

Automatic Classification by Topic Domain for Meta Data Generation, Web Corpus Evaluation, and Corpus Comparison

Roland Schäfer

Freie Universität Berlin
Habelschwerdter Allee 45
14196 Berlin, Germany

roland.schaefer@fu-berlin.de

Felix Bildhauer

Institut für Deutsche Sprache
R5, 6–13
68161 Mannheim, Germany

bildhauer@ids-mannheim.de

Abstract

In this paper, we describe preliminary results from an ongoing experiment wherein we classify two large unstructured text corpora—a web corpus and a newspaper corpus—by topic domain (or subject area). Our primary goal is to develop a method that allows for the reliable annotation of large crawled web corpora with meta data required by many corpus linguists. We are especially interested in designing an annotation scheme whose categories are both intuitively interpretable by linguists and firmly rooted in the distribution of lexical material in the documents. Since we use data from a web corpus and a more traditional corpus, we also contribute to the important field of corpus comparison and corpus evaluation. Technically, we use (unsupervised) topic modeling to automatically induce topic distributions over gold standard corpora that were manually annotated for 13 coarse-grained topic domains. In a second step, we apply supervised machine learning to learn the manually annotated topic domains using the previously induced topics as features. We achieve around 70% accuracy in 10-fold cross validations. An analysis of the errors clearly indicates, however, that a revised classification scheme and larger gold standard corpora will likely lead to a substantial increase in accuracy.

1 Introduction

In the experiment reported here, we classified large unstructured text corpora by *topic domain*. The *topic domain* of a document—along with other high-level classifications such as *genre* or

register—is among the types of meta data most essential to many corpus linguists. Therefore, the lack of reliable meta data in general is often mentioned as a major drawback of large, crawled web corpora, and the automatic generation of such meta data is an active field of research.¹ It must be noted, however, that such high-level annotations are not reliably available for many very large traditional corpora (such as newspaper corpora), either. When it comes to the automatic identification of high-level categories like *register* (such as *Opinion*, *Narrative*, *Informational Persuasion*; Biber and Egbert 2016), even very recent approaches based on very large amounts of training data cannot deliver satisfying (arguably not even encouraging) results. For instance, Biber and Egbert (2016, 23) report *accuracy*=0.421, *precision*=0.268, *recall*=0.3. It is not even clear whether categories such as *register* and *genre* can be operationalized such that a reliable annotation is possible for humans.

By contrast, automatic text categorization by *content* yielded much more promising results years ago already (Sebastiani, 2002). Furthermore, data-driven induction of topics (*topic modeling*) has proven quite successful, and it is in many respects a very objective way of organizing a collection of documents by content. Deriving topic classifications from text-internal criteria is also advocated in the EAGLES (1996) guidelines, among others. However, topic modeling usually does not come with category labels that are useful for linguistic corpus users. In our project, we explore the possibility of inferring a small, more traditional set of *topic domains* (or *subject areas*) from the topics induced in an unsupervised manner by Latent Semantic Indexing (Landauer and Dumais, 1994; Landauer and Dumais, 1997).

¹See, for example, many of the contributions in Mehler et al. (2010).

Since we classify and compare one large German web corpus and one large German newspaper corpus with respect to their distribution of topic domains, our paper also contributes to the area of corpus comparison, another important issue in corpus linguistics (Kilgarriff, 2001; Biemann et al., 2013). For the construction of crawled web corpora, such comparisons are vital because next to nothing is known about their composition.

The computational tools used in our method (unsupervised topic induction and supervised classifiers) are by now well-established and highly developed. This paper contributes to the field of applying such methods and making them usable for real-life problems of data processing and the development of suitable annotation schemes rather than to the development of the underlying mathematics and algorithms.

2 Gold Standard Data

Our gold standard corpora were prepared by manual annotation of documents from two large German corpora. The first data set consists of 870 randomly selected documents from DECOW14A, a crawled web corpus (Schäfer and Bildhauer, 2012; Schäfer, 2015), henceforth *Web*. The second data set contains 886 documents randomly selected from DeReKo, a corpus composed predominantly of newspaper texts (Kupietz et al., 2010), henceforth *News*. Our choice of corpora was motivated by fact that we expected some overlap w. r. t. to topics covered in them, but also some major differences. The documents in these gold standard corpora were classified according to a custom annotation scheme for topic domain which builds on previous work by Sharoff (2006). The design goal was to have moderate number (about 10–20) of topic domains that can be thought of as subsuming more fine-grained topic distinctions. We developed the annotation scheme in a cyclic fashion, taking into account annotator feedback after repeated annotation processes. For the experiment reported here, we used a version that distinguishes 13 topic domains, namely *Science, Technology, Medical, Public Life and Infrastructure, Politics and Society, History, Business, Law, Fine Arts, Philosophy, Beliefs, Life and Leisure, Individuals*.

3 Experiment Setup

Our general approach was to infer a topic distribution over a corpus using *unsupervised* topic mod-

eling algorithms as a first step. In the second step, rather than examining and interpreting the inferred topical structure, we used the resulting document–topic matrix to learn topic domain distinctions for the documents from their assignment to the topics in a *supervised* manner. To achieve this, supervised classifiers were used. Through permutation of virtually all available classifiers (with the appropriate capabilities) available in the Weka toolkit (Hall and Witten, 2011), LM Trees (Landwehr et al., 2005) and SVMs with a Pearson VII kernel (Üstün et al., 2006) were found to be most accurate. Due to minimally higher accuracy, SVMs were used in all subsequent experiments. Some topic domains occurred only rarely in the gold standard, and we did not expect the classifier to be able to generalize well from just a few instances. Therefore, we evaluated the results on the *full* data set and a *reduced* data set with rare categories removed.

For the first step (unsupervised topic induction), we used LSI and LDA (Latent Dirichlet Allocation, Blei et al. 2003) as implemented in the Gensim toolkit (Řehůřek and Sojka, 2010). In our first experiments, the LDA topic distribution was unstable, and results were generally unusable, possibly due to the comparatively small gold standard corpora used. We consequently only report LSI results here and will return to LDA in further experiments (cf. Section 5). However, for any topic modeling algorithm, our corpora can be considered small. Therefore, we inferred topics not just based on the annotated gold standard data sets, but also on larger datasets which consisted of the gold standard mixed with additional documents from the source corpora. For the training of the SVM classifiers, the documents that had been mixed in were removed again because no gold standard annotation was available for them. We systematically increased the number of mixed-in document in increments of roughly half as many documents as contained in the gold standard corpora.

We pre-processed both corpora in exactly the same way (tokenization, lemmatization, POS-tagging, named entity recognition). Using the lemma and the simplified POS tags (such as *kindergarten_nn*) as terms in combination with some filters (use only lower-cased purely alphabetic common and proper noun lemmas between 4 and 30 characters long) usually gave the best results.

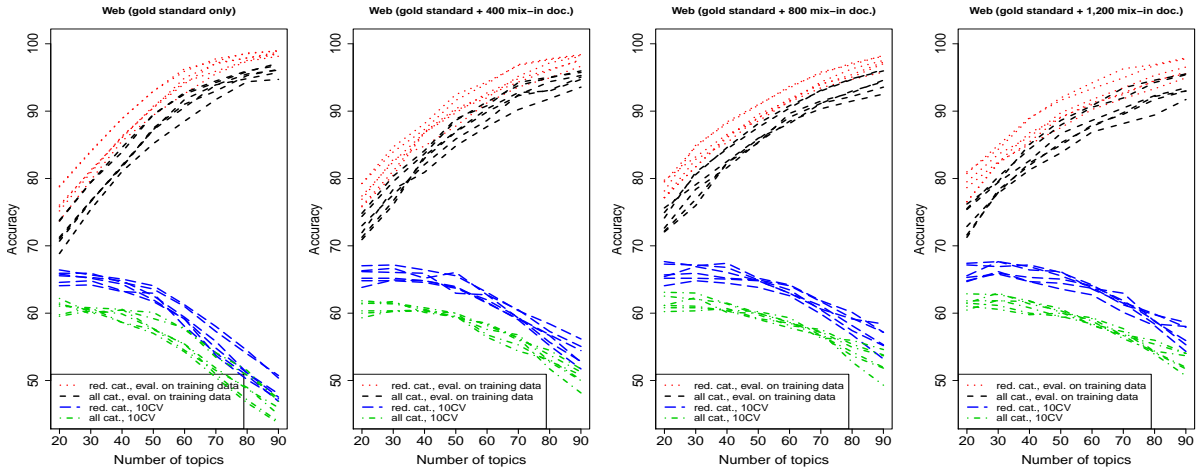


Figure 1: Accuracy with different numbers of topics for the Web dataset

4 Results

Figure 1 shows the classification accuracy using 20 to 90 LSI topics. Each line corresponds to one sub-experiment (with slightly different preprocessing options for lines of the same color and style), and the lines form well distinguishable bands. The highest accuracy is achieved with the reduced set of topic domains (minor categories removed) when the evaluation is performed on the training data. The full set of topic domains leads to a drop in accuracy of about 5%. The two lower bands show the classification accuracy in a 10-fold cross-validation (10CV), again with the reduced set of topic domains performing roughly 5% better. While a higher number of topics improves results on the training data, the accuracy in the cross-validation drops. Too large numbers of topics obviously allow the method to pick up idiosyncratic features of single documents or very small clusters of documents, leading to extreme overfitting.

The four panels show results based on different topic models. Panel (a) uses a topic model inferred only from the (more than 800) gold standard documents. Results in panel (b) through (d) are based on topic models inferred on larger data sets as described in Section 3. In the experiment reported in panel (d), for example, 1,200 documents were added to the 870 gold standard documents. While the results of the 10CV are slightly improved by mixing in more documents, the maximum achieved accuracy does not change significantly. We mixed in up to 8,000 additional documents (not all results shown here) with no significant change compared to panel (d) in Figure 1.

We consider the maximum 10CV accuracy with the reduced set of topic domains most informative w. r. t. the potential quality of our method, and we report it in Table 1.

A very similar plot for the News data is shown in Figure 3. The best results are also given in Table 1. The added accuracy (4.23% according to Table 1) is a side effect of the more skewed distribution of topic domains in the News gold standard data.

The κ statistic for the Web and Newspaper results from Table 1 is $\kappa_{\text{Web}} = 0.575$ and $\kappa_{\text{News}} = 0.582$, indicating that achieving a higher accuracy for the web data is actually slightly harder than for the newspaper data (see also the analysis of the confusion matrices below).

When the Web and News data are pooled, however, quality drops below any acceptable level, cf. Figure 3 and Table 1. Mixing in more documents (panels b–d) improves the evaluation results on the training data, but the 10CV results remains steady at around 50%. This is remarkable because larger training data sets should lead to increased, not degraded accuracy. While a deeper analysis of the LSI topic distributions remains to be undertaken, it is evident what most likely causes these below average results on the side of the SVM classifier when looking at the confusion matrices, cf. Table 2. In the Web gold standard (panel a), the dominant modal category is *Life and Leisure*. The distribution of topic domains is reasonably skewed, and the confusion is distributed roughly uniformly across categories. The News gold standard (panel b) consists mainly of two clusters of documents in the domains *Politics and Society* and *Life and Leisure*. For the pooled data set (panel c), this

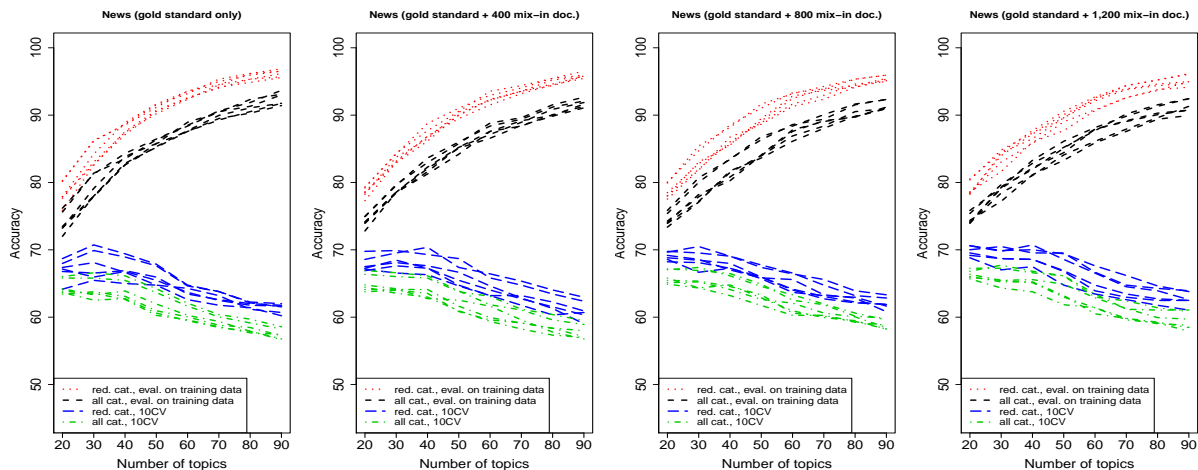


Figure 2: Accuracy with different numbers of topics for the News dataset

Corpus	Mixed-in	Attribute	Topics	Accuracy	Precision*	Recall*	F-Measure*
Web	3,200	token	20	68.765%	0.688	0.688	0.674
News	3,600	lemma + POS	40	72.999%	0.725	0.730	0.696
Web + News	0	lemma + POS	30	51.872%	0.431	0.519	0.417

Table 1: Evaluation at best achievable accuracy with the reduced set of topic domains in 10-fold cross-validation (*weighted average across all categories)

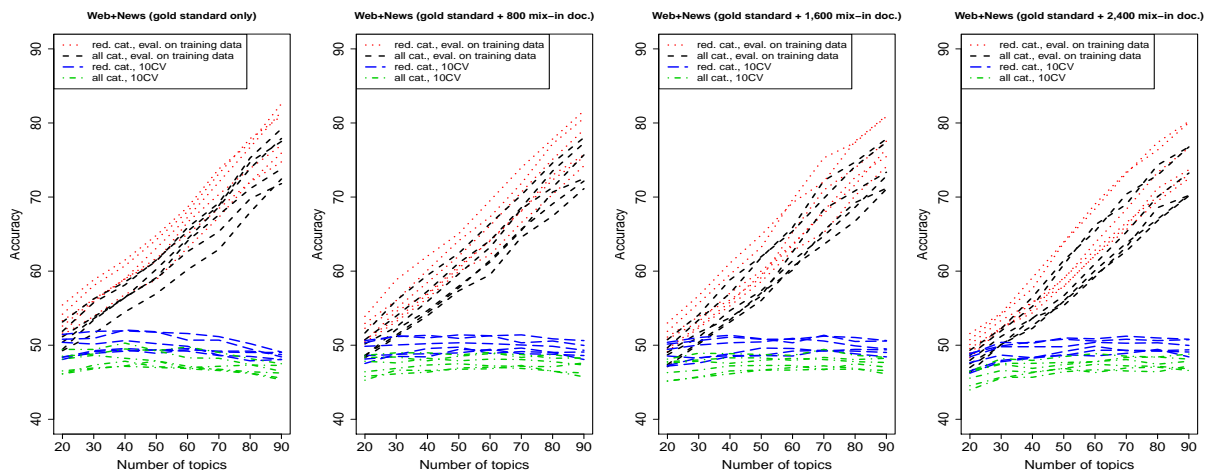


Figure 3: Accuracy with different numbers of topics for the pooled Web + News datasets

leads to a situation in which the classifier simply assigns most documents to *Life and Leisure* and the rest mostly to *Politics and Society*. This indicates that for such skewed distributions of topic domains, larger gold standard data sets are required. It is not indicative of a general failure of the method or a general incompatibility of newspaper and web data in the context of our method. The confusion matrices in Table 2 clearly indicate, however, that topic domains are represented quite differently in newspaper and web corpora.

5 Conclusions and Outlook

The results presented here are highly encouraging, and they clearly indicate the route to be taken in further experiments. First of all, there appears to be a connection between induced topic distributions and more general topic domains. The decreased performance in cross-validation experiments indicates that larger gold standard data sets are required. Such data sets are currently being annotated under our supervision.

Web		Classified							
		PolSoc	Busi	Life	Arts	Public	Law	Beliefs	Hist
Annotated	PolSoc	26	12	10	1	1	0	1	0
	Busi	5	105	40	7	1	2	1	1
	Life	3	14	286	6	4	1	1	1
	Arts	3	2	36	78	1	0	2	6
	Public	0	3	11	0	9	1	0	0
	Law	3	9	8	0	1	8	0	0
	Beliefs	4	3	11	6	1	0	30	1
	Hist	9	0	9	7	1	1	2	15

News		Classified					
		PolSoc	Busi	Life	Indiv	Arts	Public
Annotated	PolSoc	223	6	39	0	0	8
	Busi	20	24	9	0	0	0
	Life	24	1	324	0	0	1
	Indiv	5	0	17	0	0	1
	Arts	2	0	28	0	6	0
	Public	35	0	30	0	0	34

Pooled		Classified								
		PolSoc	Busi	Medical	Life	Arts	Public	Law	Beliefs	Hist
Annotated	PolSoc	199	7	0	109	0	12	0	0	0
	Busi	18	23	0	172	0	2	0	0	0
	Medical	6	0	0	29	0	1	0	0	0
	Life	25	4	0	632	0	5	0	0	0
	Arts	2	2	0	160	0	0	0	0	0
	Public	46	2	0	56	0	19	0	0	0
	Law	8	0	0	31	0	0	0	0	0
	Beliefs	0	0	0	59	0	0	0	0	0
	Hist	4	0	0	50	0	0	0	0	0

Table 2: Confusion matrices for the best achievable results on the Web (a), News (b), and pooled (c) data sets as reported in Table 1; different sets of categories are the result of excluding low-frequency topic domains (below 20 for Web and News, below 30 for pooled data)

Secondly, there appears to be a significant difference in the topic distribution and the topic/domain mapping in newspaper and web corpora. This might be one of the reasons behind the collapse of the classifier when newspaper and web data are pooled. In future experiments, it remains to be discovered whether larger gold standard corpora can alleviate such problems. This will eventually enable us to decide whether separate models or pooled models for the two kinds of corpora are more appropriate.

Thirdly, the highly skewed topic distributions in both newspaper and web corpora indicate that splitting up some topic domains might lead to a better fit. In fact, annotators have independently asked whether *Politics and Society* and *Life and Leisure*—the critical categories which make the classifier collapse (cf. Section 4)—could not be split up into at least two categories each.

Additionally, we will investigate whether alternative topic modeling algorithms lead to a better fit.² Moreover, as suggested by an anonymous reviewer, our results could be compared with a baseline classification that does not make use of topic modeling algorithms. Finally, we are currently experimenting with an extended annotation scheme that allows for multiple weighted assignments of documents to topic domains.

The ultimate goal of our project is to automatically annotate existing web corpora that are several billion tokens large with meta data such as their topic domain and to release the data freely (under a maximally permissive Creative Commons Attribution license).³ The experiments re-

²The Gensim toolkit offers a wide array of algorithms, including *doc2vec* and an alternative LDA implementation *ldamallet*.

³The software and the classifiers will be made available under permissive open source licenses allowing even their use in commercial applications.

ported here indicate that with some tweaking, it will be possible to create such free resources and achieve very high levels of quality.

Acknowledgments

Roland Schäfer’s research on the project presented here was funded by the German Research Council (Deutsche Forschungsgemeinschaft, DFG) through grant SHA/1916-1 *Linguistic Web Characterization*.

References

- Douglas Biber and Jesse Egbert. 2016. Using grammatical features for automatic register identification in an unrestricted corpus of documents from the open web. *Journal of Research Design and Statistics in Linguistics and Communication Science*, 2:3–36.
- Chris Biemann, Felix Bildhauer, Stefan Evert, Dirk Goldhahn, Uwe Quasthoff, Roland Schäfer, Johannes Simon, Leonard Swiezinski, and Torsten Zesch. 2013. Scalable construction of high-quality web corpora. *Journal for Language Technology and Computational Linguistics*, 28(2):23–60.
- David M. Blei, Andrew Y. Ng, and Michael I. Jordan. 2003. Latent dirichlet allocation. *Journal of Machine Learning Research*, 3:993–1022.
- EAGLES. 1996. Preliminary recommendations on text typology. Technical report EAG-TCWG-TTYP/P, EAGLES.
- Mark Hall and Ian H. Witten. 2011. *Data mining: practical machine learning tools and techniques*. Kaufmann, Burlington, 3rd edition.
- Adam Kilgarriff. 2001. Comparing corpora. *International Journal of Corpus Linguistics*, 6(1):97–133.
- Marc Kupietz, Cyril Belica, Holger Keibel, and Andreas Witt. 2010. The German reference corpus DeReKo: A primordial sample for linguistic research. In Nicoletta Calzolari, Khalid Choukri, Bente Maegaard, Joseph Mariani, Jan Odijk, Stelios

- Piperidis, Mike Rosner, and Daniel Tapias, editors, *Proceedings of the Seventh International Conference on Language Resources and Evaluation (LREC '10)*, pages 1848–1854, Valletta, Malta. European Language Resources Association (ELRA).
- Thomas K. Landauer and Susan T. Dumais. 1994. Latent semantic analysis and the measurement of knowledge. In R. M. Kaplan and J. C. Burstein, editors, *Educational testing service conference on natural language processing techniques and technology in assessment and education*. Educational Testing Service, Princeton, NJ.
- Thomas K. Landauer and Susan T. Dumais. 1997. A solution to plato’s problem: the latent semantic analysis theory of acquisition, induction and representation of knowledge. *Psychological Review*, 104(2):211–240.
- Niels Landwehr, Mark Hall, and Eibe Frank. 2005. Logistic model trees. *Machine Learning*, 95(1–2):161–205.
- Alexander Mehler, Serge Sharoff, and Marina Santini, editors. 2010. *Genres on the Web. Computational Models and Empirical Studies*, volume 42 of *Text, Speech and Language Technology*. Springer, Dordrecht.
- Radim Řehůřek and Petr Sojka. 2010. Software Framework for Topic Modelling with Large Corpora. In *Proceedings of the LREC 2010 Workshop on New Challenges for NLP Frameworks*, pages 45–50, Valletta, Malta. ELRA.
- Roland Schäfer and Felix Bildhauer. 2012. Building large corpora from the web using a new efficient tool chain. In Nicoletta Calzolari, Khalid Choukri, Thierry Declerck, Mehmet Uğur Doğan, Bente Maegaard, Joseph Mariani, Jan Odijk, and Stelios Piperidis, editors, *Proceedings of the Eight International Conference on Language Resources and Evaluation (LREC'12)*, pages 486–493, Istanbul. ELRA.
- Roland Schäfer. 2015. Processing and querying large web corpora with the COW14 architecture. In Piotr Bański, Hanno Biber, Evelyn Breiteneder, Marc Kupietz, Harald Lungen, and Andreas Witt, editors, *Proceedings of Challenges in the Management of Large Corpora 3 (CMLC-3)*, Lancaster. UCREL.
- Fabrizio Sebastiani. 2002. Machine learning in automated text categorization. *ACM Computing Surveys*, 34(1):1–47.
- Serge Sharoff. 2006. Creating general-purpose corpora using automated search engine queries. In Marco Baroni and Silvia Bernardini, editors, *Wacky! Working papers on the Web as Corpus*. GEDIT, Bologna.
- Bülent Üstün, Willem J. Melssen, and Lutgarde M.C. Buydens. 2006. Facilitating the application of Support Vector Regression by using a universal Pearson VII function based kernel. *Chemometrics and Intelligent Laboratory Systems*, 81:29–40.