

Gender Classification of Web Authors Using Feature Selection and Language Models

Abstract. In the present article, we address the problem of automatic gender classification of web blog authors. More specifically, we employ eight widely used machine learning algorithms, in order to study the effectiveness of feature selection on improving the accuracy of gender classification. The feature rankings are performed over a set of statistical, part-of-speech tagging and language model features. In the experiments, we employed classification models based on decision trees, support vector machines and lazy-learning algorithms. The experimental evaluation performed on blog author gender classification data demonstrated the importance of language model features for this task and that feature selection significantly improves the accuracy of gender classification, regardless of the type of the machine learning algorithm used.

Keywords: text classification, gender identification, feature selection

1 Introduction

The growth of social media and its users is impressive and the exploration of the available information in terms of topic, author, genre, etc., is a basic task. Every user of social media leaves his digital traces, makes transactions, expresses his opinion about things and describes moments of his life. Many trends have sprung up over tweets, blogposts and Facebook statuses. These trends express not only an individual user, but often an entire social group. Therefore, it is an interesting task to detect demographic characteristics, such as gender, in users' text data. Information about the gender can be derived not only from the data the user provides about himself, but also implicitly, from the linguistic choices he/she makes. The automatic extraction of information from the everyday enormously growing volumes of data related to the gender, age and other demographic characteristics of the user are essential in the e-government, security and e-commerce market.

The online user's attitude can be observed and explained from a social perspective and his/her "digital traces"[9] may be very informative about current trends in any domain. The user's online activity leaves several elements, not only about his/her preferences or transactions, but also about his/her identity. The user, consciously or not, provides information about his/her social status, gender, age even his/her educational level and profession.

In the present article, we perform feature ranking and subset selection, aiming to improve the accuracy of author gender classification on web blogs. The feature selection is performed over a large set of features, using statistical, part-of-speech tagging and language model based feature extraction methodologies. These text features were

evaluated by machine learning classification algorithms, in order to evaluate the gender classification performance for different numbers of features.

The rest of the paper is organized as follows. In Section 2 we present the related work. In Section 3 we describe the proposed methodology for gender identification of web users from blog posts. In Section 4 we demonstrate and analyze the experimental results and finally, in Section 5 we conclude this work.

2 Related Work

User's gender detection can be perceived as a text classification task, where machine learning techniques are used to identify the author's gender [4, 15]. Koppel et al. [8] propose text classification methods to extract the author's gender from formal texts, using features such as n-grams and function words that are more frequent in authorship attribution. They combine stylometric and text classification techniques, in order to extract the author's gender. In a subsequent study [2], Argamon et al., by applying factor analysis and machine learning techniques, contribute with gender and age information, in texts mined from the blogosphere. Ansari et al. [1] use frequency counting of tokens, tf-idf and POS-tags to find the gender of blog authors. For gender recognition in Twitter, Burger et al. [3] exploit the content of the tweet combined with the username and other information related to the user.

Recent studies in gender estimation [7, 11, 12, 18] deal with social media and propose methods detecting the gender and, in some cases, the age of the users. The experiments were implemented with gender-polarized words, POS tags and sentence length among other features. In Sarawgi et al. [14], they perform a comparative study of gender attribution, without taking into account the topic or the genre of the selected text. In [17] they use the Naïve Bayes classifier, combined with text features and features such as the web page color. Holgrem and Shyu [5] applied machine learning techniques, using a feature vector containing word counts, in order to detect the author's gender of Facebook statuses. Rangel and Rosso [13] introduced a set of stylistic features to extract the gender and age of authors, using a large set of documents from the social web, written in Spanish. Finally, Marquardt et al. [10] focused on detecting the best feature set towards age and gender prediction in social media.

3 Proposed Gender Identification Methodology

For the feature selection on the task of blog authors' gender identification, we adopted a standard approach followed in most of the previous related work, i.e. pre-processing, feature extraction and classification structure was utilized, as shown in Figure 1.

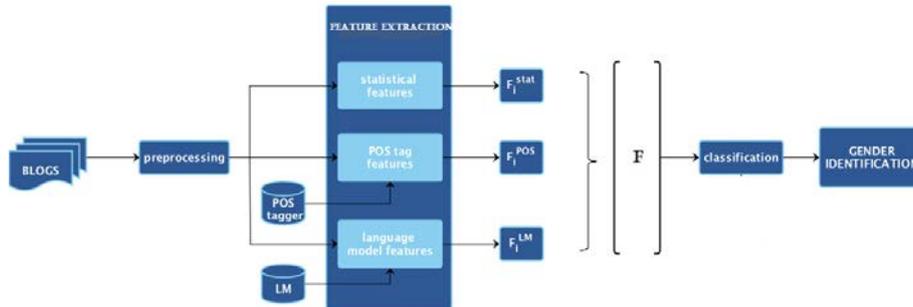


Fig.1.Block diagram of the blog post users' gender identification scheme.

Specifically, each blog post is initially preprocessed. During pre-processing, each post is split into sentences and each sentence is split into words. Afterwards, three feature extraction methodologies are applied in parallel and independently to each other to each post. In detail, statistical, part-of-speech (POS) tag and language model features are extracted constructing vectors F_i^{stat} , F_i^{POS} and F_i^{LM} , respectively. These features are consequently concatenated to a super vector $F = F_i^{stat} || F_i^{POS} || F_i^{LM}$. This results to one feature vector, F , per blog post, which is processed by a classification algorithm, in order to label each post with a gender class.

3.1 Feature Extraction

Three categories of features were computed for each blog post, namely the statistical features, the POS tag features and the language model features.

As far as the statistical features are considered, they consist of statistical values in character and word level. The statistical features that we employed are the following: the number of characters per web post; the normalized number of characters in capital; the normalized number of alphabetic characters; the normalized number of space characters; the normalized number of tab (" \backslash t") characters; the number of occurrence of each alphabetic character; the normalized number of digit characters; the normalized number of occurrence of special characters ("@", "#", "\$", "%", "&", "*", "~", "^", "-", "=", "+", ">", "<", "[", "]", "{", "}", "|", "\", "/"); the total number of words; the normalized number of words that consist of less than 4 characters (short words); the normalized number of characters per word; the average word length; the number of sentences; the number of paragraphs; the number of lines; the average number of characters per sentence; the average number of words per sentence; the normalized number of different words; the number of words that appear once in the document; the number of words that appear twice in the document; the number of punctuation symbols (".", ",", "!", "?", ":", ";", "'", "''"); the number of function words; the average number of sentences per paragraph; the average number of characters per paragraph; the normalized number of words that start with a capital letter; the normalized number of emoticons; the normalized number of words whose letters are all capital; the standard deviation of the word length; the maximum word length; the minimum word

length. All the above features compound the F_i^{stat} feature vector, which has dimensionality equal to 30.

POS tag features, which mainly represent a particular part of speech for every word in a given text, were then computed. These are: the number of nouns; the number of proper nouns; the number of adjectives; the number of prepositions; the number of verbs; the number of pronouns; the number of interjections; the number of adverbs; the number of articles. The F_i^{POS} feature vector comprises all these features and has dimensionality equal to 9.

Finally, for the language model features, we use 2 unigram, bigram and trigram language models, one for each gender class (female and male), in order to measure loglikelihood and entropy as well as their normalized values for each model. These features are language independent. The F_i^{LM} feature vector contains all the above features and, therefore, is a 24-dimensional feature vector.

The concatenation of the three feature vectors results to F , as described above. Thus, for each blog post, one final feature vector, F , is constructed, which has dimensionality equal to $30+9+24=63$. For the estimation of the above text features we used the NLTK [19] open-source toolkit.

3.2 Feature Selection

For the evaluation of the importance of the features we relied on ReliefF algorithm [6]. ReliefF evaluates the worth of a feature by repeatedly sampling an instance and considering the value of the given feature for the nearest instance of the same and different class. ReliefF can operate on both discrete and continuous class data. In the present evaluation we used the Ranker [16] search method. Ranker method ranks features by their individual evaluations. For the feature selection step we used the WEKA machine learning toolkit software [20] implementation.

3.3 Classification

For the classification stage, we used a number of dissimilar machine learning algorithms, which are well studied and have extensively been used in several text classification tasks. In particular, we used a multilayer perceptron neural network (MLP), using the back-propagation algorithm for training and three layers, and support vector machines (SVMs) using the sequential minimal optimization algorithm, which were tested using radial basis kernel (rbf) and polynomial kernel (poly). In addition, we used four tree algorithms, namely the pruned C4.5 decision tree (J48), the random tree (RandTree), constructing a tree that considers K randomly chosen attributes at each node, the random forest (RandForest) constructing a forest of random trees and the fast decision tree learner (RepTree), that builds a decision tree using information gain or variance and prunes it using reduced-error pruning with back-fitting. Finally, we employed one lazy-learning algorithm, the k -nearest neighbor, IBk. All classifiers were implemented using the WEKA toolkit [20].

4 Experimental Setup and Evaluation

The feature selection methodology for blog post authors' gender classification presented in the previous section was evaluated using the statistical, POS tag and language model features and the classifiers presented above. In order to avoid overlap between training and test subsets, a 10-fold cross validation evaluation protocol was followed.

For our evaluation we used the publically available "Blog author gender classification dataset", which was introduced in the work of Mukherjee and Liu [11]. It contains blog posts in English, 1390 written by female authors and 1546 written by male authors. The dataset includes 1319917 words and 6085202 characters. Each post was labeled as female or male, depending on the gender of its author, which was determined by the profile of each author.

At first, we tested the discriminative ability of the language model features, so we used only these features with the classifiers. The dimensionality of the feature vector is, therefore, 24. Table 1 shows the experimental results, in terms of percentages of correctly classified blog posts. SVM using polynomial kernel achieves the best accuracy, which is equal to 69.35%.

Table 1. Accuracy for gender classification using only the language model features

J48	MLP	SVMrbf	SMVpoly	RandForest	RandTree	REPTree	IBk
68.19	68.02	67.91	69.35	67.92	60.42	67.40	68.66

As a second step, we evaluated the appropriateness of all the features on the task of web blog authors' gender identification, using the ReliefF criterion. As seen in Table 2, the language model based features were found to be the most relevant to gender classification, since sixteen out of the top-20 features derive from the language model based features we trained.

Table 2. Ranking results for the top-20 features

rank	score	feature description	feature type
1	0.00769	norm. entropy of female unigram lang. model	language model
2	0.00769	norm. entropy of male unigram lang. model	language model
3	0.00769	norm. loglikel. for female unigram lang. model	language model
4	0.00769	norm. loglikel. for male unigram lang. model	language model
5	0.00668	norm. entropy of female trigram lang. model	language model
6	0.00668	norm. entropy of male trigram lang. model	language model
7	0.00668	norm. loglikel. for female trigram lang. model	language model
8	0.00668	norm. loglikel. for male trigram lang. model	language model
9	0.00666	entropy for male trigram lang. model	language model
10	0.00489	norm. entropy for male bigram lang. model	language model
11	0.00489	norm. entropy for female bigram lang. model	language model
12	0.00489	norm. log likelihood for male bigram lang. model	language model

13	0.00489	norm. loglikel. for female bigram lang. model	language model
14	0.00478	entropy for female trigram lang. model	language model
15	0.00354	norm. number of characters per word	statistical
16	0.00351	entropy for male bigram lang. model	language model
17	0.00309	entropy for female bigram lang. model	language model
18	0.00241	number of short words	statistical
19	0.00189	number of different words	statistical
20	0.00182	digits	statistical

On the other hand, the absence of POS tag features from the top-20 list implies that these features are weakly correlated to the gender of an author. It is worth mentioning that the normalized number of characters per word, short words, different words and digits were found important for classifying a blog post as male or female, considering that they are ranked in the top-15, top-18, top-19 and top-20 features respectively.

In total, 62 out of the 63 text features were found to be to some degree relevant with the gender classification problem, i.e. they demonstrated positive attribute quality value. The remaining one, namely the normalized number of tab characters, obtained a negative value, which means that it is not a relevant attribute with respect to the gender classification problem.

The classification accuracy was evaluated in terms of percentages of correct classified blog posts. The experimental results, in percentages, for the eight different machine learning techniques and for each of the 6 different subsets of top- n features with $n = \{10, 20, 30, 40, 50, 60\}$ are presented in Table 3. The last column presents the classification accuracies in the case that we do not perform feature selection ("All") and the dimensionality of the feature vector is 63. The best performing subset of features for each algorithm is indicated in bold. The top 10 features consist only of language model features. The top 20 features contain 16 language model features and so do the top 30 features. The top 40 features include 21 language model features; the top 50 features include 23 language model features and finally, the top 60 features contain all the language model features.

Table 3. Accuracy for gender classification per feature subset and classification algorithm

	top10	top20	top30	top40	top50	top60	All
J48	68.32	67.92	67.44	67.54	67.47	67.68	65.46
MLP	68.46	67.88	68.56	68.56	68.66	67.27	67.57
SVMrbf	68.02	68.26	68.32	68.42	68.36	68.32	67.91
SVMpoly	69.43	69.24	69.75	69.75	69.93	70.27	69.72
RandForest	69.52	68.46	69.41	70.50	70.03	69.96	69.69
RandTree	67.01	68.43	67.17	65.90	66.62	66.55	61.31
REPTree	67.23	67.51	66.49	66.79	66.38	67.34	66.62
IBk	68.12	68.29	68.66	68.56	68.56	68.49	63.01

Comparing the results obtained with the use of feature selection to those using all 63 features, i.e. with no feature selection, we observe that for all classification algo-

rithms the former results are better compared to the latter. This fact verifies our initial hypothesis, that the feature selection is able to improve the accuracy of gender classification. As seen in Table 3, the improvement in accuracy varies from 0.51% in SVM using rbf kernel to 7.12%, in the case of Random Tree. The IBk shows a significant improvement in performance, which is approximately 6%, followed by J48 and MLP, with improvements of 2.8% and 1% respectively.

The best classification accuracy, 70.50%, was achieved by combining the top-40 features with Random Forest, followed by the 70.27%, achieved using the top-60 features with support vector machine using polynomial kernel. It is worth mentioning that SVM using polynomial kernel performs better with a larger number of features. This is owed to the curse of dimensionality phenomenon from which SVMs do not suffer. On the other hand, J48 performs best when using only the top 10 features. In general, the top 10 and top 20 features obtain quite high results and specifically, Random Tree and REP Tree achieve their best performance using only the top 20 features. This shows the great impact of the language model features on the gender classification problem, since the top 20 features consist mainly of them. Finally, as can be seen in Table 1, the use of the language model features only, succeeded a quite high performance, having accuracy equal to 69.35%.

5 Conclusion

In this article, we investigated the effectiveness of feature selection for improving the accuracy of web blog authors' gender identification. A combination of statistical, part of speech tagging and language model based features was used to represent each web post. We evaluated the gender identification performance using eight classification algorithms for different subsets of features according to ranking scores, produced by the ReliefF algorithm. Random Forest algorithm outperformed all the evaluated algorithms, when using the top 40 text-based features, with accuracy equal to 70.50%. The experimental results showed that the use of subsets of features improved the overall gender identification accuracy for all the evaluated algorithms. Finally, the superiority of the language model features was proved, since the best classification results were obtained, mainly, due to their contribution.

References

1. Ansari, Y.Z., Azad, S.A., Akhtar, H. 2013. Gender Classification of Blog Authors. International Journal of Sustainable Development and Green Economic, ISSN No.: 2315-4721, V-2, I-1, 2, 2013.
2. Argamon, S., Koppel, M., Pennebaker, W., and Schler, J. 2007. Mining the Blogosphere: Age, gender and the varieties of self-expression. First Monday. 12, 9 (September 2007). DOI=[Http://www.uic.edu/htbin/cgiwrap/bin/ojs/index.php/fm/article/view/2003](http://www.uic.edu/htbin/cgiwrap/bin/ojs/index.php/fm/article/view/2003).
3. Burger, J., Henderson, J., Kim, G., Zarrella, G. 2011. Discriminating gender on Twitter. In *Proceedings of the Conference on Empirical Methods in Natural Language*

- Processing* (EMNLP '11). Association for Computational Linguistics, Stroudsburg, PA, USA, 1301-1309.
4. Cheng, N., Chandramouli, R., Subbalakshmi, K.P. 2011. Author gender identification from text. In *The International Journal of Digital Forensics & Incident Response*, Volume 8 Issue 1, July 2011, pp. 78-88.
 5. Holgrem, J., Shyu, E. 2013. Gender Classification of Facebook Posts.
 6. Kenji Kira and Larry A. Rendell. 1992. A practical approach to feature selection. In *Proceedings of the ninth international workshop on Machine learning* (ML92), Derek Sleeman and Peter Edwards (Eds.). Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, 249-256.
 7. Kobayashi, D., Matsumura, N., and Ishizuka, M. 2007. Automatic Estimation of Bloggers' Gender. In *Proceedings of International Conference on Weblogs and Social Media*. Boulder: Omnipress.
 8. Koppel, M., Argamon, S., and Shmuni, A. R. 2003. Automatically categorizing written texts by author gender. In *Literary and Linguistic Computing*, 17: 401-12
 9. Lazer, D., Pentland, A., Adamic, L., Aral, S., Barabási, A.-L., Brewer, D., Christakis, N., Contractor, N., Fowler, J., Gutmann, M., Jebara, T., King, G., Macy, M., Roy, D., and Van Alstyne, M. 2009. Computational Social Science. In *Science*, vol. 323, February 2009: 721-723.
 10. Marquardt, J., Farnadi, G., Vasudevan, G., Moens, M., Davalos, S., Teredesai, A., De Cock, M. 2014. Age and Gender Identification in Social Media. *Author Profiling Task at PAN 2014*.
 11. Mukherjee, A. and Liu, B. 2010. Improving gender classification of blog authors. In *Proceedings of the 2010 Conference on Empirical Methods in Natural Language Processing*. EMNLP'10. DOI=<http://www.aclweb.org/anthology/D10-1021>.
 12. Peersman, C., Daelemans, W., and Van Vaerenbergh, L. 2011. Predicting Age and Gender in Online Social Networks. In *Proceedings of the 3rd Workshop on Search and Mining User-Generated Contents*. SMUC'11, Glasgow, UK.
 13. Rangel, F., Rosso, P. 2013. Use of Language and Author Profiling: Identification of Gender and Age. In *Proceedings of the Tenth International Workshop on Natural Language Processing and Cognitive Science* (Marseille, France, October 2013).
 14. Sarawgi, R., Gajulapalli, K., and Choi, Y. 2011. Gender Attribution: Tracing Stylometric Evidence Beyond Topic and Genre. In *Proceedings of the Fifteenth Conference on Computational Natural Language Learning* (Portland, USA, 9 - 24 June, 2011). Association for Computational Linguistics, Stroudsburg, PA, USA, 78 - 86.
 15. Soler, J., Wanner, L. 2014. How to Use Less Features and Reach Better Performance in Author Gender Identification. In *Proceedings of LREC 2014*.
 16. Witten, I.H., & Frank, E. (2005). *Data mining: practical machine learning tools and techniques* (2nd ed, Morgan-Kaufman Series of Data Management Systems). San Francisco: Elsevier.
 17. Yan, X., Yan, L. 2006. Gender Classification of Weblog Authors. *Computational Approaches to Analyzing Weblogs*, AAAI.
 18. Zhang, C., and Zhang, P. 2010. *Predicting gender from blog posts*. Technical Report. University of Massachusetts Amherst, USA.
 19. NLTK: <http://www.nltk.org/>.
 20. WEKA: <http://www.cs.waikato.ac.nz/ml/weka/>