
PyTAIL: Interactive and Incremental Learning of NLP Models with Human in the Loop for Online Data

Shubhanshu Mishra*

shubhanshu.com
mishra@shubhanshu.com

Jana Diesner

University of Illinois at Urbana-Champaign
jdiesner@illinois.edu

Abstract

Online data streams make training machine learning models hard because of distribution shift and new patterns emerging over time. For natural language processing (NLP) tasks that utilize a collection of features based on lexicons and rules, it is important to adapt these features to the changing data. To address this challenge we introduce PyTAIL, a python library, which allows a human in the loop approach to actively train NLP models. PyTAIL enhances generic active learning, which only suggests new instances to label by also suggesting new features like rules and lexicons to label. Furthermore, PyTAIL is flexible enough for users to accept, reject, or update rules and lexicons as the model is being trained. Finally, we simulate the performance of PyTAIL on existing social media benchmark datasets for text classification. We compare various active learning strategies on these benchmarks. The model closes the gap with as few as 10% of the training data. Finally, we also highlight the importance of tracking evaluation metric on remaining data (which is not yet merged with active learning) alongside the test dataset. This highlights the effectiveness of the model in accurately annotating the remaining dataset, which is especially suitable for batch processing of large unlabelled corpora. PyTAIL will be available at <https://github.com/socialmediaie/pytail>.

1 Introduction

Analysis of large scale natural language corpora often requires annotation of dataset in a given domain with pre-trained models. Generally, these models are pre-trained on a fixed training dataset which is often different from the domain of the dataset under consideration. This often leads to poor performance of the model on this new domain. One way to address this gap is to utilize domain adaptation [Sarawagi, 2008, Daumé III, 2007] to improve the model accuracy. However, efficient domain adaptation requires labeled training data from the new domain, which is costly to acquire. The problem gets compounded for social media data, for which the vocabulary and language usage continuously evolve over time. Take the example of sentiment classification, where the ways of expressing the same opinion also change with time. For example, the opinion label of the phrase “you are just like *subject*”, will depend on the general opinion about “*subject*” when the phrase was expressed. Similarly, many new words are coined on social media [Eisenstein, 2013, Gupta et al., 2010]. This poses a challenge for maintaining these models retain their accuracy over time. In this work, we propose an approach to alleviate this issue by creating a system based on active human-in-the-loop learning which incrementally updates an existing classifier by requiring a user to provide few new examples from the new data. Traditionally, this setup, called active learning [Settles, 2009] only deals with suggesting new training examples to annotate. However, since many NLP models consume feature based on existing rules or lexicons, with changing data characteristics it

*Work done while at University of Illinois at Urbana-Champaign.

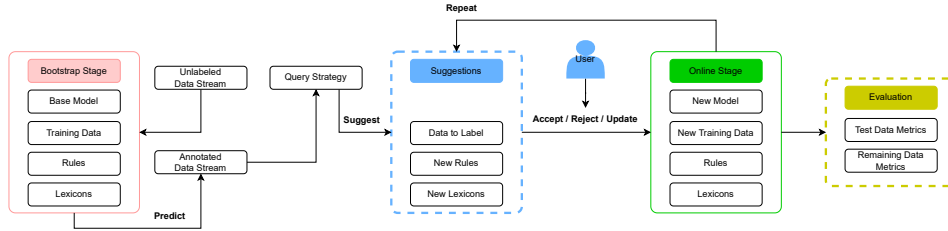


Figure 1: **PyTAIL Workflow**: Given a user and an unlabeled data stream, along with some bootstrapping artifacts, PyTAIL suggests data instances, rules, and lexicons which can be merged with bootstrapping artifacts to continuously create new model.

may be more desirable to also suggest rule and lexicon updates in the model. Our system PyTAIL (Python Text Analysis and Incremental Learning) addresses the issues highlighted here by allowing human-in-the-loop active learning systems to integrate new data points, rules, and lexicons. Our main contributions are as follows: (i) Introduce PyTAIL, an open source tool with an active learning workflow which uses new data, rules, and lexicons to continuously train NLP models. (ii) Introduce a social media text classification benchmark for active learning research. (iii) Introduce an evaluation setup on unconsumed data in active learning to quantify how quickly a corpus can be fully annotated with a reasonable accuracy.

2 Incremental learning of models with human in the loop

In this section we describe PyTAIL (Python Text Analysis and Incremental Learning). PyTAIL’s goal is to enable efficient construction of training data using active learning, while supporting incremental learning of models using the most recent data. A description of PyTAIL workflow is shown in figure 1. PyTAIL is built with the following features in mind: (i) Low cost of continuous training data acquisition (ii) Incorporation of domain knowledge using lexicon and rules (iii) Efficient update of model using only the newly acquired training data.

Overview As shown in figure 1, the user starts with a collection of artifacts in the Bootstrap Stage. This can include a pre-trained model, a small seed training dataset, existing rules, and lexicons. Next, the user introduces their unlabeled data stream from their domain of interest, e.g. social media corpora. The bootstrap artifacts are used to predict this data stream. These predictions are then fed to the query strategy (described below) to identify artifacts for the suggestion stage. The user can then accept, reject, update these suggestions or even introduce new suggestions. Next, the model is updated using updated artifacts such that the rules and lexicons are used for updating the model features and the annotated data is used for updating the model. Finally, PyTAIL shows continuous evaluation metrics which include metric on a test set, user accepted training set, and unobserved data stream. This process is repeated till a stopping criteria is met, e.g. the exhaustion of data stream or achieving reasonable evaluation score. PyTAIL supports two modes for training, one is human in the loop (HITL) mode, and another is simulation mode. The simulation model uses pre-defined heuristics to simulate human actions based on model prediction scores. The default model when applied to benchmark datasets is the simulation mode.

Human in the loop (HITL) mode In the HITL mode, PyTAIL uses the pre-trained model to suggest top K instances to the user. The user can sort the instances using the scoring criterion. In order to reduce the cognitive work of labeling an instance from scratch, the user is shown the model predictions (as well as the label probability). The user is only required to edit the labels if they disagree. Model predictions for all the unlabeled instances from the top suggestions are now used as gold labels and fed to the model during the update process (this is similar to self-supervision with the possibility of human intervention). The user is also shown the prominent features for that instance, the user can select these features and mark them as useful or useless. Lexicon matches with the annotations are also shown, along with prominent key phrases in the unlabeled data stream. The user can choose to update the lexicon with these new suggestions. Once the model update has happened, the user is provided feedback on the change in model evaluation on a held out data.

3 PyTAIL for Social Media Text Classification

Benchmark for social media active learning We introduce an active learning benchmark of 10 social media text classification datasets consisting of 200K posts. These datasets cover sentiment classification, abusive content identification, and uncertainty indicator. Details in appendix section A.

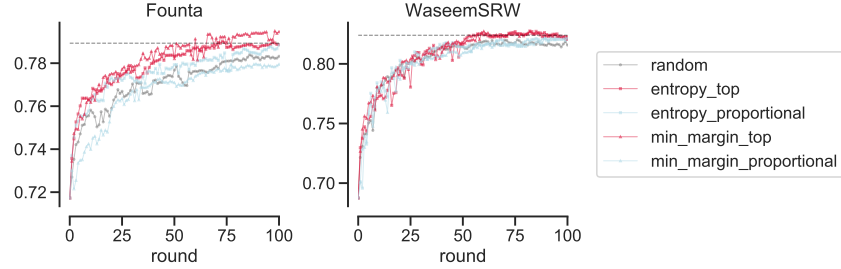
Model We use a logistic regression model with L_2 regularization. The regularization parameter is selected for each model using cross validation. We track the model scores on the held out test as well as validation data. Each text is represented using a set of features as described in section B.

Query selection strategies Active learning algorithms [Settles, 2009] identify most informative instances from unlabeled data that can be used to construct a high quality training dataset. The process of identifying informative instances is called **query selection**. Top instances $X_{selected}$ from the unlabeled data $X_{unlabeled}$ are identified based on a score. We consider two types of score: (i) *entropy* = $\sum_i p_i * \log(p_i)$ - higher is better (ii) *min-margin* = $\max_{i \neq *}\{p_i - p_* \mid p_* = \max_j p_j\}$ - lower is better. The entropy based scoring favors model predictions with highest randomness. The min-margin based scoring is useful in ensuring that the difference between the top prediction score and the second top prediction score is less. The selection is done using three strategies: (i) *Rand*: Instances are selected randomly without considering their scores, this acts as a baseline. (ii) X_{top} : Top K instances are selected based on their scores (X). (iii) X_{prop} : K instances are sampled proportional to their scores (X). This adds a degree of randomness to the top k strategy. These new instances are then added to the existing training instances $X_{train} = X_{train} \cup X_{selected}$, and the model is retrained.

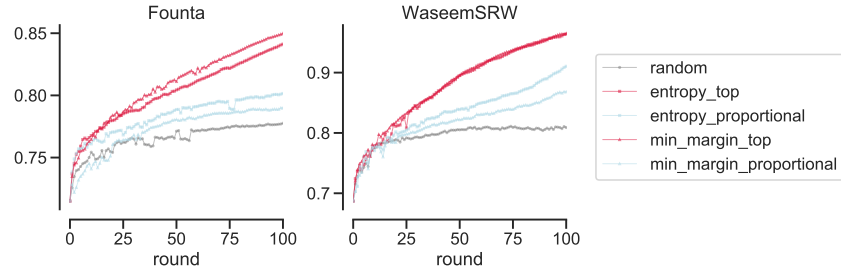
Table 1: Performance of query strategies across datasets using around 10% training dataset.

task	dataset	round	N	N_{left}	$\%used$	Full	Rand	E_{top}	E_{prop}	M_{top}	M_{prop}
Test Dataset											
ABUSIVE	Founta	42	41,861	37,661	0.10	0.79	0.77	0.78	0.78	0.79	0.77
	WaseemSRW	14	13,072	11,672	0.11	0.82	0.79	0.78	0.77	0.78	0.76
SENTIMENT	Airline	9	8,725	7,825	0.10	0.82	0.76	0.78	0.79	0.77	0.77
	Clarin	45	44,299	39,799	0.10	0.66	0.63	0.61	0.62	0.63	0.63
	GOP	8	7,121	6,321	0.11	0.67	0.63	0.64	0.63	0.62	0.64
	Healthcare	1	590	490	0.17	0.59	0.64	0.60	0.61	0.60	0.60
	Obama	2	1,777	1,577	0.11	0.63	0.56	0.60	0.58	0.59	0.57
	SemEval	13	12,145	10,845	0.11	0.65	0.59	0.60	0.61	0.58	0.61
UNCERTAINTY	Riloff	2	1,201	1,001	0.17	0.78	0.77	0.76	0.77	0.76	0.79
	Swamy	1	555	455	0.18	0.39	0.39	0.40	0.39	0.34	0.31
Remaining Dataset											
ABUSIVE	Founta	42	41,861	37,661	0.10	NaN	0.77	0.80	0.78	0.81	0.78
	WaseemSRW	14	13,072	11,672	0.11	NaN	0.78	0.79	0.77	0.80	0.76
SENTIMENT	Airline	9	8,725	7,825	0.10	NaN	0.75	0.79	0.79	0.80	0.78
	Clarin	45	44,299	39,799	0.10	NaN	0.62	0.62	0.62	0.64	0.63
	GOP	8	7,121	6,321	0.11	NaN	0.62	0.64	0.62	0.63	0.63
	Healthcare	1	590	490	0.17	NaN	0.53	0.56	0.53	0.47	0.50
	Obama	2	1,777	1,577	0.11	NaN	0.54	0.56	0.57	0.56	0.56
	SemEval	13	12,145	10,845	0.11	NaN	0.61	0.62	0.62	0.63	0.62
UNCERTAINTY	Riloff	2	1,201	1,001	0.17	NaN	0.80	0.82	0.84	0.82	0.81
	Swamy	1	555	455	0.18	NaN	0.37	0.40	0.40	0.33	0.36

Evaluation on remaining dataset Active learning systems often just track the test dataset performance. However, we observe another dataset which is not used for training, it is the left over dataset X_{left} after selecting the examples in each round. X_{left} is continuously decreasing and tracking the performance of the model on X_{left} can reveal how fast can an in-distribution dataset be accurately annotated using the specific querying strategy. This is suitable for simulation mode where the whole dataset ($X_{left} = X_{unlabeled}$) is already annotated.



(a) Progression of active learning classifier performance (micro f1-score) on the test set across 100 rounds of active learning. The annotation budget for each round is 100 instances, and the model is warm started with 100 random samples of the training data. Black dotted line is the classifier performance when trained on all of the training data. Data ordered alphabetically and X and Y axes are not shared.



(b) Active learning performance on unselected data across multiple classification tasks: Progression of active learning classifier performance (micro f1-score) on the unselected data set across 100 rounds of active learning. The annotation budget for each round is 100 instances, and the model is warm started with 100 random samples of the training data. Data ordered alphabetically and X and Y axes are not shared.

Figure 2: Abusive content detection results. Detailed results in section C.

Simulation Experiments Human annotation for `PyTAIL` can be simulated. First, X_{train} is set to $N = 100$ random samples from $X_{unlabeled}$. In each round, X_{select} is K ($K=100$) instances from $X_{unlabeled}$ based on the scoring criterion described above. We conduct 100 rounds of active learning (200 for Clarin as it is a very large dataset) and evaluate the models using the micro-f1 score. We also compare against a model trained on the full data (Full). The experimental results on the test split of each data are shown in figure 3 and table 1. We observe that the top K strategy is usually the best followed by the proportional strategy across all data. For larger datasets we see that the model closes the gap very soon. We also show experimental results on the X_{left} part of the training data in figure 4. We observe that the top K strategy is consistently the best, followed by the proportional strategy across all data. The increase in performance on the X_{left} is indicative of the fact that active learning ensures that the remaining data is actually easy to annotate without human correction. This evaluation presents a more practical usage pattern of ML models. This usage pattern requires annotating pre-selected and large $X_{unlabeled}$. In reality, once the dataset is selected, one is interested in reducing the size of X_{train} to efficiently annotate the data. We think, it is in this setting that the active learning is most beneficial. If the user can achieve high labeling accuracy by annotating few samples, then the user’s job is done.

4 Conclusion

We described experiments for evaluating active learning approaches for text classification tasks on tweet data. We introduced, `PyTAIL`, a user interface for active learning of NLP models by only requiring the user to update the labels for the model prediction if required. One limitation of our work is that our experiments are only conducted using simple linear model as they are easier to experiment with for sparse text features which we used for feature importance. However, the API does not place any restriction on the type of model. `PyTAIL` will publicly available as an open source tool.

References

- Hal Daumé III. Frustratingly Easy Domain Adaptation. *Association for Computational Linguistics (ACL)s*, (June):256–263, 2007. ISSN 0736587X. doi: 10.1.1.110.2062. URL <http://arxiv.org/abs/0907.1815>.
- Jacob Eisenstein. What to do about bad language on the internet, 6 2013. URL <https://aclanthology.coli.uni-saarland.de/papers/N13-1037/n13-1037https://www.aclweb.org/anthology/N13-1037>.
- Antigoni-Maria Founta, Constantinos Djouvas, Despoina Chatzakou, Ilias Leontiadis, Jeremy Blackburn, Gianluca Stringhini, Athena Vakali, Michael Sirivianos, and Nicolas Kourtellis. Large Scale Crowdsourcing and Characterization of Twitter Abusive Behavior. In *International AAAI Conference on Web and Social Media*, 2 2018. URL <http://arxiv.org/abs/1802.00393https://www.aaai.org/ocs/index.php/ICWSM/ICWSM18/paper/view/17909>.
- Manish Gupta, Rui Li, Zhijun Yin, and Jiawei Han. Survey on social tagging techniques. *ACM SIGKDD Explorations Newsletter*, 12(1):58, 2010. ISSN 19310145. doi: 10.1145/1882471.1882480. URL <http://portal.acm.org/citation.cfm?doid=1882471.1882480>.
- Shubhanshu Mishra and Jana Diesner. Detecting the Correlation between Sentiment and User-level as well as Text-Level Meta-data from Benchmark Corpora. In *Proceedings of the 29th on Hypertext and Social Media - HT '18*, pages 2–10, New York, New York, USA, 2018. ACM Press. ISBN 9781450354271. doi: 10.1145/3209542.3209562. URL <http://dl.acm.org/citation.cfm?doid=3209542.3209562>.
- Igor Mozetič, Miha Grčar, Jasmina Smailović, H Alani, Igor Mozetič, and A Scala. Multilingual Twitter Sentiment Classification: The Role of Human Annotators. *PLOS ONE*, 11(5):e0155036, 5 2016. ISSN 1932-6203. doi: 10.1371/journal.pone.0155036. URL <http://dx.plos.org/10.1371/journal.pone.0155036>.
- Ellen Riloff, Ashequl Qadir, Prafulla Surve, Lalindra De Silva, Nathan Gilbert, and Ruihong Huang. Sarcasm as Contrast between a Positive Sentiment and Negative Situation. In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, pages 704–714, Seattle, Washington, USA, 2013. Association for Computational Linguistics. URL <https://aclweb.org/anthology/papers/D/D13/D13-1066/>.
- Sunita Sarawagi. Information extraction. *Foundation and Trends in Databases*, 1(3):261–377, 1 2008. ISSN 1541-1672. doi: 10.1561/1500000003. URL <http://dl.acm.org/citation.cfm?id=234209http://portal.acm.org/citation.cfm?doid=234173.234209http://dl.acm.org/citation.cfm?id=1498844.1498845http://dl.acm.org/citation.cfm?id=1498845>.
- Burr Settles. Active Learning Literature Survey. Technical report, University of Wisconsin–Madison, 2009. URL <http://burrsettles.com/pub/settles.activelearning.pdf>.
- Sandesh Swamy, Alan Ritter, and Marie-Catherine de Marneffe. "i have a feeling trump will win.....": Forecasting Winners and Losers from User Predictions on Twitter. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 1583–1592, Stroudsburg, PA, USA, 2017. Association for Computational Linguistics. doi: 10.18653/v1/D17-1166. URL <http://aclweb.org/anthology/D17-1166>.
- Zeeraq Waseem and Dirk Hovy. Hateful Symbols or Hateful People? Predictive Features for Hate Speech Detection on Twitter. In *Proceedings of the NAACL Student Research Workshop*, pages 88–93, Stroudsburg, PA, USA, 2016. Association for Computational Linguistics. doi: 10.18653/v1/N16-2013. URL <http://aclweb.org/anthology/N16-2013>.

Checklist

The checklist follows the references. Please read the checklist guidelines carefully for information on how to answer these questions. For each question, change the default **[TODO]** to **[Yes]**, **[No]**, or **[N/A]**. You are strongly encouraged to include a **justification to your answer**, either by referencing the appropriate section of your paper or providing a brief inline description. For example:

- Did you include the license to the code and datasets? **[Yes]** See Section ??.
- Did you include the license to the code and datasets? **[No]** The code and the data are proprietary.
- Did you include the license to the code and datasets? **[N/A]**

Please do not modify the questions and only use the provided macros for your answers. Note that the Checklist section does not count towards the page limit. In your paper, please delete this instructions block and only keep the Checklist section heading above along with the questions/answers below.

1. For all authors...
 - (a) Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope? **[Yes]**
 - (b) Did you describe the limitations of your work? **[Yes]**
 - (c) Did you discuss any potential negative societal impacts of your work? **[N/A]**
 - (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? **[Yes]**
2. If you are including theoretical results...
 - (a) Did you state the full set of assumptions of all theoretical results? **[N/A]**
 - (b) Did you include complete proofs of all theoretical results? **[N/A]**
3. If you ran experiments...
 - (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? **[No]** No URL for anonymity. Our code is public.
 - (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? **[Yes]**
 - (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? **[No]**
 - (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? **[No]** No GPU used. All experiments on CPU.
4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
 - (a) If your work uses existing assets, did you cite the creators? **[Yes]**
 - (b) Did you mention the license of the assets? **[No]**
 - (c) Did you include any new assets either in the supplemental material or as a URL? **[No]**
 - (d) Did you discuss whether and how consent was obtained from people whose data you're using/curating? **[N/A]**
 - (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? **[Yes]**
5. If you used crowdsourcing or conducted research with human subjects...
 - (a) Did you include the full text of instructions given to participants and screenshots, if applicable? **[N/A]**
 - (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? **[N/A]**
 - (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? **[N/A]**

A Datasets

A.1 Sentiment classification

For sentiment classification we use the same data as in [Mishra and Diesner, 2018]. A description of these data is shown in table 2. Clarin Mozetič et al. [2016] and SemEval are the two largest corpora. However, SemEval has a larger test set. All the sentiment datasets use the traditional labels of positive, neutral, and negative for labeling the tweets.

Table 2: Description of sentiment classification datasets. Datasets clustered together are enclosed between horizontal lines. Labels are *negative, neutral, positive*.

data	split	tokens	tweets	vocab
Airline	dev	20079	981	3273
	test	50777	2452	5630
	train	182040	8825	11697
Clarin	dev	80672	4934	15387
	test	205126	12334	31373
	train	732743	44399	84279
GOP	dev	16339	803	3610
	test	41226	2006	6541
	train	148358	7221	14342
Healthcare	dev	15797	724	3304
	test	16022	717	3471
	train	14923	690	3511
Obama	dev	3472	209	1118
	test	8816	522	2043
	train	31074	1877	4349
SemEval	dev	105108	4583	14468
	test	528234	23103	43812
	train	281468	12245	29673

A.2 Abusive content classification

The second task we consider is abusive content classification. This task has recently gained prominence, owing to the the growth of abusive content on social media platforms. We utilize two datasets of abusive content. The first data is Founta from Founta et al. [2018], which tags tweets as *abusive, hateful, normal, spam*. The second dataset is WaseemSRW from Waseem and Hovy [2016]. It tags the data as *none, racism, sexism*. The rationale for including both these data under the same task is the core idea of identifying abusive content either direct or using racist or sexist variation. A description of these data is shown in table 3.

Table 3: Description of abusive content classification datasets. Datasets which are clustered together are enclosed between horizontal lines. Labels for Founta are *abusive, hateful, normal, and spam*. Labels for WaseemSRW are *none, racism, and sexism*.

data	split	tokens	tweets	vocab
Founta	dev	102534	4663	22529
	test	256569	11657	44540
	train	922028	41961	118349
WaseemSRW	dev	25588	1464	5907
	test	64893	3659	10646
	train	234550	13172	23042

A.3 Uncertainty indicators

Finally, we consider a collection of datasets for the task of identifying uncertainty indicators. Uncertainty indicators are defined as indicators in text which capture a level of uncertainty about the text, e.g., veridictality or sarcasm (uncertainty in intended meaning). We consider two datasets for this task as well. The first dataset is Riloff from Riloff et al. [2013]. This dataset consists of tweets annotated for sarcasm and non-sarcasm. The second dataset is Swamy from Swamy et al. [2017]. This dataset tries to identify the level of veridictality or degree of belief expressed in the tweet. The label set for this data is *definitely no*, *probably no*, *uncertain*, *probably yes*, *definitely yes*. A description of these data is shown in 4.

Table 4: Description of uncertainty indicators dataset. Datasets which are clustered together are enclosed between horizontal lines. Labels for Riloff are *sarcasm* and *not sarcasm*. Labels for Swamy are *definitely no*, *definitely yes*, *probably no*, *probably yes*, and *uncertain*.

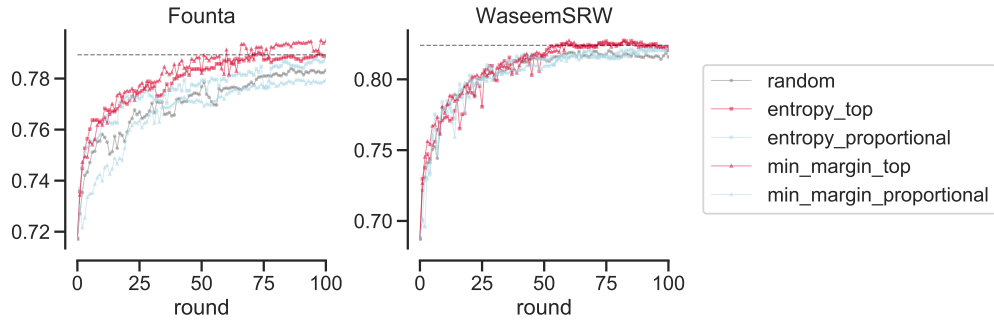
data	split	tokens	tweets	vocab
Riloff	dev	2126	145	1002
	test	5576	362	1986
	train	19652	1301	5090
Swamy	dev	1597	73	738
	test	3909	183	1259
	train	14026	655	2921

B Data Pre-processing

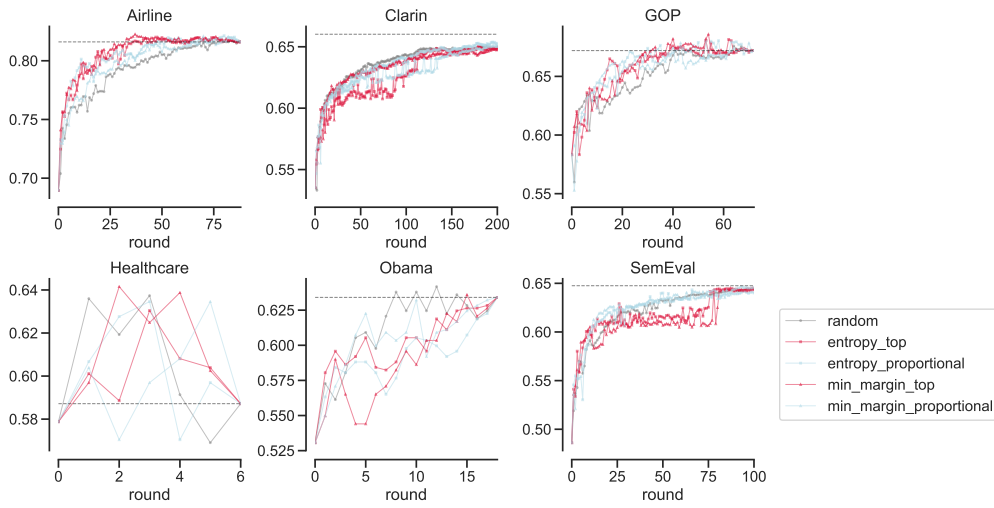
Each tweet is tokenized and pre-processed by normalizing all mentions of hashtags, URLs, and mentions. We also use a large sentiment lexicon². Furthermore, we suggest including a domain specific negative filter, i.e., words which should not be used to identify classification signals. For sentiment classification this can be entities in the corpora which should not bias the model.

C Results

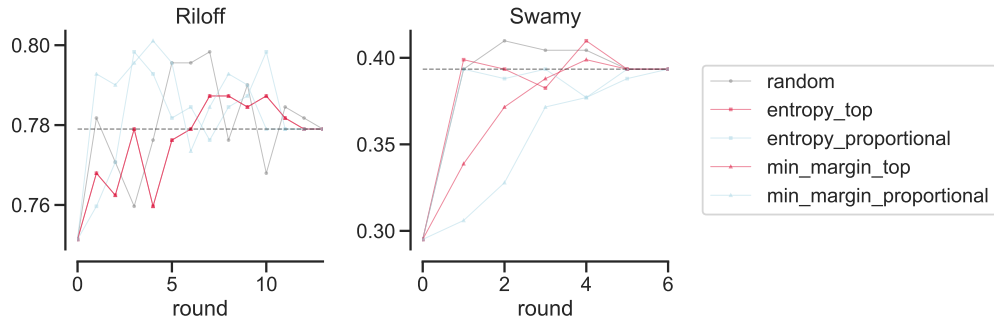
²<https://github.com/juliasilge/tidytext/blob/master/data-raw/sentiments.csv>



(a) Abusive content detection

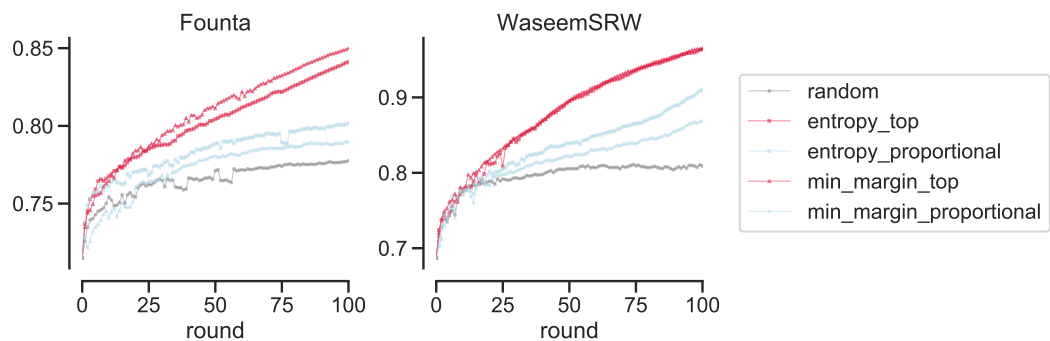


(b) Sentiment classification

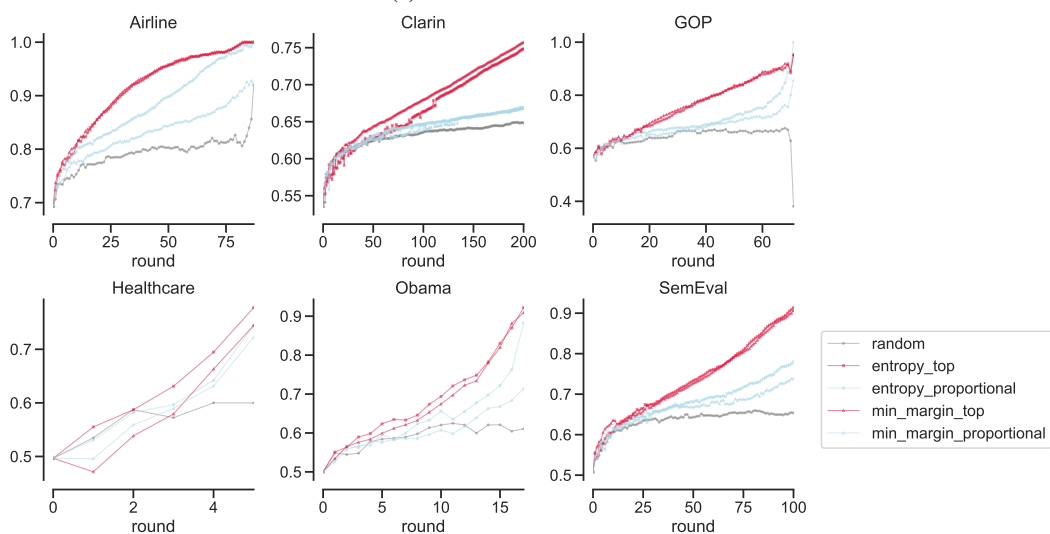


(c) Uncertainty indicators

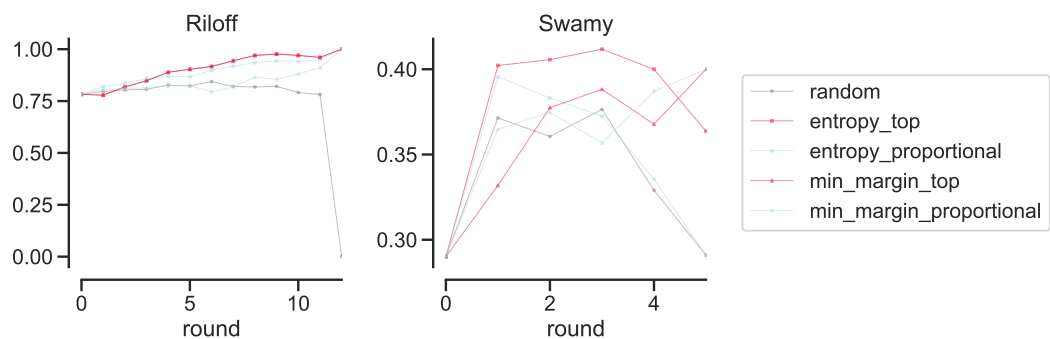
Figure 3: Progression of active learning classifier performance (micro f1-score) on the respective test set across 100 rounds of active learning (200 for Clarin). The annotation budget for each round is 100 instances, and the model is warm started with 100 random samples of the training data. Black dotted line is the classifier performance when trained on all of the training data. Data ordered alphabetically and X and Y axes are not shared.



(a) Abusive content detection



(b) Sentiment classification



(c) Uncertainty indicators

Figure 4: Progression of active learning classifier performance (micro f1-score) on the respective unselected data set across 100 rounds of active learning (200 for Clarin). The annotation budget for each round is 100 instances, and the model is warm started with 100 random samples of the training data. Data ordered alphabetically and X and Y axes are not shared.