

# Deep Learning Distributional Features for Noise Handling in Open-set Web-genre Classification

Dimitrios Pritsos and Efstathios Stamatatos

<sup>1</sup> Dimitrios Pritsos University of the Aegean  
Karlovasi, Samos – 83200, Greece.  
dpritsos@aegean.gr

<sup>2</sup> Efstathios Stamatatos University of the Aegean  
Karlovasi, Samos – 83200, Greece.  
stamatatos@aegean.gr

**Abstract.** Web genre detection is a task that can enhance information retrieval systems by providing rich descriptions of documents and enabling more specialized queries. Most of previous studies in this field adopt the closed-set scenario where a given palette comprises all available genre labels. However this is not a realistic setup since web genres are constantly enriched with new labels and existing web genres are evolving in time. Open-set classification, where some pages used in the evaluation phase do not belong to any of the known genres, is a more realistic setup for this task. In this case, all pages not belonging to known genres can be seen as noise. This paper focuses on systematic evaluation of open-set web genre identification when the noise is either structured or unstructured. Two open-set methods combined with alternative text representation schemes and similarity measures are tested based on two benchmark corpora. Moreover, we adopt the openness test for web genre identification that enables the observation of effectiveness for a varying number of known/unknown labels.

**Keywords:** Web Genre Identification · Information Retrieval · Natural Language Processing

## 1 Introduction

## 2 Relevant Work

## 3 Distributional Features Learning

Gensim description of my methods (see HTML2VEC)

## 4 Open-set Classification Methods

### 4.1 Nearest Neighbors Distance Ratio

The Nearest Neighbors Distance Ratio (NNRD) algorithm is our variant implementation of the proposed open-set algorithm of Mendes et al. (8). In the original approach

euclidean distance has been used because of the variation of data set on which the algorithm has been evaluated. In our approach we are using cosine distance, because in text classification is being confirmed to be the proper choice in hundreds of publications. Moreover, the cosine distance is comparable to the results of the *Random Feature Subspacing Ensemble* algorithm found in (9) where cosine similarity is used for the WGI evaluation.

The NNRD algorithm is an extension of the simple *Nearest Neighbors* NN algorithm where additionally to the sets of training vectors (one set for each class) a threshold is selected by maximizing the *Normalized Accuracy* (NA) as shown in equation 1) on the *Known* and the *Marked as Unknown samples*.

$$NA = \lambda A_{KS} + (1 - \lambda) A_{MUS} \quad (1)$$

where  $A_{KS}$  is the *Known Samples Accuracy* and  $A_{MUS}$  is the *Marked as Unknown Samples Accuracy*. The balance parameters  $\lambda$  regulates the mistake trade-off on the known and marked-unknown samples prediction.

The optimally selected threshold is the the *Distance Ratio Threshold* (DRT) where NA is maximized. Equation 2 is used for calculating the Distance Ratio (DR) of the two nearest class samples, say  $s_{c_a}$  and  $u_{c_b}$ , to a random sample  $r_x$  under the constrain  $c_a c_b$ , where  $c_g$  is the sample's class.

It is very important to note that the  $c_g$  is trained in an open-set framework, therefore, the samples pairs selected for comparison might either be from the known or the marked as unknown samples. Thus  $g \in 1, 2, \dots, N$  and  $g = \emptyset$  when samples is marked as unknown.

$$DR = \frac{D(r_x, s_{c_a})}{D(r_x, s_{c_b})} \quad (2)$$

where  $D(x, y)$  is the distance between the samples where in this study is the *Cosine Distance*.

Therefore, the fitting function of the NN algorithm, described in pseudocode 1.1, is the optimization procedure to find the DRT values for classes respective sets of training samples where NA is maximized.

**Algorithm 1.1:** *Nearest Neighbor Distance Ratio* training data fitting function

---

**Data:**  $G$  the set of genre class tags  $\{1, 2, \dots, N\}$ ,  $p$  the hyper-parameter regulates the percentage of  $G$  tags will be marked as unknown,  $k$  the hyper-parameter regulates the percentage of known  $G$  tags that will be kept for validation only,  $T$  the *Distance Ratio* thresholds set than will test for finding the one which is minimizing the *Normalized Accuracy*,  $\lambda$  regulates the mistakes trade-off on the known and marked-unknown samples prediction (see eq.2),  $C[g]$  the matrix of class vector sets one for every genre class tag  $g \in G$

**Result:**  $DRT$  the *Distance Ratio Threshold* calculated by the NNRD algorithm's fitting function,  $C[g]$

- 1  $K_i^G, K_{validation}^G, U_{validation}^G, I^G = Split(G, p, k)$  splitting the  $G$  tags in to known/unknown samples combinations using the  $p$  and  $k$  hyper-parameters. The amount of split combinations is calculated by the equations 3 and 4.;
- 2  $V^G = U_{validation}^G \cup K_{validation}^G$  the validation set is the union of the  $I$  splits of the known-validation and the marked-as-unknown sets, of the whole training set;
- 3 **for each**  $i \in I$  **do**
- 4      $D_{VK}^{cos}[i] = COS_D(V_i^G, K_i^G)$  calculating all the Cosine Distances between the web-page of  $K^G$  and  $V^G$  sets for every  $I$  split combination;
- 5 **end**
- 6  $C_A^{min} = argmin(D_{VK}^{cos})$  getting the indices of the closest classes from  $V$ ;
- 7  $C_B^{min} = argmin(D_{VK}^{cos})$  getting the indices of the *second closest* classes from  $V$ ;
- 8  $R_V = D_{VK}^{cos}[C_A^{min}] / D_{VK}^{cos}[C_B^{min}]$  calculating the Distance Ratios  $R$  for all the vectors in  $V$
- 9  $NA^{max} \leftarrow 0$  initializing *Maximized Normalized Accuracy* with 0 value.  $DRT \leftarrow 0$  initializing *Distance Ratio Threshold* with 0 value.
- 10 **for each**  $drt \in T$  **do**
- 11     **for each**  $r, i \in \{R_V, count(R_V)\}$  **do**
- 12         **if**  $r < drt$  **then**
- 13              $vi = C_A^{min}[i]$  keep the respective index;
- 14              $Y[i] = G[vi]$  setting the genre's class tag as prediction for this random vector of set  $V$ ;
- 15         **else**
- 16              $Y[i] = \emptyset$  setting as none of the known genres or "I don't know";
- 17         **end**
- 18     **end**
- 19      $NA_V = NormalizedAccuracy(Y, R_V)$  calculating the *Normalized Accuracy* as shown in equation 1 for tested threshold  $drt$ ;
- 20     **if**  $NA_V > NA^{max}$  **then**
- 21          $NA^{max} \leftarrow NA_V$  keeping the maximum  $NA$  until the outer for-loop finishes;  $DRT \leftarrow drt$  keeping the *Distance Ratio Threshold* maximizes the *Normalized Accuracy*;
- 22     **else**
- 23     **end**
- 24 **end**

---

In the optimization procedure the training samples are splited based on their class tags  $c_x$ . Then some class tags are *marked as unknown* and some are left being known. Therefore, all the samples of the marked as unknown are used only in the validation subset while the known class tags samples are farther splited into the classes sets (one for each class) and into the known validation set. Then, samples of the validation sets, both then known and then marked as unknown, are used seamlessly for calculating the set of Distance Ratios (one for each class). Afterwards, a set of DRT values are tested given a range of values  $R \in t_1, t_2, t_n$  beforehand where the  $t_x$  is selected which is maximizing the NA of the validation set.

The splitting procedure the of the training set is regulated by a hyper-parameter  $p$  which defines the percentage of the class tags set  $g \in 1, 2, \dots, N$  where they will be marked as unknown. Then the total number of all possible splitting combination are calculated and these split-sets are used for finding the DRT. The combination are found using equations 3 and 4, where eq.4 is the *Binomial Coefficient*.

$$U_{num} = \text{int}(N * p) \quad (3)$$

where  $N$  is the size of the class tags set  $1, 2, \dots, N$  and  $p$  is the percentage regulation paramter for keeping the number of tags to be marked as unknown.

$$S_{num} = \frac{N!}{U_{num}!(N - U_{num})!} \quad (4)$$

The NNDR is a open-set classification algorithm, therefore, every random sample will be classified to one of the classes the NNDR has been fitted or to the unknown when its DR is greater then DRT. While training as explained above the DRT values are tested incrementally until the optimal data fitting for the training function.

In prediction phase the DRT is passed to the NNDR prediction function together with the random samples and the training samples as shown in pseudocode 1.2.

**Algorithm 1.2:** *Nearest Neighbor Distance Ratio* prediction function

---

**Data:**  $W$  the vector set of the random web-page to be classified,  $C[g]$  the matrix of class vector sets one for every genre class tag  $g \in G$ ,  $DRT$  the *Distance Ration Threshold* calculated by the NNRD algorithms fitting function

**Result:**  $Y \in \{G, \emptyset\}$ ,  $R$  the Distance Ratio scores vector, one score for every input vector of the random set  $W$

```

1 for each  $g \in G$  do
2    $D_{C_g X}^{cos} = COS_D(C[g], X)$  calculating all the Cosine Distances between the
   random web-page vectors and the class vectors of class  $g$ ;
3 end
4  $C_A^{min} = argmin(D_{C_g W}^{cos})$  getting the indices of the closest classes from  $W$ ;
5  $C_B^{min} = argmin(D_{C_g W}^{cos})$  getting the indices of the second closest classes from  $W$ ;
6  $R_W = D_{C_g W}^{cos}[D_A^{min}] / D_{C_g W}^{cos}[D_B^{min}]$  calculating the Distance Ratios  $R$  for all the
   vectors in  $W$ 
7 for each  $r, i \in \{R_W, count(R_W)\}$  do
8   if  $r < DRT$  then
9      $vi = C_A^{min}[i]$  keep the respective index;
10     $Y[i] = G[vi]$  setting the genre's class tag as prediction for this random
    vector fo set  $W$ ;
11  else
12     $Y[i] = \emptyset$  setting as none of the known genres or "I don't know";
13  end
14 end

```

---

**4.2 Random Feature Sub-spacing Ensemble**

The Random Feature Sub-spacing Ensemble (RFSE) algorithm is a variation of the method presented by Koppel et al. (5) for the task of *author identification*. In the original approach, there is only one training example for each author and a number of simple classifiers is learned based on random feature subsampling. Each classifier uses the cosine distance to estimate the most likely author. The key idea is that it is more likely for the true author to be selected by the majority of the classifiers since the used subset of features will still be able to reveal that high similarity. That is, the style of the author is captured by many different features so a subset of them will also contain enough stylistic information. Since WGI is also a style-based text categorization task, this idea should also work for it.

In our study we adopt the RFSE method as introduced in (10) shown in *Algorithm ??*. There are multiple training examples (documents) for each available genre. To maintain simplicity of classifiers, we have used a *centroid vector* for each genre. In the training phase, a centroid vector is formed, for every class, by averaging all the Term-Frequency (TF) vectors of the training examples of web pages for each genre.

The class centroids are all formed for a given feature type. Then, an evaluation document is compared against every centroid and this process is repeated  $I$  times. Every time a different feature sub-set is used. Then, the scores are ranked from highest to

lowest and we measure the number of times the document is top-matched with every class. The document is assigned to the genre with maximum number of matches given that this score exceed a predefined  $\sigma$  threshold. In the opposite case, the document remains unclassified, the RFSE responds "I Don't Know".

With respect to the similarity function, we examine cosine similarity (similar to (10)) and MinMax similarity (inspired by (6)). Moreover, in this paper we introduce a measure that combines these two similarity functions and selects the one that is most confident in each iteration. More specifically, since cosine and MinMax may have different mean and standard deviation for the set of all evaluation documents and all iterations per document, we first normalize their value. Then, for each evaluation document and each iteration we select the one with maximum normalized value. We call this similarity measure *Combo*.

## 5 Open-set Evaluation Methodology

Macro-F1 Macro-P Macro-R Macro-PR-Curves

Measuring the effect the marked-as-unknown or marked-as-noise.

Precision-Recall curve is a standard method to visualize the performance of classifiers. In this paper, the Precision-Recall curve is calculated in 11-standard recall levels  $[0, 0.1, \dots, 1.0]$ . Precision values are interpolated based on the following formula:

$$P(r_j) = \max_{r_j \leq r \leq r_{j+1}} (P(r)) \quad (5)$$

where  $P(r_j)$  is the precision at  $r_j$  standard recall level.

To compensate the potentially unbalanced distribution of web pages over the genres, we are using the macro-averaged precision and recall measures. In more detail, we use the modified version of precision and recall for open-set classification tasks proposed by (8). This modification calculates precision and recall only for the known classes (available in the training phase) while the unknown samples (belonging to classes not available during training) affect false positives and false negatives. To find parameter settings that obtain optimal evaluation performances we use 2 scalar measures, the Area Under the Precision-Recall Curve (AUC) and  $F_1$ . We will show that the appropriate selection of the optimization measure is highly significant in the presence of noise.

## 6 Experimental Setup

### 6.1 Corpora

In this paper we study the performance of the open-set classification models on the WGI task. In particular, the two open-set algorithms described above are analytically tested on benchmark corpora. In particular, our experiments are based on the following corpora already used in previous work in WGI (2; 11; 4):

1. *SANTINIS* (7): This is a corpus comprising 1,400 English web pages evenly distributed into 7 genres as well as 80 BBC web pages evenly categorized into 4 additional genres. In addition, it comprises a random selection of 1,000 English web

pages taken from the SPIRIT corpus (3). The latter can be viewed as noise in this corpus. Details are given in table 1.

2. *KI-04* (2): This is a collection of 1,205 English web pages unevenly categorized into 8 genres. Details can be seen in table 1.

## 6.2 Settings

Our text representation features are based exclusively on textual information from web pages excluding any structural information, URLs, etc. Based on the good results reported in (13; 10; 1) as well as some preliminary experiments, the following document representation schemes are examined: Character 4-grams (C4G), Word unigrams (W1G), and Word 3-grams (W3G). We use the Term-Frequency (TF) weighting scheme and the feature space is defined by a *Vocabulary* which is extracted based on the terms appearing at training set only.

As concerns OCSVM model, two parameters have to be tuned: the number of features  $F$  and  $v$ . For the former, we used  $F = \{1k, 5k, 10k, 50k, 90k\}$ , of most frequent terms of the vocabulary. Following the reports of previous studies (12) and some preliminary experiments, we examined  $v = \{0.05, 0.07, 0.1, 0.15, 0.17, 0.3, 0.5, 0.7, 0.9\}$ . In comparison to (10), this set of parameter values is more extended. With respect to RFSE, four parameters should be set: the vocabulary size  $F$ , the number of features used in each iteration  $fs$ , the number of iterations  $I$ , and the threshold  $\sigma$ . We examined  $F = \{5k, 10k, 50k, 100k\}$ ,  $fs = \{1k, 5k, 10k, 50k, 90k\}$ ,  $I = \{10, 50, 100\}$  (following the suggestion in (5) that more than 100 iterations does not improve significantly the results) and  $\sigma = \{0.5, 0.7, 0.9\}$  (based on some preliminary tests). Additionally, in this work we are testing three document similarity measures: cosine similarity, MinMax similarity, and combined cosine similarity and MinMax. Finally, to extract the best possible parameter settings for each classification method we apply grid-search over the space of all parameter value combinations.

SANTINIS		KI-04	
Genre	Pages	Genre	Pages
Blog	200	Article	127
Eshop	200	Discussion	127
FAQ	200	Download	152
Frontpage	200	Help	140
Listing	200	Link Collection	208
Personal Home Page	200	Portrayal-Non Private	179
Search Page	200	Portrayal- Private	131
DIY Mini Guide (BBC)	20	Shop	175
Editorial (BBC)	20		
Features (BBC)	20		
Short Bio (BBC)	20		
Noise (Spirit1000)	1000		

**Table 1.** Corpora descriptions and amount of pages per genre.

## 7 Experiments

### 7.1 WGI with Unstructured Noise

We initially examine the performance of OCSVM and RFSE models based on SANTINIS corpus. In the training phase, only the 11 known genres are considered. In the testing phase, the noise pages coming from the SPIRIT corpus are also used. Note that information about the true genre of these pages is not available. Therefore, we have to deal with unstructured noise. We perform 10-fold cross validation and in each fold we include the full set of 1,000 pages of noise. This evaluation strategy is giving a more realistic evaluation framework since the size of the noise is much greater than the size of any genre included in the given palette.

Figures ?? and ?? depict the Precision-Recall curves (PRC) of OCSVM and RFSE models, respectively. For each model and each one of the three document representations, the parameters that maximize performance with respect to the  $F_1$ -measure are used. Note that when recall does not reach 1.0 this means that some pages belonging to known classes were classified as unknown. In all cases, RFSE outperforms OCSVM. Moreover, for both methods, W3G seems to be the best feature type for this corpus, followed by C4G. OCSVM performance is only comparable with RFSE when W3G is used.

We further explore the performance of the open-set WGI methods by selecting parameter settings with different optimization criteria. Tables ?? and ?? show the combination of parameters that optimize performance of OCSVM and RFSE based on AUC,  $F_1$  and  $F_{0.5}$ . Moreover, in the tables we show the values of all three performance measures where one of them is maximized. It is clear that the performance in all cases is maximized when W3G document representation is used. In previous studies based on a closed-set framework, C4G was the document type of features to maximize performance (14). This indicates that contextual and content information is important for this corpus (1).

In addition, in almost all cases, RFSE models are far more effective than OCSVM. Another important conclusion is that the optimization criterion plays a crucial role for the properties of the model especially for RFSE. When AUC is maximized, recall is favoured. On the other hand, while  $F_1$  is maximized, precision is substantially increased. Fig. ?? shows the performance of OCSVM and RFSE models when AUC and  $F_1$  criteria are used to select parameter settings. As can be seen, the RFSE model based on  $F_1$  maximization avoids to make wrong decisions and leaves a large number of web pages unclassified. On the other hand, the model optimized by AUC prefers to make a lot of errors in order to recognize more web pages of known genres. OCSVM models seem not significantly affected. Note that choosing between WGI models that prefers precision over recall and vice versa is an application-specific task.

As it was highlighted in the previous section, according to the properties of the application in which WGI is involved, precision may be more important than recall or vice-versa. In figure ?? the macro-precision of RFSE is depicted for W3G, W1G and C4G features. MinMax similarity is used since it increases significantly the perfor-



mance of RFSE in respect with precision. As concerns text representation, W1G is the best choice when precision is at more importance than recall. On the other hand, W3G features seem to be more stable because the standard error is lower than that of the other features and also the W3G model is not affected too much when openness surpasses 0.5 (actually it improves).

In the case of C4G and W1G where the openness level is 0.646 the standard error in both case is very high. Since, we observe this problem only in the case where the problems has been reduced to binary, we are interested to see whether it is caused by choice of the document representation or by the choice of the similarity measure.

Despite OCSVM's improvement when structured noise is used, it can only be competitive to RFSE on a high openness level, where all genre labels but one are considered unknown. This can be better viewed in figure ?? where OCSVM is compared with RFSE models based on MinMax and Combo similarity measures for a varying openness level. These curves correspond to W1G features, so they are not the optimal models. However, they provide a fair comparison between examined methods. As standard error bars indicate, the performance of RFSE models with respect to the  $F_1$  measure is significantly better than that of OCSVM while openness is less than 0.5. Beyond that level, OCSVM is significantly better than RFSE models. Note also Combo measure helps RFSE in while openness is relatively low and MinMax seems to be a better choice when openness increases.

In this paper we presented an experimental study on WGI focusing on open-set evaluation for this task. In contrast to vast majority of previous work in this area, we adopt the open-set scenario that is more realistic for WGI since it is not feasible to construct a genre palette with all available genres and appropriate samples for each one of them. Moreover, we examined two open-set classification methods and several feature types and similarity measures. To the best of our knowledge, this is the first time the performance of WGI models is evaluated using performance measures and tests specifically designed for open-set classification tasks.

The presented evaluation of open-set WGI covers two basic scenarios. The first is when noise is unstructured, i.e., information about the true genre of pages not belonging to the known genre palette is not available. The second scenario applies when noise is structured, i.e., we actually know the true genre of pages not included in the training classes. For both cases, we propose appropriate evaluation methodologies and present comparative results for the tested models.

In almost all examined cases, RFSE models outperformed the corresponding OCSVM models. This verifies previous work findings about the appropriateness of RFSE for WGI (10). RFSE is able to provide effective models and additionally it is possible to manage preference on recall or precision, an application-dependent choice, by focusing on optimizing AUC or  $F_1$  respectively. On the other hand, OCSVM proved to be the best-performing method in extreme cases when openness is high. Actually, the restrictions of the available corpora did not allow us to examine cases where openness approaches 1.0. However, it seems that when openness is more than 0.5 OCSVM outperforms RFSE.

As concerns the feature types, in most of the cases W3G and C4G provided the best results. However, the selection of text representation features is a crucial choice that

affects performance and it seems to be corpus-dependent. Another crucial parameter of RFSE is the similarity measure. Among the examined measures, MinMax and its combination with cosine similarity provide the most robust results. The choice of similarity measure correlates with feature types. It seems that the combo measure is more effective than MinMax in low openness conditions.

To enhance the evaluation of WGI models in open-set conditions, we need larger corpora including multiple genre labels. New enhanced open-set WGI methods are needed and they should be evaluated using the proposed paradigm. Otherwise, using an evaluation paradigm more appropriate for closed-set tasks, the performance may be over-estimated.

## Bibliography

- [1] Asheghi, N.R.: Human Annotation and Automatic Detection of Web Genres. Ph.D. thesis, University of Leeds (2015)
- [2] Meyer zu Eissen, S., Stein, B.: Genre classification of web pages. *KI 2004: Advances in Artificial Intelligence* pp. 256–269 (2004)
- [3] Joho, H., Sanderson, M.: The spirit collection: an overview of a large web collection. In: *ACM SIGIR Forum*. vol. 38, pp. 57–61. ACM (2004)
- [4] Kanaris, I., Stamatatos, E.: Learning to recognize webpage genres. *Information Processing & Management* **45**(5), 499–512 (2009)
- [5] Koppel, M., Schler, J., Argamon, S.: Authorship attribution in the wild. *Language Resources and Evaluation* **45**(1), 83–94 (2011)
- [6] Koppel, M., Winter, Y.: Determining if two documents are written by the same author. *Journal of the Association for Information Science and Technology* **65**(1), 178–187 (2014)
- [7] Mehler, A., Sharoff, S., Santini, M.: *Genres on the Web: Computational Models and Empirical Studies*. Text, Speech and Language Technology, Springer (2010)
- [8] Mendes Júnior, P.R., de Souza, R.M., Werneck, R.d.O., Stein, B.V., Pazinato, D.V., de Almeida, W.R., Penatti, O.A., Torres, R.d.S., Rocha, A.: Nearest neighbors distance ratio open-set classifier. *Machine Learning* pp. 1–28 (2016)
- [9] Pritsos, D., Stamatatos, E.: Open set evaluation of web genre identification. *Language Resources and Evaluation* pp. 1–20 (2018)
- [10] Pritsos, D.A., Stamatatos, E.: Open-set classification for automated genre identification. In: *Advances in Information Retrieval*, pp. 207–217. Springer (2013)
- [11] Santini, M.: Automatic identification of genre in web pages. Ph.D. thesis, University of Brighton (2007)
- [12] Scholkopf, B., Platt, J., Shawe-Taylor, J., Smola, A., Williamson, R.: Estimating the support of a high-dimensional distribution. Technical Report MSR-TR-99-87 (1999)
- [13] Sharoff, S., Wu, Z., Markert, K.: The web library of babel: evaluating genre collections. In: *Proceedings of the Seventh Conference on International Language Resources and Evaluation*. pp. 3063–3070 (2010)
- [14] Sharoff, S., Wu, Z., Markert, K.: The web library of babel: evaluating genre collections. In: *LREC*. Citeseer (2010)