	0.6188	0.6188	0.4983	0.4983
1000	0.6188	0.6188	0.4983	0.4983
700	0.6188	0.6188	0.4921	0.4921
500	0.6188	0.6188	0.4921	0.4921
300	0.6188	0.6188	0.4921	0.4921
200	0.6188	0.6125	0.4921	0.4921
100	0.6188	0.3950	0.4921	0.4921

Table 1. Accuracies that result from using different number of important features ranked by information gain and gain ratio. SVM linear kernel with $C=1$ was used.

4.3 Learning algorithm

After getting the features, different two different kernels were tested with different parameters - polynomial kernel and radial basis function kernel. For polynomial kernel, we set the gamma to 1 and for age classification, the C parameter is given to be either 0.0001, 0.001, 0.1, 1, 10, 1000, 10000. The degree varies between 1, 2, and 3. For the polynomial kernel in gender classification, the gamma is also set to 1, C parameter is either 0.0001, 0.1, 1, 10, 10000. For the radial basis function parameters in both tasks, the C and gamma vary between 0.0001, 0.001, 0.01, 1, 100, 1000, 10000.

Num Features	Age		Gender	
	info gain	gain ratio	info gain	gain ratio
26263	1.0000	1.0000	1.0000	1.0000
10000	0.0757	0.0757	0.1212	0.1212
9000	0.0757	0.0757	0.1124	0.1124
8000	0.0211	0.0211	0.1306	0.1306
7000	0.0211	0.0211	0.1859	0.1859
5000	0.0211	0.0211	0.0002	0.0002
2000	0.0211	0.0211	0.0012	0.0012
1000	0.0211	0.0211	0.0015	0.0015
700	0.0211	0.0211	0.0006	0.0006
500	0.0211	0.0211	0.0006	0.0006
300	0.0211	0.0211	0.0006	0.0006
200	0.0211	0.0211	0.0006	0.0006
100	0.0211	0.0002	0.0006	0.0006

Table 2. Wilcoxon signed rank test results comparing accuracies give by the highest number of features and those with lower number of features.

5 Results and Discussion

The results for age classification using support vector machines with polynomial kernel is given by table 3. The highest accuracy obtained was 80.92%. This was obtained from three different settings, degree 3 with C to be either 10, 1000, 10000. However, checking on Wilcoxon signed rank test, these results are not statistically significant to 80.25% given by settings with degree 2 and C to be either 10, 1000, 10000. These results are also not statistically significant against 75% with degree 3 and C to be 1.

C	degree		
	1	2	3
0.0001	0.6321	0.6192	0.5121
0.001	0.6321	0.6192	0.5121
0.1	0.6321	0.6254	0.5992
1	0.6321	0.7104	0.7500
10	0.6321	0.8025	0.8092
1000	0.6321	0.8025	0.8092
10000	0.6321	0.8025	0.8092

Table 3. Accuracy results of using SVM with polynomial kernel on age classification task with different C and degree parameters using top 9000 informative features ranked by information gain.

The results for age classification using support vector machines with radial basis function kernel is given by table 4. Results show that the highest achieved

accuracy was 80.92% from the kernel with gamma to be 0.001 and with a C to be 10000. This wasnt statistically different from the settings that gave 80.25%.

C	gamma						
	0.0001	0.001	0.01	1	100	1000	10000
0.0001	0.3950	0.3950	0.3950	0.3950	0.3950	0.3950	0.3950
0.001	0.3950	0.3950	0.3950	0.3950	0.3950	0.3950	0.3950
0.1	0.3950	0.3950	0.3950	0.3950	0.3950	0.3950	0.3950
1	0.3950	0.3950	0.3950	0.6517	0.3950	0.3950	0.3950
10	0.3950	0.6121	0.6908	0.8025	0.3950	0.3950	0.3950
1000	0.6188	0.6971	0.8025	0.8025	0.3950	0.3950	0.3950
10000	0.6971	0.8092	0.8025	0.8025	0.3950	0.3950	0.3950

Table 4. Accuracy results of using SVM with radial basis function kernel for age classification with different C and gamma parameters using top 9000 informative features ranked by information gain.

The results for gender classification using support vector machines with polynomial function kernel is given by table 5. Results show that the highest achieved accuracy was 74.29%. This comes from three different settings. The first is when the degree is 2 and C was chosen to be 10000. The second and third settings were when degree is 3 with C to be either 10 or 10000. However, these arent statistically different from settings which gave 71.71% and 67.04% accuracy.

C	degree		
	1	2	3
0.0001	0.6183	0.3950	0.3950
0.1	0.6183	0.6183	0.3950
1	0.6183	0.6446	0.6704
10	0.6183	0.7171	0.7429
10000	0.6183	0.7429	0.7429

Table 5. Accuracy results of using SVM with polynomial kernel for gender classification with different C and degree parameters using the top 7000 informative features ranked by information gain.

The results for gender classification using support vector machines with radial basis function kernel is given by table 6. Results show that the highest achieved accuracy was 74.29%. This comes from the settings with gamma to be 0.01 and C to be 10000. However, these arent statistically different from settings which gave 73.63% and 71.71% accuracy.

C	gamma						
	0.0001	0.001	0.01	1	100	1000	10000
0.0001	0.3950	0.3950	0.3950	0.3950	0.3950	0.3950	0.3950
0.001	0.3950	0.3950	0.3950	0.3950	0.3950	0.3950	0.3950
0.1	0.3950	0.3950	0.3950	0.3950	0.3950	0.3950	0.3950
1	0.3950	0.3950	0.3950	0.6250	0.3950	0.3950	0.3950
10	0.3950	0.3950	0.6317	0.7363	0.3950	0.3950	0.3950
1000	0.3950	0.6317	0.7171	0.7363	0.3950	0.3950	0.3950
10000	0.6317	0.7171	0.7429	0.7363	0.3950	0.3950	0.3950

Table 6. Accuracy results of using SVM with radial basis function kernel for gender classification with different C and gamma parameters using the top 7000 informative features ranked by information gain.

6 Conclusions and Future Work

Given the different experiments, it could be shown that the highest accuracy that could be achieved were 80.92% and 74.29% for age and gender respectively. The highest accuracies could be attained could either be from polynomial kernels or radial basis function kernels. For instance, for age classification, the highest accuracy could either be obtained with a polynomial kernel with degree 3 and C to be 10 or it could also be a through a radial basis function with gamma to be 0.001 and C to be 10000. For gender classification, the similarity exists as well. A polynomial kernel could be used with the second degree and C to be 10000 and it will yield the same accuracy as the one with a radial basis function with gamma to be 0.01 and and C to be 10000.

We have to note that in this set of experiments, the features used were only *tfidf*. However, using other features in conjunction with parameter tuning, can also yield better results. Other features can include POS tags, text and character ngrams, length of words, ratios between uppercase and lowercase, non-dictionary words, emoticons, and lexical diversity. Furthermore, different types of preprocessing can also be done for this thing.

References

1. Shlomo Argamon, Moshe Koppel, James W Pennebaker, and Jonathan Schler. Automatically profiling the author of an anonymous text. *Communications of the ACM*, 52(2):119–123, 2009.
2. Corinna Cortes and Vladimir Vapnik. Support-vector networks. *Machine learning*, 20(3):273–297, 1995.
3. A Pastor López-Monroy, Manuel Montes-y Gómez, Hugo Jair Escalante, and Luis Villaseñor-Pineda. Using intra-profile information for author profiling.
4. Adrian Pastor Lopez-Monroy, Manuel Montes-y Gomez, Hugo Jair Escalante, Luis Villaseñor-Pineda, and Esaú Villatoro-Tello. Inaoe’s participation at pan’13: Author profiling task. In *CLEF 2013 Evaluation Labs and Workshop*, 2013.
5. Suraj Maharjan, Prasha Shrestha, and Tamar Solorio. A simple approach to author profiling in mapreduce.

6. James Marquardt, Golnoosh Farnadi, Gayathri Vasudevan, Marie-Francine Moens, Sergio Davalos, Ankur Teredesai, and Martine De Cock. Age and gender identification in social media. *Proceedings of CLEF 2014 Evaluation Labs*, 2014.
7. Michał Meina, Karolina Brodzinska, Bartosz Celmer, Maja Czoków, Martyna Patera, Jakub Pezacki, and Mateusz Wilk. Ensemble-based classification for author profiling using various features. *Notebook Papers of CLEF*, 2013.
8. F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.
9. Francisco Rangel, Paolo Rosso, Irina Chugur, Martin Potthast, Martin Trenkmann, Benno Stein, Ben Verhoeven, and Walter Daelemans. Overview of the 2nd author profiling task at pan 2014. In *Proceedings of the Conference and Labs of the Evaluation Forum (Working Notes)*, 2014.
10. Francisco Rangel, Paolo Rosso, Moshe Moshe Koppel, Efstathios Stamatatos, and Giacomo Inches. Overview of the author profiling task at pan 2013. In *CLEF Conference on Multilingual and Multimodal Information Access Evaluation*, pages 352–365. CELCT, 2013.
11. Francisco Rangel, Paolo Rosso, Martin Potthast, Benno Stein, and Walter Daelemans. Overview of the 3rd Author Profiling Task at PAN 2015. In L Cappellato, N Ferro, J Gareth, and E San Juan, editors, *CLEF 2015 Labs and Workshops, Notebook Papers*. CEUR-WS.org, 2015.
12. K Santosh, Romil Bansal, Mihir Shekhar, and Vasudeva Varma. Author profiling: Predicting age and gender from blogs. *Notebook Papers of CLEF*, 2013.
13. Jonathan Schler, Moshe Koppel, Shlomo Argamon, and James W Pennebaker. Effects of age and gender on blogging. In *AAAI Spring Symposium: Computational Approaches to Analyzing Weblogs*, volume 6, pages 199–205, 2006.
14. Helmut Schmid. Probabilistic part-of-speech tagging using decision trees. In *Proceedings of the international conference on new methods in language processing*, volume 12, pages 44–49. Citeseer, 1994.
15. Julio Villena-Román and José Carlos González-Cristóbal. Daedalus at pan 2014: Guessing tweet author’s gender and age.
16. Edson RD Weren, Viviane P Moreira, and José PM de Oliveira. Exploring information retrieval features for author profiling—notebook for pan at clef 2014. *Cappellato et al.[6]*.