

# **Active Learning in Natural Data Processing**

Author: **Marko Sahan**

# Contents

<b>Motivation</b>	<b>3</b>
<b>I Active Learning</b>	<b>5</b>
<b>1 Single Instance Active Learning</b> <sup>1</sup>	<b>5</b>
1.1 Introduction . . . . .	5
1.2 Methods . . . . .	6
Deep Ensemble Filter (DEnFi) . . . . .	7
Dropout MC . . . . .	7
SGLD . . . . .	7
Softmax uncertainty . . . . .	8
1.3 Experiments . . . . .	8
1.3.1 Hyperparameter Tuning . . . . .	8
1.3.2 Text Category Classification . . . . .	8
1.3.3 Fake News Classification . . . . .	12
1.4 Conclusion . . . . .	12
<b>2 Batch Active Learning</b> <sup>2</sup>	<b>13</b>
2.1 Introduction . . . . .	13
2.2 Batch Active Learning Methods . . . . .	14
1a. Uncertainty representation . . . . .	14
1b. Initialization of the training . . . . .	14
2. Acquisition function . . . . .	15
3. Batch selection strategy . . . . .	15
2.3 Experiment Setup . . . . .	15
2.4 Results . . . . .	17
2.5 Conclusion . . . . .	18
<b>3 Objectives</b>	<b>18</b>
<b>II Coreference Resolution</b>	<b>19</b>
<b>1 Introduction to Coreference Resolution Problem</b>	<b>19</b>
<b>2 Objectives</b>	<b>20</b>
<b>Summary</b>	<b>21</b>

---

<sup>1</sup>The major part of the content in this part was taken and restructured from [43].

<sup>2</sup>The major part of the content in this part was taken from the paper sent to COLING conference.

## Motivation

The generalized aim of Natural Language Processing (NLP) is a decision making process given the contextual understanding. Over the last decade, the field made significant progress thanks to both spotlight attention of the commercial sector and many scientists and also to the immense investment and improvements in computational power. The current list of NLP tasks includes but is not limited to such problems like: topic classification [58], machine translation [52], named entity recognition [59], language modeling [50], coreference resolution [13], etc.. Astonishingly, a vast amount of task-specific solutions show almost human-like semantic understanding [53] that brought even more attention to the field. The increased research activity resulted in a brain gain and a strong push of the boundaries in a majority of NLP problems, e.g representing language in a vector space [56]. The newer and more complex models' architecture approaches, accompanied by the motto "the more parameters, the better", made the retraining more complicated and costly. However, the innovative fine tuning approaches present efficient solutions for the problem of the full retraining paradigm [43]. The non arguable outcome of these improvements is an increase in the number of training instances i.e language models are usually trained with billions of text data [16].

Our research contribution lies in optimal models learning both from scratch or fine tuning of the pre-trained parameters. The training data collection is a highly time consuming and complicated procedure. In terms of supervised learning, the target labels acquisition is sometimes exceedingly expensive, e.g legal documents relevancy labeling [46]. The idea of the research is to generalize the optimal learning and present its functionality to one of the most complex problems - coreference resolution. We simultaneously explore various correlated sectors such as i) active learning, ii) models' uncertainty representation and iii) coreference resolution. The topics of *active learning* and *models uncertainty representation* are codependent. We propose a more granular study and extension of the human-machine communication. Among a huge quantity of unlabeled data the machine is expected to tell to an annotator (subject matter expert or so called oracle) which labels from the set of unlabeled data will bring forth the most information about the studied problem. Hence, the learning involves less data annotations and the model is trained with less effort. The following task is formulated as a supervised learning technique trained with the data obtained by the sequential labels querying from a human expert. This branch of research involves integration of enhanced models uncertainty measurement algorithms e.g deep ensembles [29], MC Dropout [18], SGLD [61] or Vadam [26] to the NLP problem for the empirical model weights distribution estimate. The additional information, given the distribution of predicted labels and the model prediction uncertainty [43], allows for faster model learning (hot and warm start methods) with a lower number of training data. The architecture agnostic generalization of the empirical model weights distribution estimate will grant us more freedom in NLP models selection while preserving the precise insight into the model processes given prediction-based uncertainty. The model uncertainty measurement and representation are done through the empirical estimate and sampling from the model weights distribution. The uncertainty representation approach provides the model with an expanded vision of both model learning and inference. The described technique has shown that such algorithms may enhance the learning process significantly [40].

The subsequent branch of the research is the active learning and model uncertainty integration into the coreference resolution problem. Coreference resolution (CR) combines detection and linking various mentions of entities within the text: linking noun phrases with their counterparts and pronouns, anaphora disambiguation, linking words with their pro-forms, etc.. Hence, CR-solving models significantly impact the quality of the text mining algorithms. The state-of-art CR models' architectures [13] are not as massive as the language models and not as trivial as different NLP tasks. Thus, our inclination towards more complex neural network structures, tremendous requirements for context understanding, vast usage potential, and personal interest make this type of model the perfect candidate for uncertainty algorithms integration and improvement. A deeper insight into the models' decision making process given the uncertainty is expected to allow for a new coreference resolution state-of-art threshold to be set.

An exemplary use case where coreference resolution can be applied is categorizing entities and their pronouns to provide one with a broader spectrum of information for future decision making. Based on the extracted data, it is possible to unify all knowledge in the form of a Knowledge Graph (KG) [57] which can be further utilized for linking concepts represented by textual spans. Dependencies and connections between the entities can enrich the feature space with highly discriminative samples for other tasks. For example, assume that we have the following two consecutive sentences: "John Smith and Amanda Brown are accountants at XYZ company. Amanda's colleague was accused of drunk driving". Based on these sentences, one would wish to classify if some of the entities from the text can be charged with a misdemeanor. For a human reader, it is evident that Amanda's colleague refers to John. However, for a machine, that is a challenging task. Therefore, proper identification of entity clusters like

John Smith, Amanda’s colleague, Amanda Brown, XYZ would significantly improve the machine’s understanding of the text. Another potential application of coreference resolution lies within the problem of opinion mining in media resources, where people frequently express their views and opinions. For example, heated discussions may emerge under political news articles. In these discussions, participants refer to subjects of the particular article with, for instance, pronouns. Therefore, proper CR may provide better traction of the audience’s attitude towards entities from the article by linking comment mentions to them.

The research paper covers the topic of active learning in the first part. The active learning part introduces a reader to single instance and batch active learning problem study, where the content is taken from our articles ([43] best PhD student article and [44] <sup>3</sup>). This part was well studied with the additional contribution to the field. In the end of the first part we highlight what is our objective for the further work. The second part gives a general introduction to a coreference resolution problem with a detailed breakdown of the field’s history and state-of-the-art. We also define our objective for further research that follows with a summary.

---

<sup>3</sup>The paper is in review for COLING conference

## Part I

# Active Learning

In this section we show two vast comparative studies that cover not only single instance and batch active learning problems but also a newer approach on fine tuning the algorithms without a full retraining. The research involves a detailed state-of-the-art algorithms study in the active learning field with the additional novelty, modification and extension of learning algorithms.

## 1 Single Instance Active Learning <sup>4</sup>

Supervised classification of texts relies on the availability of reliable class labels for the training data. However, the process of collecting data labels can be complex and costly. A standard procedure is to add labels sequentially by querying an annotator until reaching satisfactory performance. Active learning is a process of selecting unlabeled data records for which the knowledge of the label would bring the highest discriminability of the dataset. In this part, we provide a comparative study of various active learning strategies for different embeddings of the text on various datasets. We focus on Bayesian active learning methods that are used due to their ability to represent the uncertainty of the classification procedure. We compare three types of uncertainty representation: i) SGLD, ii) Dropout, and iii) deep ensembles. The latter two methods in cold- and warm-start versions. The texts were embedded using Fast Text, LASER, and RoBERTa encoding techniques. The methods are tested on two types of datasets, text categorization (Kaggle News Category and Twitter Sentiment140 dataset) and fake news detection (Kaggle Fake News and Fake News Detection datasets). We show that the conventional dropout Monte Carlo approach provides good results for the majority of the tasks. The ensemble methods provide more accurate representation of uncertainty that allows to keep the pace of learning of a complicated problem for the growing number of requests, outperforming the dropout in the long run. However, for the majority of the datasets the active strategy using Dropout MC and Deep Ensembles achieved almost perfect performance even for a very low number of requests. The best results were obtained for the most recent embeddings RoBERTa

### 1.1 Introduction

The development of a text classifier on a new problem requires the availability of the training data and their labels. Labeling involves human annotators and a common practice is to label as many text documents as possible, train a classifier and search for new data and labels if the performance is unsatisfactory. Random choice of the documents for the data set extension can be costly because the new documents may not bring new information for the classification. Active learning strategy aims to select among available unlabeled documents those that the classifier is most uncertain about and queries an annotator for their labels. Therefore, it has the potential to greatly reduce the effort needed for the development of a new system. While it was introduced almost two decades ago, recent improvements in deep learning motivate our attempt to revisit the topic. For example, SVM-based active learning approaches for text classification date back to 2001 [54], where the superiority of active learning over random sampling was demonstrated. Since deep recurrent and convolutional neural networks achieve better classification results, Bayesian active learning methods for deep networks gained popularity especially in image classification [18], [34].

The Bayesian approach is concerned with querying labels for data for which the classifier predicts the greatest uncertainty. The uncertainty is quantified using the so-called acquisition function, such as predictive variance or predictive entropy [47]. While different acquisition functions often provide similar results, different representations of predictive distribution yield much more diverse results. The most popular approach using Dropout MC [18] has been tested on text classification [3] and named entity recognition [48], [34], however other techniques such as Langevin dynamics [61] and deep ensembles [30] are available. Deep ensembles often achieve better performance [6], [51] but require higher computational cost since they train an ensemble of networks after each extension of the data set. One potential solution of this problem has been recently proposed in [55], where the ensemble is not trained from a fresh random initialization after each query but initialized randomly around the position of the ensembles from the previous iteration. In this contribution, we test this approach and compare it with the dropout MC and Langevin dynamics representations. We also provide sensitivity study for the choice of the hyperparameters. Active learning for fake news detection was considered in [7] using uncertainty based on

---

<sup>4</sup>The major part of the content in this part was taken and restructured from [43].

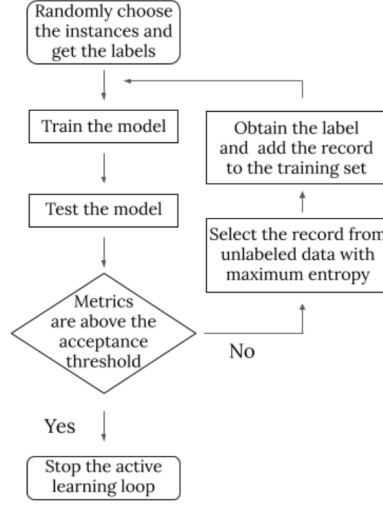


Figure 1: Flowchart of the active learning algorithm

probability of classification. It was later extended to a context aware approach [8]. An entropy based approach has been presented in [21] using an ensemble of three different kinds of networks.

## 1.2 Methods

Throughout the paper, we will use three different embedding algorithms such as Fast Text [37], LASER [4] and RoBERTa [33]. RoBERTa is a modified BERT transformer model [12] based on multi-head attention layers [56]. Transformer models provide state of the art results in context understanding and it is advantageous to compare behavior of active learning algorithms with respect to different embedding techniques. Representation of the  $i$ -th text document  $\mathbf{x}_i$  is calculated as the mean value from sentence embeddings of all sentences in the text

$$\mathbf{x}_i = \frac{1}{|D_i|} \sum_{j \in D_i} f_{\text{sentence embed}}(C^{(j)})$$

where  $D_i$  is the set of vectors where each vector represents a sentence in the  $i$ -th document,  $|D_i|$  is a cardinality of  $D_i$ ,  $C^{(j)}$  is a matrix of  $j$ -th sentence where words are encoded with one hot or byte pair encoding [49] technique and  $f_{\text{sentence embed}}$  is a function that creates sentence embeddings with respect to the given one-hot or byte pair encoded words. Fast Text encoding is made with respect to one-hot encoded words. LASER and transformer based models take byte-per encoded words as an input. Sentence embeddings for LASER and RoBERTa are output of deep neural network models. Sentence embedding for the Fast Text model is calculated as a mean value of all embeddings of words in the sentence. For supervised classification, each document embedding  $\mathbf{x}_i$  has to have an associated label  $\mathbf{y}_i$ . We are concerned with binary classification for simplicity, however, an extension to multiclass is straightforward. We assume that for the full corpus of text documents  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]$ , only an initial set of  $l_0 \ll n$  labels is available,  $\mathbf{Y}^{(0)} = [\mathbf{y}_0, \dots, \mathbf{y}_{l_0}]$ , splitting the full set  $\mathbf{X}$  to the labeled,  $\mathbf{X}^{(0)} = [\mathbf{x}_1, \dots, \mathbf{x}_{l_0}]$ , and unlabeled parts,  $\mathbf{X} \setminus \mathbf{X}^{(0)}$ . Active learning is defined as a sequential extension of the training data set. In each iteration,  $l = 1, \dots, L$ , the algorithm computes entropy of the predictive probability distribution for each document in the unlabeled dataset and selects the index of the document with the highest entropy (entropy acquisition function), formally:

$$k_l = \arg \max_{k \in K} \mathbb{E}_{p(\theta|\mathbf{X}^{(l-1)}, \mathbf{Y}^{(l-1)})} (H(\mathbf{y}_k|\theta)) \quad (1)$$

$$H(\mathbf{y}_k|\theta) = \mathbb{E}_{p(\mathbf{y}_k|\theta, \mathbf{x}_k)} (-\log p(\mathbf{y}_k|\theta, \mathbf{x}_k)) \quad (2)$$

where  $K$  is the set of indexes of all unlabeled documents,  $\mathbb{E}$  is the expectation operator over the posterior probability of the classifier parameters  $\theta$  trained on all labeled data  $p(\theta|\mathbf{X}^{(l-1)}, \mathbf{Y}^{(l-1)})$  and  $H(\mathbf{y}_k|\theta)$  is conditional entropy of the predicted class for  $\mathbf{x}_k$ . The document of the selected index is sent to the human annotator with a request for labeling. When the selected text is annotated, the text is added with its label to the labeled data set  $\mathbf{X}^{(l)} = [\mathbf{X}^{(l-1)}, \mathbf{x}_{k_l}]$ ,  $\mathbf{Y}^{(l)} = [\mathbf{Y}^{(l-1)}, \mathbf{y}_{k_l}]$ . The procedure is repeated  $L$  times. The active learning process is visualized in Figure 1.

The key component of the method is a representation of the posterior distribution of the parameter  $\theta$ . Due to the complexity of the neural networks it is always represented by samples, with a different method of their generation. We will compare the following methods: i) SGLD: Stochastic Gradient with Langevin dynamics [61], which adds additional noise to the gradient in stochastic gradient descent, ii) Dropout MC: is an extension of the ordinary dropout that samples binary mask multiplying output of a layer, hence stopping propagation through all neurons where zeros is sampled through the network. The extension applies the sampled mask even for predictions generating samples from the predictive distribution [18], iii) Deep ensembles: consist of  $N$  networks trained in parallel from different initial conditions [30], and iv) Softmax uncertainty: is the most simple approach that uses only one network to estimate a single value  $\theta$  and maximizing entropy  $H(\mathbf{y}_k|\theta)$  instead of the expectation 1 [7]. Ensembles based approach is the current state-of-the-art in active learning [6]. While many of these have been tested in active learning, the majority of authors assumed that after each step of active learning, the network training starts from the initial conditions. This is clearly suboptimal, since the information from the previous training is lost. A simple solution was presented in [55], where it was argued that estimated results from the previous step can be used as centroids around which the new initial point is sampled. Since this is a form of a warm-start, we also test warm-start strategies for Dropout. The methods for representation of parametric uncertainty are:

### Deep Ensemble Filter (DEnFi):

is a deep ensemble method with 10 neural networks in the ensembles and warm-start training strategy [30] using weights of the ensemble members in the previous iteration as initial conditions for the new ensemble. Each weight is perturbed by an additive Gaussian noise of variance  $q$  which is a hyperparameter. In our experiments, the ensemble is trained with parameters `initialization_epochs` = 2000 on the initial data and with additional `warm_start_epochs` = 700 epochs after each extension of the learning data set. DEnFi algorithms are displayed in Algorithms 1.

---

#### Algorithm 1 DEnFi active learning training algorithm

---

**Initialize:** ensemble classifier structure  $\mathbf{y} = \mathbf{y}(\mathbf{x}, \theta_1^{(0)}, \dots, \theta_N^{(0)})$ , where  $\theta_j^{(0)}$ ,  $j \in \{1, \dots, N\}$  are different initializations of classifier parameters. Initial data  $\mathbf{Y}^{(0)}, \mathbf{X}^{(0)}$ . Noise perturbation variance is  $q$ . Set iteration counter  $i = 0$

**Iterate** until a stopping condition:

1. For each parameters' set  $\theta_j^{(i)}$ ,  $j \in \{1, \dots, N\}$  provide
    - (a) noise perturbation as  $\theta_j^{(i)} = \theta_j^{(i)} + \gamma_j^{(i)}$ ,  $\gamma_j^{(i)} \sim \mathcal{N}(0, qI)$
    - (b) Train a classifier for  $\mathbf{X}^{(i)}, \mathbf{Y}^{(i)}$  and initial number of epochs  $\gg$  warm start epochs if  $i = 0$ , else warm start epochs
  2. Sample a new instance from Equation 1 and update  $\mathbf{X}^{(i+1)} = [\mathbf{X}^{(i)}, \mathbf{x}_{k_{i+1}}]$ ,  $\mathbf{Y}^{(i+1)} = [\mathbf{Y}^{(i)}, \mathbf{y}_{k_{i+1}}]$ ,  $i = i + 1$
- 

### Dropout MC:

is the standard algorithm [18] that trains only a single network with sampled dropout indices and uses the sampling even in the prediction step. Generation of the Monte Carlo prediction is obtained by sampling different values of the dropout binary variable and one forward pass of the network for each sample. We study three versions of the algorithm: i) cold-start with 3000 epochs after each request, ii) warm-start with weights from the previous iteration perturbed by an additive noise of variance  $q$  with 700 epochs and iii) hot-start, with 50 epochs after each request without perturbation (hot start is a warm-start method with  $q = 0$ ). Dropout rate is 0.5.

### SGLD:

variance of the noise added to the gradient descent:

$$\hat{\theta}_{n+1} = \hat{\theta}_n - \frac{|\mathbf{X}^{(l)}|}{|\mathbf{X}^{(l)}|_{\text{minibatch}}} \frac{\eta_n}{2} \left( \nabla L(\mathbf{X}^{(l)}, \mathbf{Y}^{(l)}, \hat{\theta}_n) \right) + \epsilon_n, \quad \epsilon_n \sim \mathcal{N}(0, \eta_n \mathbf{I})$$

where  $L$  is a loss function and  $\epsilon_n$  is noise with covariate matrix  $\eta_n \mathbf{I}$  where  $\eta_n$  is the learning rate of the SGD algorithm. We choose initial value of  $\eta_1$  to be 0.01 and update rule  $\eta_{n+1} = \frac{\eta_n}{b-3000} + 0.05$ . The noise  $\epsilon_n$  is added to the gradient only after 3000 of initial training epochs. Then we draw 50 samples with 100 epochs between consecutive samples to avoid correlation.

### Softmax uncertainty:

The simplest approach to uncertainty representation is a single neural network with a softmax output layer that considers uncertainty as the output of the softmax score. We add it to comparison since it has been applied to active learning in [7]. The model is trained to run 2000 epochs on the initial data with additional 200 epochs after each extension of the learning data set.

## 1.3 Experiments

The methods were compared on the positive/negative tweets from the Tweets Dataset [19], 5 pairs of categories from the News Category Dataset [38] and two types of the news datasets for fake news detection [1], [2]. Documents from the News Category dataset were downloaded using links provided in the dataset. The names of the tested categories are shown in table 2, Figure 2a and in Figure 2b.

Specifically, we compared active learning and random sampling strategies for different settings of document embeddings and different representations of uncertainty. Each experiment was initiated by random choice of the initial training set of  $l_0 = 10$  samples from 1000 text documents (500 text documents per category), which were the initial 1000 documents of the datasets. For each experiment  $L = 200$  requests for annotation are simulated. The document selection follows the  $\epsilon$ -greedy approach [60], i.e. the sample with maximum acquisition function 1 is accepted with probability  $\epsilon = \frac{\exp(l-40)}{\exp(l-40)+1}$ . A random document is selected for labeling if not accepted. After each request, the classification performance is evaluated on the remaining part of the selected dataset (i.e. on the 990 text documents in the first evaluation) using the area under the ROC curve (AUC) metrics [15]. In order to make the results statistically valid, we repeat the described simulation loop 10 times for Twitter and News Category datasets and 5 times for Fake News and Fake News Detection datasets.

### 1.3.1 Hyperparameter Tuning

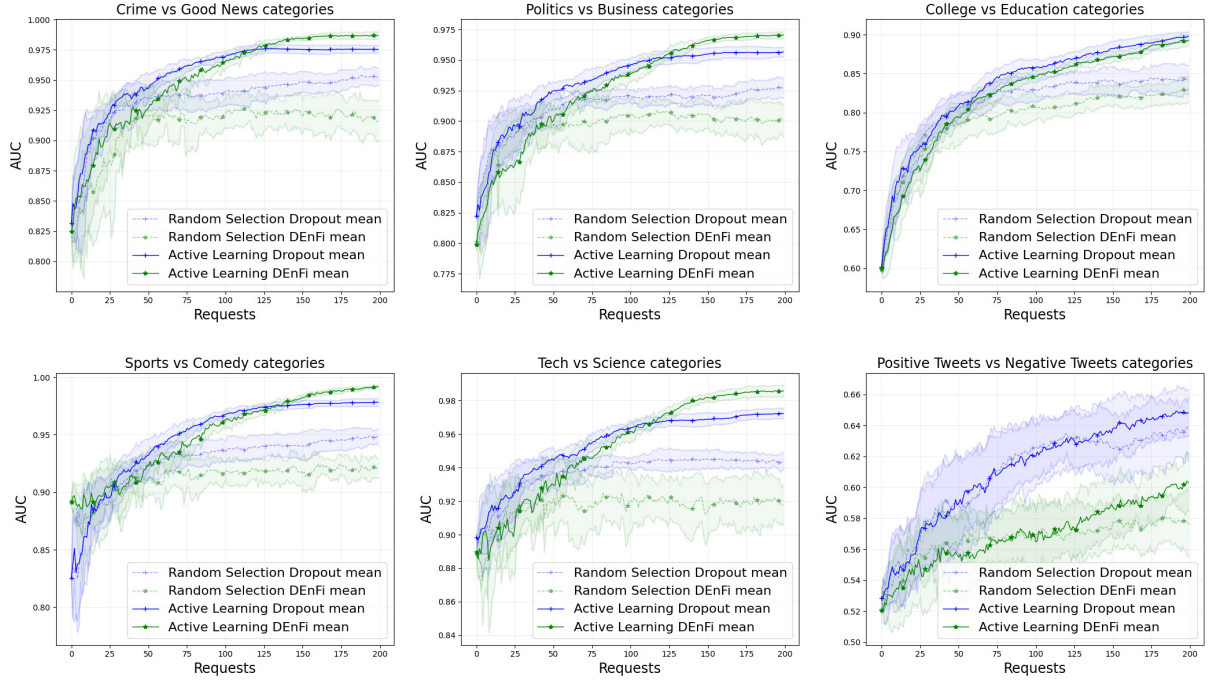
The classification network was designed as a feed-forward NN with one dense layer of 100 neurons with sigmoid activation functions, and softmax output layer. Warm start versions of both Dropout MC and DEnFi have hyperparameter  $q$  that governs the perturbation of the previous result before training on the extended dataset. Tuning of this hyperparameter was performed by a grid search for both the Fast Text and RoBERTa text encoding techniques. Since the main focus of the paper is on the effect of the warm-start strategy, the results of the effect of the noise variance  $q$  for DEnFi and Dropout MC is displayed in Table 1a and Table 1f for active learning strategy on the Tech vs Science task and Fake News Detection dataset, respectively.

Based on Table 1a it is clearly seen that the variance of the perturbation noise related to the best result is increasing with the number of requests. We conjecture that the variance has the role of a selector of the exploration/exploitation tradeoff. Low variance favors exploitation and improves quickly, higher variance implies less accurate guesses in the initial iterations but better performance in the long run. The results from Table 1a indicate that higher values of  $q$  are not performing well thus the range of  $q$  was reduced for tuning of the hyperparameter for the fake news datasets. The increase of the optimum value of the noise variance with the number of requests is also visible in Table 1f, especially for the RoBERTa encoded data. Since calibration of the variance for all methods and all datasets would be too computationally expensive, we ran all remaining experiments with  $q = 0.3$  for Fast Text encoding and  $q = 0.1$  for RoBERTa and LASER encodings. However, tuning of this hyperparameter for an application scenario or an adaptive tuning strategy offers clearly a potential for further improvement.

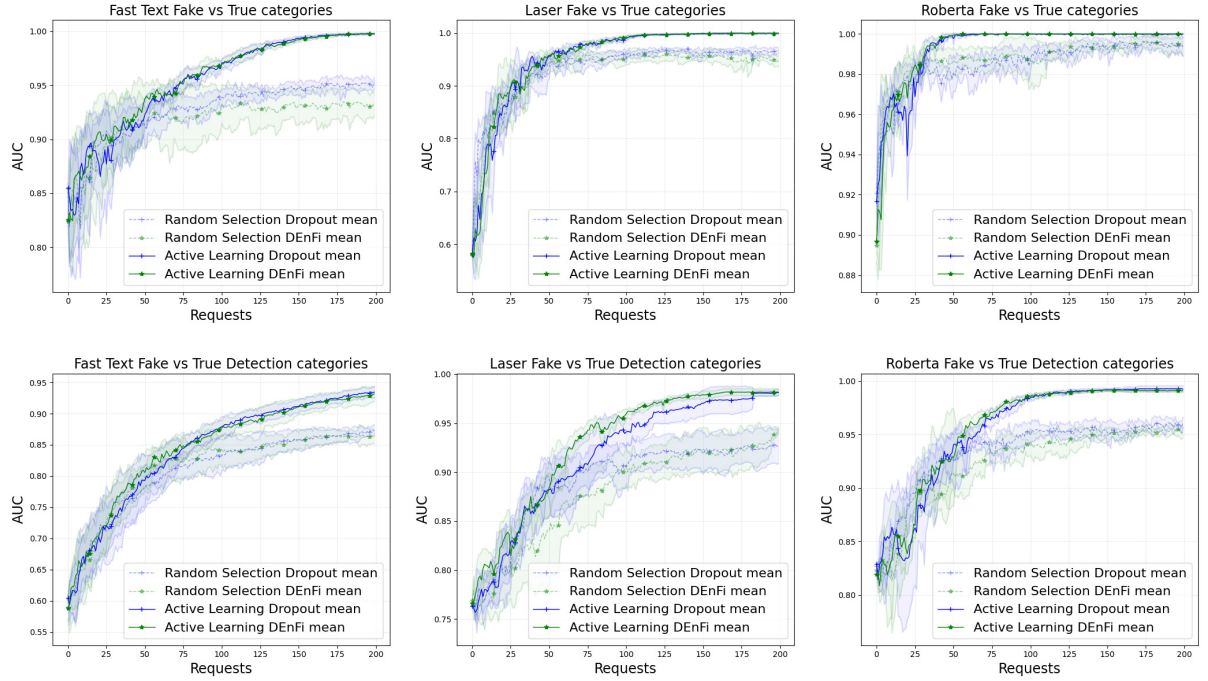
### 1.3.2 Text Category Classification

A comparison of AUC of the active learning strategy after 200 requests for all tested algorithms and Fast Text encoding with respect to Tweets and News Category datasets is reported in Table 2. This experiment was designed





(a) Results for six pairs of categories



(b) Results for two datasets: Fake News (top row) and Fake News Detection (bottom row), and three embeddings: Fats Text (left column), Laser (middle column) and RoBERTa (right column)

Figure 2: Evolution of the mean AUC with growing number of requests for DEnFi and Dropout MC warm-start algorithms. The uncertainty bounds are illustrated as one standard deviation from the mean value with respect to 10 runs for top figure and 5 runs for bottom figure. Both DEnFi and Dropout MC were initially trained on 10 labeled text documents before sequential learning strategies were initialized

Noise	AUC after # of iterations			
Variance	50	100	150	200
0.1	<b>0.945*</b>	<b>0.968*</b>	0.974	0.976
0.2	0.932	0.964	0.976	0.978
0.3	0.930	0.961	<b>0.982*</b>	0.986
0.4	0.909	0.948	0.976	<b>0.990*</b>
0.6	0.874	0.921	0.952	0.979
1	0.805	0.871	0.906	0.941

(b) Fast Text DEnFi

Noise	AUC after # of iterations			
Variance	50	100	150	200
0.1	<b>0.936</b>	<b>0.966*</b>	0.972	0.976
0.2	<b>0.938*</b>	<b>0.966*</b>	<b>0.981*</b>	0.983
0.3	0.920	0.956	<b>0.982</b>	<b>0.989*</b>
0.4	0.917	0.955	0.976	<b>0.988</b>
0.6	0.894	0.948	0.972	<b>0.986</b>
1	0.859	0.914	0.941	0.970

(c) Fast Text Dropout

Noise	AUC after # of iterations			
Variance	50	100	150	200
0.1	<b>0.935</b>	<b>0.970*</b>	<b>0.982</b>	0.988
0.2	<b>0.940*</b>	0.954	<b>0.983*</b>	<b>0.990*</b>
0.3	0.903	0.930	0.967	0.987
0.4	0.883	0.911	0.945	0.976
0.6	0.799	0.853	0.885	0.945
1	0.661	0.750	0.818	0.875

(d) RoBERTa DEnFi

Noise	AUC after # of iterations			
Variance	50	100	150	200
0.1	<b>0.930*</b>	<b>0.970</b>	0.982	0.986
0.2	<b>0.923</b>	<b>0.971*</b>	<b>0.986*</b>	<b>0.990*</b>
0.3	0.917	0.959	0.982	<b>0.990*</b>
0.4	0.878	0.949	0.974	<b>0.990*</b>
0.6	0.822	0.908	0.952	0.980
1	0.643	0.741	0.862	0.921

(e) RoBERTa Dropout

(e) Results for Tech vs Science categories

Noise	AUC after # of iterations			
Variance	50	100	150	200
0.1	<b>0.794</b>	<b>0.886*</b>	<b>0.910</b>	<b>0.942*</b>
0.2	<b>0.806*</b>	<b>0.877</b>	<b>0.911*</b>	<b>0.932</b>
0.3	<b>0.806*</b>	0.866	<b>0.901</b>	<b>0.911</b>

(g) Fast Text DEnFi

Noise	AUC after # of iterations			
Variance	50	100	150	200
0.1	<b>0.806</b>	<b>0.879*</b>	<b>0.916*</b>	<b>0.949*</b>
0.2	<b>0.804</b>	<b>0.884*</b>	0.914	0.939
0.3	<b>0.809*</b>	<b>0.878</b>	<b>0.912</b>	<b>0.944</b>

(h) Fast Text Dropout

Noise	AUC after # of iterations			
Variance	50	100	150	200
0.1	<b>0.929*</b>	<b>0.986*</b>	<b>0.992</b>	0.991
0.2	<b>0.923</b>	0.975	0.990	<b>0.998*</b>
0.3	<b>0.891</b>	0.935	0.955	0.975

(i) RoBERTa DEnFi

Noise	AUC after # of iterations			
Variance	50	100	150	200
0.1	<b>0.941*</b>	<b>0.984*</b>	<b>0.992*</b>	<b>0.993</b>
0.2	0.917	0.974	<b>0.992*</b>	<b>0.995*</b>
0.3	0.917	0.959	0.982	0.990

(j) RoBERTa Dropout

(j) Results for Fake and True news categories (Fake News Detection dataset)

Table 1: **Fast Text** and **RoBERTa** encoding based AUC of text classification after selected number of requests of the active learning using DEnFi and Dropout MC warm-start for various selection of the perturbation noise  $q$ . Average over 10 runs for top tables and 5 runs for bottom tables. The best result is denoted by a star, results within one standard deviation of the winner are displayed in bold font.

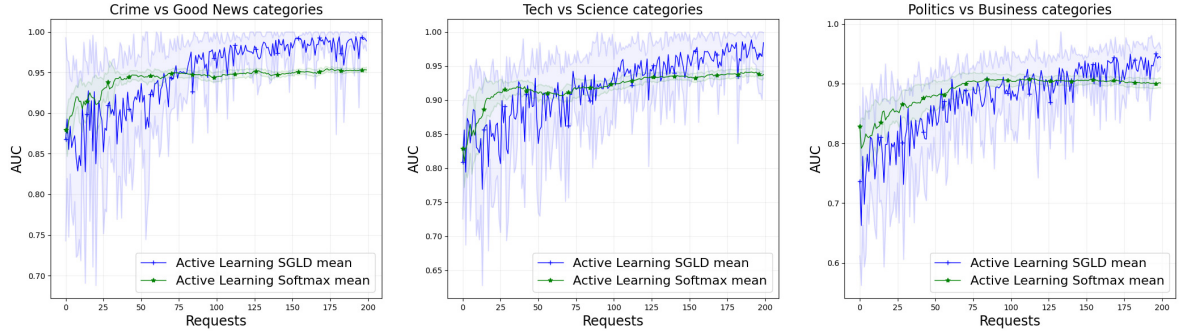


Figure 3: Evolution of the mean AUC with growing number of requests for SGLD, Softmax Uncertainty algorithms, and three pairs of categories. The uncertainty bounds are illustrated as one standard deviation from the mean value with respect to 10 runs. Both SGLD and Softmax Uncertainty were initially trained on 10 labeled text documents before sequential learning strategies were initialized

Method	Crime/ Good News	Sports/ Comedy	Politics/ Business	Tech/ Science	Education/ College	Pos./Neg. Tweets
SGLD	<b>0.989*</b>	0.968	0.944	0.984	0.881	0.621
DEnFi, $q = 0.3$	0.987	<b>0.992*</b>	<b>0.971*</b>	<b>0.986</b>	<b>0.893</b>	0.603
MC Dropout cold-start	0.975	0.978	0.957	0.972	<b>0.898*</b>	<b>0.648</b>
MC Dropout warm-start, $q = 0.3$	0.978	0.979	0.954	0.973	0.877	<b>0.657*</b>
MC Dropout warm-start, $q = 0$	0.978	0.951	0.944	<b>0.989*</b>	0.824	0.561
Softmax Uncertainty	0.953	0.939	0.901	0.938	0.800	0.609

Table 2: Average AUC over 10 runs for five different algorithms after 200 iterations of active learning and six different datasets. The best result is denoted by a star, results within one standard deviation of the winner are displayed in bold font.

for comparison of various uncertainty representations. While all methods except for the Softmax Uncertainty achieved best results on some datasets, the most consistent results were delivered by the DEnFi and Dropout MC methods. The SGLD performance results exhibited the highest variance within the 10 runs of the method. The results for three pair categories and an active learning strategy are visualized in figure 4. Since it is the least reliable method, we will not study it any further. It is clear from Table 2 and Figure 3 that the simple Softmax Uncertainty provides the worst results. Since Dropout MC has comparable computational complexity, the use of simple Softmax Uncertainty is always suboptimal.

Detailed analysis of the DEnFi and Dropout MC strategies is provided in Figure 2a for Tweets and News Category datasets. For better insight, we also display the performance of the random sampling learning strategy where training data are extended using randomly sampled documents. Note that for random sampling, Dropout MC outperforms consistently DEnFi on all tasks. We conjecture that this is due to the robustness of the dropout regularization. However, the power of DEnFi becomes apparent with the increasing number of requests. It is improving slower than Dropout MC at the beginning, but improves faster, thus outperforming Dropout MC in the long run. We conjecture that this is due to better exploration capability of the DEnFi while Dropout MC excels at exploitation. The speed of improvement depends on the complexity of the learning task. For simpler tasks (such as crime vs. Good News), AUC over 0.98 is achieved quickly. However, for more complex tasks, such as Positive vs Negative Tweets, the number of data needed for improvement is much higher. The active learning is on par with the random sampling up to 125 requests and even after 200 requests, the AUC is below 0.7 indicating poor performance. Note that the active learning strategy of DEnFi starts improving over the random sampling sooner than Dropout MC with a sharper slope which indicates a high probability of obtaining the same profile as the other datasets in the long run.

### 1.3.3 Fake News Classification

The fake news data experiment is made under the same initial conditions as in Section 1.3.2 but for a wider choice of the document encoding techniques. Based on results from Section 1.3.2 we show comparison of only the DEnFi and Dropout warm-start for three different types of encodings (Fast Text, LASER, RoBERTa) on the two fake news datasets in Figure 2b. The results for fake news datasets are similar to those for the Twitter and News Category datasets. Advantage of the DEnFi approach is significant only for the Fake News Detection dataset with Laser embeddings. For other cases the difference is negligible. Under the exploration/exploitation hypothesis, this means that the data sets are well separable and more sophisticated exploration capabilities of DEnFi are not needed.

The best performance of the active learning strategy was achieved for the RoBERTa embedding. Note that the embedding has much greater influence on the speed of learning than uncertainty representation.

## 1.4 Conclusion

We have studied the suitability of various uncertainty representations for the task of active text classification. The established Dropout MC methodology was compared against deep ensembles, SGLD and simple softmax strategy. To reduce the computational cost, we studied warm-start strategies for both ensembles (called DEnFi) and Dropout MC, that resulted in significant reduction of training epochs per the active learning iteration. The resulting methods exhibit a different tradeoff between exploration and exploitation. While Dropout MC has been found to be more reliable in random sampling strategy and improving faster at the beginning of active learning, the DEnFi was found to prefer exploration sacrificing performance at the beginning but outperforming Dropout MC in the longer run of active learning. The highest deviation of the random strategies between DEnFi and Dropout MC resulted in 0.04 AUC in favor of Dropout MC. However, for the active learning strategy both Dropout MC and DEnFi converged to the same numbers.

Sensitivity of the active learning to the choice of embedding was evaluated on the fake news datasets. It was observed that the most recent embedding method (RoBERTa) facilitated the fastest learning and that the choice of embedding was more important than the uncertainty representation. However, methods of active learning achieved significantly faster learning than the random sampling approach (the difference between the random and the active learning strategies differs from 0.015 AUC up to 0.080 AUC) on all tested datasets of various complexity.

## 2 Batch Active Learning<sup>5</sup>

Batch active learning is a generalization of a single instance active learning by selecting a batch of documents for labeling. This task is much more demanding because plenty of different factors come into consideration (i. e. batch size, batch evaluation, etc.). In this part, we provide a large scale study by decomposing the existing algorithms into building blocks and systematically comparing meaningful combinations of these blocks with a subsequent evaluation on different text datasets. While each block is known (warm start weights initialization, Dropout MC, entropy sampling, etc.), many of their combinations like Bayesian strategies with agglomerative clustering are first proposed in our paper with excellent performance. Particularly, our extension of the warm start method to batch active learning is among the top performing strategies on all datasets. We studied the effect of this proposal comparing the outcomes of varying distinct factors of an active learning algorithm. Some of these factors include initialization of the algorithm, uncertainty representation, acquisition function, and batch selection strategy. Further, various combinations of these are tested on selected NLP problems with documents encoded using RoBERTa embeddings. Datasets cover context integrity (Gibberish Wackerow), fake news detection (Kaggle Fake News Detection), categorization of short texts by emotional context (Twitter Sentiment140), and sentiment classification (Amazon Reviews). Ultimately, we show that each of the active learning factors has advantages for certain datasets or experimental settings.

### 2.1 Introduction

Supervised learning of classifiers relies on the availability of class labels which often involves a human annotator for a majority of NLP tasks. This can be costly for large datasets. Active learning is a strategy designed to minimize this cost by automatic selection of those unlabeled documents that are expected to bring useful information for the classifier. Advantages of this approach have been demonstrated even for classical methods such as SVM [54]. The most conventional active learning strategies select only one unlabeled document after each training round to query due to the simplicity of its selection. The next query document is selected only after the first one is labeled and the model retrained, which means that the annotator has to wait for retraining. This impractical strategy can be avoided if the active learning algorithm selects a batch of documents. Novel methods for batch active learning appear frequently, each demonstrating advantages on their benchmark data.

Various comparative studies have been performed recently with various focus and results. Batch active learning was studied in [14] for only one size of the batch (50 documents per query). Large sensitivity to the type of dataset was reported in [41], where different methods won for different data. Large variability of the results was also observed in [24]. In [45], the min-margin strategy was shown to be competitive with the prediction entropy-based method on a range of embeddings. The comparative studies shared similar properties, such as a fixed network for embeddings (improvement with retraining can be expected [35] but may be too costly). All studies also assume a cold start, i.e. completely new initialization of the classifier after each round of querying. This is motivated by the fear of overfitting, which was demonstrated in [23] for hot start, i.e. continuation of training of the classifier. A compromise in the form of warm-start, i.e. adding noise to the weights of the previous classifier, was proposed in [43].

In this contribution, we take a different approach to benchmarking of the batch active learning algorithms. Specifically, we decompose the algorithms into their building blocks: i) the size of the minibatch, ii) acquisition function, iii) representation of uncertainty of the classifier, and iii) initialization of the network. This approach allows us to quantify contribution of each of the building-block and combine them in previously untested versions. This allows us to demonstrate the following contribution:

1. We present an extension of the Hierarchical Agglomerative Clustering (HAC) approach [10] to the Bayesian setting by replacing the min-margin with a Bayesian acquisition function, such as BALD [22]. This is a novel combination that has not been tested before.
2. We show that performance of various methods is clustered based on particular building blocks of the method. Thus indicating that active learning methods may be tailored for each target application. A good example of Bayesian methods lies in estimating the distribution of neural networks which performed the best on Fake news but showed the same results on other datasets.
3. In a large scale study we demonstrate that warm start is often beneficial and even simple methods (such as single neural network with entropy acquisition) provide results competitive to, or better than, complex active learning schemes. We also show that cold start approaches reach the same or even worse results than

---

<sup>5</sup>The major part of the content in this part was taken from the paper sent to COLING conference.

the aforementioned warm start techniques. This is encouraging for practitioners that are interested in the methodology.

The paper is organized as follows. In Section 2.2, we briefly review all tested factors of batch active learning. The experimental setup of the sensitivity study is described in Section 2.3 and the results are reported in Section 2.4.

## 2.2 Batch Active Learning Methods

Throughout the paper, we will use the RoBERTa embedding [33] to represent documents in the feature space. RoBERTa is a modified BERT transformer model [12] that achieved comparable performance to BERT in [45] and outperformed all other embeddings in [43]. Representation of the  $k$ -th text document  $\mathbf{x}_k$  is calculated as the mean value from sentence embeddings of all sentences in the text.

The aim of document classification is to find a classifier  $\hat{\mathbf{y}} = \mathbf{y}(\theta, \mathbf{x})$  predicting the class label for each document representation  $\mathbf{x}$ . In a supervised setting, the classifier parameters are found on a training set  $\{\mathbf{x}_k, \mathbf{y}_k\}_{k=1}^K$  by matching the prediction  $\mathbf{y}(\theta, \mathbf{x}_k)$  with the provided label  $\mathbf{y}_k$  for each document. We are concerned with binary classification for simplicity, however, an extension to multiclass is straightforward.

We assume that for the full corpus of text documents  $\mathbf{X}$ , only a small initial set of labels  $\mathbf{Y}^{(0)}$ , is available. The full set  $\mathbf{X}$  is thus split into the labeled,  $\mathbf{X}^{(0)}$ , and unlabeled parts,  $\mathbf{X}_u^{(0)} = \mathbf{X} \setminus \mathbf{X}^{(0)}$ , the training set in the first round is then  $\mathbf{X}^{(0)}, \mathbf{Y}^{(0)}$ . Active learning is defined as a sequential extension of the training data set following a simple iterative strategy in Algorithm 2.2.

---

### Algorithm 2 General batch active learning

---

**Initialize:** set classifier structure  $\mathbf{y} = \mathbf{y}(\mathbf{x}, \theta)$ , iteration counter  $i = 0$ , initial data  $\mathbf{Y}^{(0)}, \mathbf{X}^{(0)}, \mathbf{X}_u^{(0)}$

**Iterate** until a stopping condition:

1. Train a classifier parameter  $\theta^{(i)}$  on  $\mathbf{Y}^{(i)}, \mathbf{X}^{(i)}$ , starting from  $\theta_{\text{init}}^{(i)}$
  2. Compute the value of a label for all documents in the unlabeled dataset,  $a_l = A(\mathbf{x}_l, \theta^{(i)}), \forall \mathbf{x}_l \in \mathbf{X}_u^{(i)}$
  3. Select a batch of documents,  $\tilde{\mathbf{X}} \subset \mathbf{X}$ , for labeling using  $a_l$
  4. Query labels  $\tilde{\mathbf{y}}$  for  $\tilde{\mathbf{X}}$  and extend the training set  $\mathbf{X}^{(i+1)} = \mathbf{X}^{(i)} \cup \tilde{\mathbf{X}}, \mathbf{y}^{(i+1)} = \mathbf{y}^{(i)} \cup \tilde{\mathbf{y}}, i = i + 1$ .
- 

The general algorithm can be specialized to many variants depending on various factors as specified next. We will introduce several choices labeled by the step in which they appear in Algorithm 2.2.

### 1a. Uncertainty representation:

The uncertainty can be represented by a maximum likelihood estimate, represented by a single network, or a Bayesian probabilistic estimate, represented typically by an ensemble of networks. We will consider the following options: **Single network** with a softmax output layer predicting the normalized probability of each class in one hot encoding. This probability is conditioned on the parameter, and thus captures only aleatoric uncertainty. Uncertainty in parameters is not represented. **Ensemble of networks**, represent uncertainty in parameters by different parameter value in each ensemble thus capturing both aleatoric and epistemic uncertainty. We consider two methods for generating the ensemble members: i) *MC dropout* [18], where ensemble members are generated by random draws of the dropout layers, and ii) *deep ensembles* [30] where ensembles are trained independently. Note that MC dropout is computationally much cheaper.

### 1b. Initialization of the training:

Each training in step 1 is a new task. However, the data set typically overlaps with the one from the previous iteration, which motivates the following strategies of reusing results from the previous iteration. The **Cold start** strategy is not reusing any information, the networks are initialized by random numbers,  $\theta_{\text{init}}^{(i)} = \mathcal{N}(0, \sigma)$ , where  $\sigma$  is given by the standard network init strategy, used most often [14, 10, 45]. The **Hot start** strategy reuses all information, setting the estimate from the previous iteration as a starting point,  $\theta_{\text{init}}^{(i)} = \theta^{(i-1)}$ , criticized in [23]. The **Warm start** strategy a combination of the above,  $\theta_{\text{init}}^{(i)} = \theta^{(i-1)} + \mathcal{N}(0, \sigma)$ , where  $\sigma$  is a hyper-parameter [43].

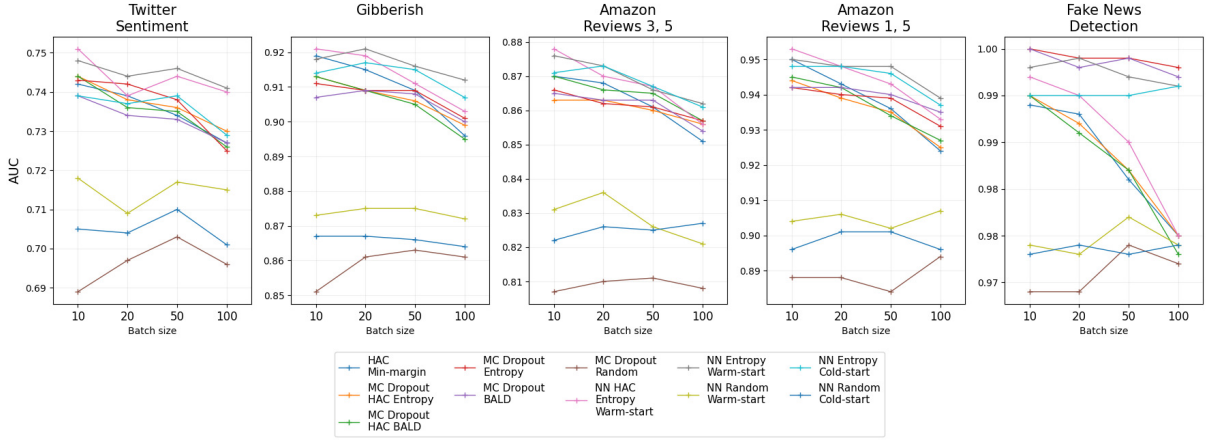


Figure 4: AUC metrics for seven active learning and two random strategies after 1000 acquired samples given datasets and batch size. Prediction of the MC dropout classifiers is an average over ensemble members.

## 2. Acquisition function:

Is a measure of the expected utility of the knowledge label,  $y_l$ , for each document,  $x_l$ , in the unlabeled data set. Different running index  $l$  is used to indicate that we operate on the unlabeled set. While many different utilities are proposed, we will study only the most popular ones. **Entropy** exists in two forms, entropy of the prediction  $a_l = \mathbb{H}(y|x_l, \theta)$  for a single network, or expected entropy  $a_l = \mathbb{E}_\theta \mathbb{H}(y|x_l, \theta)$  for ensembles. **BALD** is a mutual information metric,  $a_l = \mathbb{E}_\theta \mathbb{H}(y|x_l, \theta) - \mathbb{H}(y|x_l)$  that is meaningful only for the ensembles. **Min-margin** is a minimum difference between class predictions  $a_l = -\min_{c,d \in [1,C]} (y_c - y_d)$ , where  $C$  is the number of classes. Note carefully that the extreme of this criteria is equivalent to maximum entropy for binary classification with a single network.

## 3. Batch selection strategy:

When only one sample is to be selected, it is optimal to choose the one with maximum utility given by the acquisition function. However, the complexity of the maximum utility grows exponentially when the strategy has to select a batch of  $b$  documents for off-line labeling. Strategies that try to approximate this selection using greedy search [27] are still too computationally expensive for large batches. Therefore, we select two batch selection strategies that scale well with  $b$ . **Top** selects top  $b$  samples from sorted values of  $a_l$ . This approach may select samples close to each other, thus being redundant. **HAC** is a strategy based on the hierarchical clustering of  $a_l$  proposed in [10], and selecting top  $b$  samples from different clusters.

### Tested algorithm variants:

From the range of all possibilities, we will study the combinations that exist in the literature: HAC min-margin using cold start [10], MC dropout with Entropy and BALD criteria using warm start [18], warm start ensemble learning with Entropy and BALD called DEnFi [43], and conventional single-network with prediction Entropy a with warm start. If HAC is not in the name, the Top strategy is used.

Since HAC strategy is an orthogonal factor to the remaining ones, we propose its combination with other approaches, giving rise to: HAC Entropy for the single neural network and both ensemble methods (MC dropout and DEnFi) and HAC BALD for the ensembles.

## 2.3 Experiment Setup

The methods were compared on different datasets and different batch sizes. The used datasets are positive/negative tweets from the Tweets [19], Fake News Detection [2], two pairs of categories from Amazon Reviews [?], and Gibberish [?] datasets. From all datasets, we select from 10000 text documents (5000 text documents per category, selecting only two categories for binary classification, e.g. 1 and 5 in Amazon reviews), which were the

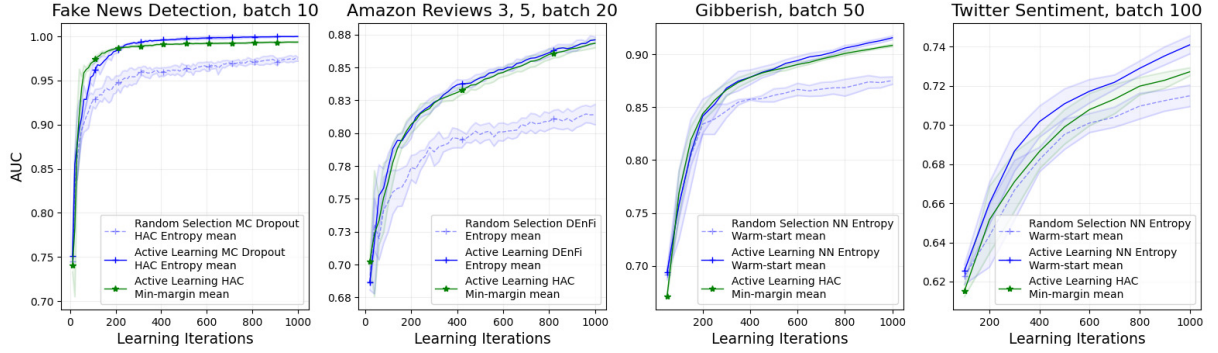


Figure 5: Evolution of the mean AUC with a growing number of requests for the best algorithms representative vs HAC Min-margin given the batch size and dataset. The uncertainty bounds are illustrated as one standard deviation from the mean value with respect to 5 runs. All algorithms were initially trained on 10 labeled text documents before sequential learning strategies were initialized

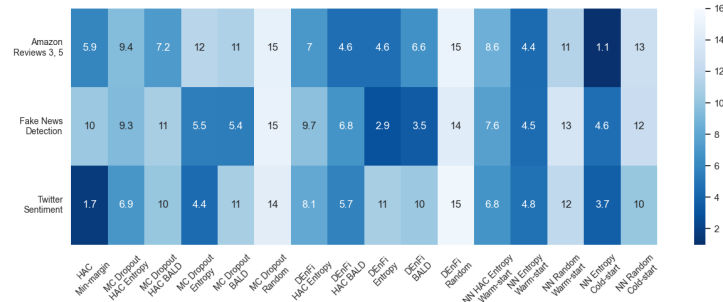


Figure 6: Aggregated mean rank for 14 tuples of learning algorithms and acquisition functions given Amazon Reviews 3,5, Fake News Detection, and Twitter Sentiment for 50 active learning iterations with batch size 20.



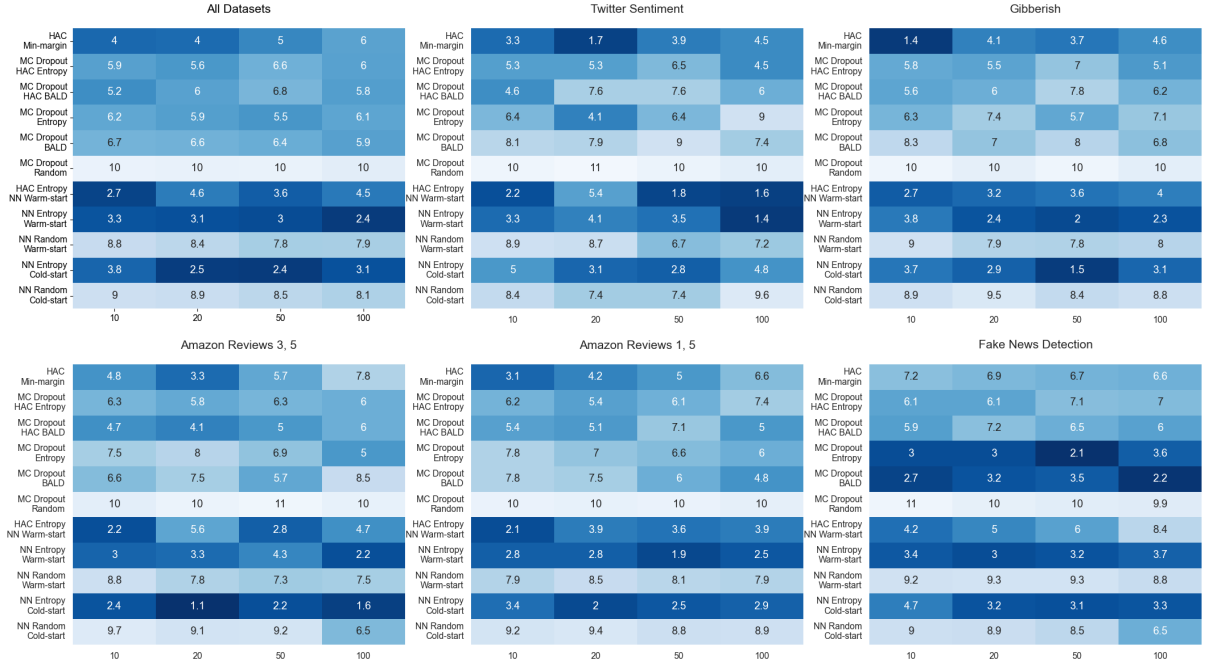


Figure 7: Aggregated rank for 7 active learning algorithms and two random strategies averaged over datasets as a function of different batch sizes.

initial 10000 documents of the datasets given categories. The only exception is Fake News Detection where only 4000 documents are available (2000 text documents per category).

#### Experiment Parameters:

Each active learning experiment was initialized by the training set  $\mathbf{X}^{(0)}, \mathbf{Y}^{(0)}$  of 10 samples. The active learning strategy was set to sample  $b$  samples with a discrete set of variants,  $b = 10, 20, 50, 100$ . The active learning was run until 1000 samples were labeled, i.e. making a different number of steps for each batch size (10 iterations for  $b = 100$ , 20 for  $b = 50$ , etc.). The batch selection follows the  $\epsilon$ -greedy approach [60], i.e. the samples selected by the acquisition function are accepted with probability  $\epsilon = \frac{\exp(l-3)}{\exp(l-3)+1}$ . A batch of random documents is selected for labeling if not accepted. The AUC is evaluated on the remaining part of the selected dataset (i.e. on the 9990 text documents in the first evaluation). The reported AUC values are averaged over 5 independent runs.

The initial number of epochs for the first iteration is 2500 for all algorithms. The same number is used for the cold start strategy in each iteration. The training of the warm start strategies is run for 150 epochs, with weights perturbation noise of variance  $\sigma = 0.3$  for both MC dropout and DENFi. Both DENFi and MC Dropout generate 5 ensemble members. The key difference is in computational complexity, while DENFi has to tune the parameters for each ensemble member, the MC dropout does it for only one network and generated ensemble members by 5 different realizations of the dropout mask.

#### Evaluation:

All algorithms were compared on the area under the curve (AUC) [15] on the test data (i.e. the documents not present in the training set). The algorithms were compared after 1000 acquired labels. The algorithms with smaller batch sizes thus benefited from higher number of retrains. To reduce the influence of stochastic initialization and training, the AUCs were run 5 times and averaged. Even then, the difference between the algorithms was sometimes marginal. To show the effect of various factors on the performance, we sorted the AUC and assigned a rank of each method accordingly. I.e. the best performing method has rank 1, second rank 2, etc. This approach allows the comparison of various methods across multiple datasets [11] using order statistics.

## 2.4 Results

**Parameter uncertainty:** The effect of parameter uncertainty (Bayesian approach) is the most costly to evaluate, due to the high computational demand of the ensemble approach (DENFi). Therefore, we have evaluated all algorithm variants only for batch size  $b = 20$ . The results are displayed in Figure 6. The advantage of the Bayesian approach is evident only for the Fake News dataset. This behavior is a result of a good neural network parameters distribution estimate. However, in other datasets, DENFi performed as good as a single neural network, and is not

worth the computational cost. As a result, the algorithm was omitted from subsequent large-scale studies.

A summary of the performance of all tested methods for various batch sizes is displayed in Figure 4 via AUC after 1000 samples for all methods, and via relative rank for all methods in Figure 7 averaged over ranks after each 100 label requests. Note that the datasets follow a similar pattern, with the exception of the Fake News data sets, where the parametric uncertainty (now represented only by the warm start MC dropout strategy) is beneficial, and HAC batch selection strategy has a negative effect (probably due to preference of large clusters).

**Acquisition functions:** Due to binary classification, the min-margin and entropy approaches coincide for a single network function. The difference between our generalization of Entropy and BALD methods for the ensemble techniques seems insignificant, Figure 7. HAC methods perform well for some cases, but a simple entropy approach shows more stable results for the majority of problems, batch sizes, and initializations.

**Initialization of the training:** The proposed modification on warm start strategies (HAC Entropy, NN Entropy, and NN Random warm start) are better or comparable in performance to the cold start (HAC Min-margin, NN Entropy, NN Random); this is achieved at a fraction of the training cost. This indicates that the additive noise is sufficient to avoid overfitting of the hot start [23].

**Batch selection strategy:** The HAC batch selection has a clear advantage for smaller batch sizes (10 and 20). This is consistent when comparing HAC and Top variants of all methods except Fake News Detection. Smaller batch sizes and the proposed generalization to warm start HAC outperforms the cold start approach in most of cases. The batch advantage diminishes for sizes of 50 and 100 where the top selection strategy achieves comparable (ensembles) or better (single NN) results. We project that the most informative samples in our datasets are clustered in small groups, hence the selection of a batch with a large enough size contains all important samples.

The contribution of this paper lies in a comparative study where we decomposed different algorithms into building blocks and generalized various approaches. For a better understanding of the methodology, we selected some approaches for demonstration. In Figure 5 a comparison of various active learning sequences is displayed. More specifically, the methods proposed by us to well studied HAC min-margin approach.

## 2.5 Conclusion

We have studied the influence of various factors (i. e. acquisition functions, batch sizes, neural networks initialization) of active learning algorithms and their performance on cover context integrity, fake news detection, and sentiment classification tasks. While complex algorithms such as deep ensembles (DEnFi and Dropout MC) sometimes achieve good performance (Fake News detection), the winner, on average, is the classical prediction entropy of a single neural network with a few proposed modifications like warm start. Although the performance of the warm start method can sometimes be the same as the cold start, the undeniable benefit is a lower computational cost. The selection of the batch size for annotation is also important. The agglomerative clustering improves performance for smaller batch sizes and may show better results than a more general method like entropy sampling.

## 3 Objectives

Despite the fact that we were able to push the limits of both single and batch active learning problems, there are still plenty of nuances that must be solved. The studied techniques work well on rather simpler neural networks. Hence, a clear further research objective is algorithms generalization to more complex neural networks. The larger the network is, the harder it is to perform a training/fine-tuning sequence, an additional training data selection and the uncertainty representations. Some extensions of the studied algorithms exist e.g. Variational Dropout [17] is a better extension of MC Dropout [18] for more complex methods. We also expect our research to follow with more sophisticated and architecture agnostic warm start neural network training.

## Part II

# Coreference Resolution

In this part we introduce a coreference resolution problem and a state of art description. In addition, we formulate our CR objective and the active learning to CR vision of the further research.

The grand research from Google [40] shows that the point-wise estimate of model parameters does not usually result in an optimal approach. The models uncertainty measurement through the estimate of the empirical model weights distribution has already shown remarkable results in the active learning fields both in Computer Vision [18], and NLP in such tasks like NER [48], [34], text classification [3] and other applications. We aim to utilize the knowledge of extended uncertainty algorithms (as it was shown in [44], [17], ...) throughout single/multi instances learning, hot/warm/cold start training and adapt it to the CR problem to improve the existing framework.

## 1 Introduction to Coreference Resolution Problem

Modern Coreference Resolution (CR) algorithms are combinations of sophisticated vector embeddings representing context and deep neural network superstructures that perform the coreference resolution itself. The set of existing models for natural language understanding (NLU) is vast. Arguably, one of the most prominent points in the history of such models is when the continuous bag-of-words and skipgram approaches were introduced [36]. At that point, machines started to learn the context surrounding particular words. Their vector representations acquired the ability to represent this context, meaning proximity of such vectors in terms of a metric of choice (L2, cosine/angular similarity) veritably described similarity of words or contexts. Still, models of these types were far from perfect, as they provided one with constant vectors per word for a pre-set vocabulary. Context-dependent representations with flexible vocabularies became available thanks to the introduction of the Transformer architecture [56] applied to the vocabulary formed not only by words but also by character ngrams constructed as meaningful parts of words. The power of the Transformer architecture lies in its encoding and decoding capability, improved by the self-attention mechanism, which learns to put stress on parts of text sequences. This gave birth to a lot of transformer-based language models such as the Bidirectional Encoder Representations from Transformers (BERT) [12], its fine-tuned variations [31], [33] and further models [42], [9]. To this date, SpanBERT [25] has proven to be one of the most efficient architectures for coreference resolution. Its crucial difference from the standard BERT model is that it learns to predict the content of masked spans of text, taking into account their beginnings and endings, omitting the ability of the base BERT model to predict preceding sentences. In contrast, BERT learns to predict the following sentence for each preceding one and attempts to infer individual masked tokens. Another model worthy of mention is a Longformer [5], whose architecture is based on transformers capable of processing long documents up to 4096 tokens [5]. The first end-to-end coreference resolution model was introduced in [32]. Its crucial difference from its predecessors was that it did not require preprocessing in the form of syntactic parsing or rule-based mention detection since the model can learn mention dependencies on its own to a forerunner-outperforming extent. The main idea of the model is to learn to score pairs of textual spans in such a way that takes into account, firstly, if these spans are entity mentions and, secondly, whether the pair is of type antecedent-descendant in terms of coreference. The NLU model of choice provides span representations. The goal is to assign to each span an antecedent span. [25] belongs to the state-of-the-art approaches which utilize the same structure on top of SpanBERT. One of the crucial drawbacks of the scoring approach is the choice of spans: sizes of relevant spans can be different, so a constant width of the window may not always be the right choice; spans can either overlap or be disjoint; if they overlap, the value of the overlap also becomes a hyperparameter. In addition to that, the number of scoring procedures is quadratic in complexity: each span has to be scored against every its counterpart. If the length of the document is large, the memory needed to store all entity mentions may become an issue (in [62] authors propose an incremental structure for the CR model, which needs a lot less memory for the price of a slight decrease in performance). While previous models are able to achieve decent results, their memory footprint is significant. The authors of [28] bypassed the need to create span representations, relying on a combination of bilinear functions applied on endpoint token representations. In addition to that, the new model is built on top of a Longformer encoder capable of processing long documents. Another work that reaches current state-of-the-art performance [13] also avoids span based approach. Even though the author uses RoBERTa embedding [33], the performance of the model is better than the one with Longformer. Author relies on attention mechanism that encodes tokens with the additional context based information from the surrounding. Hence, this is the argument of avoiding span based representation. When the CR classification is done, an additional span detection model is applied. Described approach helps to reduce the

computational complexity from  $O(n^4)$  to  $O(n^2)$ .

For languages other than English, the state of art for the CR is arguably even farther behind. For example, coreference resolution for Czech was attempted on the PCEDT dataset [20] in [39]. However, the overall performance of the approach did not reach the mark of 0.7 in terms of F-score (whereas English models show more than 0.8 F1 score [13]). In addition to that, no transformer-based NLU model was available at the time.

## 2 Objectives

Coreference resolution is a complex problem that involves both strong coding skills and a significant computational power access. Thus, the priority is to finalize state-of-the-art implementation that will be a basis for future uncertainty models integration. The subsequent action is to define the set of algorithms for uncertainty representation. In section 2 we showed that the uncertainty algorithms work well only in particular cases. When the parameters distribution is well-estimated the training and prediction performance reaches tremendous results. Further part of the work will be measured on combining the active learning strategy with such a heavy model as coreference resolution. The contribution in this part will be crucial. We expect to introduce hot and warm start approaches for heavy weight models that will immensely reduce learning time and training labels acquisition process.

The output of the study will allow us to use the architecture agnostic empirical model weights distribution estimate for better noise measurement, faster learning, and label distribution prediction, specifically for different types of CR task embedding superstructure.

## Summary

The research purposes and objectives are clear to us and lie in the final goal of the generalization of active learning methodology towards the coreference resolution problem. We underline that the choice of coreference resolution is based on more complex neural network structures, tremendous requirements for context understanding, vast usage potential, and personal interest. All these criteria make this type of model the perfect candidate for uncertainty algorithms integration and improvement. A deeper insight into the models' decision making process given the uncertainty is expected to allow for a new coreference resolution state-of-art threshold to be set. We believe that when the result of our work shows good model uncertainty representation for CR, the extension to other problems will be straightforward. We are also convinced that the old fashioned cold-start approach (retraining from scratch) became obsolete. Since our results of fine tuning the model with new labels showed incredible performance, the research in warm and hot start branches must be continued.

In this doctoral minimal thesis we showed a significant and strong progress in a study of uncertainty methods representations for active learning strategies. We were able to outperform state-of-art methods for different problems both for single instance and batch active learning. Moreover, we have successfully introduced text classification active learning approach with warm start methodology that speeds up massively model's training. We also explored the field of latest CR techniques and are aware of current top performing solutions. Described results and contributions let us operate with a sufficient amount of knowledge for the extension of CR problem by enriching it with uncertainty algorithms. Interestingly, the problem has no ultimate solution. It can be infinitely generalized (e.g choice of embedding, choice of CR superstructures, language independency, etc..)

We have submitted the application for GACR project. If the project application evaluation is successful, the research will be extended specifically to Czech language coreference resolution and another coreference resolution dataset that will be collected with newly studied active learning approaches.

## References

- [1] Fake news dataset, 2018. <https://www.kaggle.com/c/fake-news>, accessed on 20 february 2021.
- [2] Fake news detection dataset, 2018. <https://www.kaggle.com/jruvika/fake-news-detection>, accessed on 20 february 2021.
- [3] Bang An, Wenjun Wu, and Huimin Han. Deep active learning for text classification. In *Proceedings of the 2nd International Conference on Vision, Image and Signal Processing*, pages 1–6, 2018.
- [4] Mikel Artetxe and Holger Schwenk. Massively multilingual sentence embeddings for zero-shot cross-lingual transfer and beyond. *Transactions of the Association for Computational Linguistics*, 7:597–610, 2019.
- [5] I. Beltagy, M. E. Peters, and A. Cohan. Longformer: The Long-Document Transformer. *CoRR*, abs/2004.05150, 2020.
- [6] William H Beluch, Tim Genewein, Andreas Nürnberger, and Jan M Köhler. The power of ensembles for active learning in image classification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 9368–9377, 2018.
- [7] Sreyasee Das Bhattacharjee, Ashit Talukder, and Bala Venkatram Balantrapu. Active learning based news veracity detection with feature weighting and deep-shallow fusion. In *2017 IEEE International Conference on Big Data (Big Data)*, pages 556–565. IEEE, 2017.
- [8] Sreyasee Das Bhattacharjee, William J Tolone, and Ved Suhas Paranjape. Identifying malicious social media contents using multi-view context-aware active learning. *Future Generation Computer Systems*, 100:365–379, 2019.
- [9] D. Cer, Y. Yang, S. Kong, N. Hua, N. Limtiaco, R. St. John, N. Constant, M. Guajardo-Cespedes, S. Yuan, C. Tar, Y.-H. Sung, B. Strope, and R. Kurzweil. Universal sentence encoder, 2018.
- [10] Gui Citovsky, Giulia DeSalvo, Claudio Gentile, Lazaros Karydas, Anand Rajagopalan, Afshin Rostamizadeh, and Sanjiv Kumar. Batch active learning at scale. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems*, volume 34, pages 11933–11944. Curran Associates, Inc., 2021.
- [11] Janez Demšar. Statistical comparisons of classifiers over multiple data sets. *The Journal of Machine Learning Research*, 7:1–30, 2006.
- [12] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- [13] Vladimir Dobrovolskii. Word-level coreference resolution. *arXiv preprint arXiv:2109.04127*, 2021.
- [14] Liat Ein Dor, Alon Halfon, Ariel Gera, Eyal Shnarch, Lena Dankin, Leshem Choshen, Marina Danilevsky, Ranit Aharonov, Yoav Katz, and Noam Slonim. Active learning for bert: an empirical study. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 7949–7962, 2020.
- [15] Tom Fawcett. An introduction to roc analysis. *Pattern recognition letters*, 27(8):861–874, 2006.
- [16] Luciano Floridi and Massimo Chiriatti. Gpt-3: Its nature, scope, limits, and consequences. *Minds and Machines*, 30(4):681–694, 2020.
- [17] Yarin Gal, Jiri Hron, and Alex Kendall. Concrete dropout. *Advances in neural information processing systems*, 30, 2017.
- [18] Yarin Gal, Riashat Islam, and Zoubin Ghahramani. Deep bayesian active learning with image data. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 1183–1192. JMLR.org, 2017.
- [19] Alec Go, Richa Bhayani, and Lei Huang. Twitter sentiment classification using distant supervision. Technical report, Stanford, 2009. <http://help.sentiment140.com/for-students/> accessed on 26 june 2020.

- [20] Jan Hajič, Eva Hajičová, Jarmila Panevová, Petr Sgall, Ondřej Bojar, Silvie Cinková, Eva Fučíková, Marie Mikulová, Petr Pajas, Jan Popelka, Jiří Semecký, Jana Šindlerová, Jan Štěpánek, Josef Toman, Zdeňka Urešová, and Zdeněk Žabokrtský. Announcing prague czech-english dependency treebank 2.0. In *Proceedings of the 8th International Conference on Language Resources and Evaluation (LREC 2012)*, pages 3153–3160, Istanbul, Turkey, 2012. ELRA, European Language Resources Association.
- [21] Md Saqib Hasan, Rukshar Alam, and Muhammad Abdullah Adnan. Truth or lie: Pre-emptive detection of fake news in different languages through entropy-based active learning and multi-model neural ensemble. 2020.
- [22] Neil Houlsby, Ferenc Huszár, Zoubin Ghahramani, and Máté Lengyel. Bayesian active learning for classification and preference learning. *arXiv preprint arXiv:1112.5745*, 2011.
- [23] Peiyun Hu, Zachary C Lipton, Anima Anandkumar, and Deva Ramanan. Active learning with partial feedback. *arXiv preprint arXiv:1802.07427*, 2018.
- [24] Pieter Floris Jacobs, Gideon Maillette de Buy Wenniger, Marco Wiering, and Lambert Schomaker. Active learning for reducing labeling effort in text classification tasks. In *Benelux Conference on Artificial Intelligence*, pages 3–29. Springer, 2021.
- [25] M. Joshi, D. Chen, Y. Liu, D. S. Weld, L. Zettlemoyer, and O. Levy. Spanbert: Improving pre-training by representing and predicting spans. *CoRR*, abs/1907.10529, 2019.
- [26] Mohammad Khan, Didrik Nielsen, Voot Tangkaratt, Wu Lin, Yarin Gal, and Akash Srivastava. Fast and scalable bayesian deep learning by weight-perturbation in adam. In *International Conference on Machine Learning*, pages 2611–2620. PMLR, 2018.
- [27] Andreas Kirsch, Joost Van Amersfoort, and Yarin Gal. Batchbald: Efficient and diverse batch acquisition for deep bayesian active learning. *Advances in neural information processing systems*, 32, 2019.
- [28] Y. Kirstain, O. Ram, and O. Levy. Coreference Resolution without Span Representations. *CoRR*, abs/2101.00434, 2021.
- [29] Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable predictive uncertainty estimation using deep ensembles. *arXiv preprint arXiv:1612.01474*, 2016.
- [30] Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable predictive uncertainty estimation using deep ensembles. In *Advances in neural information processing systems*, pages 6402–6413, 2017.
- [31] Z. Lan, M. Chen, S. Goodman, K. Gimpel, P. Sharma, and R. Soricut. Albert: A lite bert for self-supervised learning of language representations, 2020.
- [32] K. Lee, L. He, M. Lewis, and L. Zettlemoyer. End-to-end neural coreference resolution. *CoRR*, abs/1707.07045, 2017.
- [33] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*, 2019.
- [34] David Lowell, Zachary C Lipton, and Byron C Wallace. Practical obstacles to deploying active learning. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 21–30, 2019.
- [35] Katerina Margatina, Loic Barrault, and Nikolaos Aletras. Bayesian active learning with pretrained language models. *arXiv preprint arXiv:2104.08320*, 2021.
- [36] Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*, 2013.
- [37] Tomas Mikolov, Edouard Grave, Piotr Bojanowski, Christian Puhersch, and Armand Joulin. Advances in pre-training distributed word representations. In *Proceedings of the International Conference on Language Resources and Evaluation (LREC 2018)*, 2018.

- [38] Rishabh Misra. News category dataset, 06 2018. <https://www.kaggle.com/rmisra/news-category-dataset>, accessed on 26 june 2020.
- [39] Michal Novák. Coreference resolution system not only for czech. In Jaroslava Hlaváčová, editor, *Proceedings of the 17th Conference on Information Technologies - Applications and Theory (ITAT 2017), Martinské hole, Slovakia, September 22-26, 2017*, volume 1885 of *CEUR Workshop Proceedings*, pages 193–200. CEUR-WS.org, 2017.
- [40] Yaniv Ovadia, Emily Fertig, Jie Ren, Zachary Nado, David Sculley, Sebastian Nowozin, Joshua V Dillon, Balaji Lakshminarayanan, and Jasper Snoek. Can you trust your model’s uncertainty? evaluating predictive uncertainty under dataset shift. *arXiv preprint arXiv:1906.02530*, 2019.
- [41] Sumanth Prabhu, Moosa Mohamed, and Hemant Misra. Multi-class text classification using bert-based active learning. *arXiv preprint arXiv:2104.14289*, 2021.
- [42] A. Radford and K. Narasimhan. Improving language understanding by generative pre-training. 2018.
- [43] Marko Sahan, Vaclav Smidl, and Radek Marik. Active learning for text classification and fake news detection. In *2021 International Symposium on Computer Science and Intelligent Controls (ISCSIC)*, pages 87–94. IEEE, 2021.
- [44] Marko Sahan, Vaclav Smidl, and Radek Marik. Batch active learning for text classification and sentiment analysis. 2022.
- [45] Christopher Schröder, Andreas Niekler, and Martin Potthast. Uncertainty-based query strategies for active learning with transformers. *CoRR*, abs/2107.05687, 2021.
- [46] Erich Schweighofer, Andreas Rauber, and Michael Dittenbach. Automatic text representation, classification and labeling in european law. In *Proceedings of the 8th international conference on Artificial intelligence and law*, pages 78–87, 2001.
- [47] Claude E Shannon. A mathematical theory of communication. *Bell system technical journal*, 27(3):379–423, 1948.
- [48] Yanyao Shen, Hyokun Yun, Zachary C Lipton, Yakov Kronrod, and Animashree Anandkumar. Deep active learning for named entity recognition. *arXiv preprint arXiv:1707.05928*, 2017.
- [49] Yusuxke Shibata, Takuya Kida, Shuichi Fukamachi, Masayuki Takeda, Ayumi Shinohara, Takeshi Shinohara, and Setsuo Arikawa. Byte pair encoding: A text compression scheme that accelerates pattern matching. Technical report, Citeseer, 1999.
- [50] Shaden Smith, Mostofa Patwary, Brandon Norick, Patrick LeGresley, Samyam Rajbhandari, Jared Casper, Zhun Liu, Shrimai Prabhunoye, George Zerveas, Vijay Korthikanti, et al. Using deepspeed and megatron to train megatron-turing nl-g 530b, a large-scale generative language model. *arXiv preprint arXiv:2201.11990*, 2022.
- [51] Jasper Snoek, Yaniv Ovadia, Emily Fertig, Balaji Lakshminarayanan, Sebastian Nowozin, D Sculley, Joshua Dillon, Jie Ren, and Zachary Nado. Can you trust your model’s uncertainty? evaluating predictive uncertainty under dataset shift. In *Advances in Neural Information Processing Systems*, pages 13969–13980, 2019.
- [52] Harold Somers. An introduction to machine translation. 1992.
- [53] Romal Thoppilan, Daniel De Freitas, Jamie Hall, Noam Shazeer, Apoorv Kulshreshtha, Heng-Tze Cheng, Alicia Jin, Taylor Bos, Leslie Baker, Yu Du, et al. Lamda: Language models for dialog applications. *arXiv preprint arXiv:2201.08239*, 2022.
- [54] Simon Tong and Daphne Koller. Support vector machine active learning with applications to text classification. *Journal of machine learning research*, 2(Nov):45–66, 2001.
- [55] Lukas Ulrych and Vaclav Smidl. Deep ensemble filter for active learning. Technical Report 2383, Institute of Information Theory and Automation, 2020.



- [56] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. *arXiv preprint arXiv:1706.03762*, 2017.
- [57] Q. Wang, Z. Mao, B. Wang, and L. Guo. Knowledge graph embedding: A survey of approaches and applications. *IEEE Transactions on Knowledge and Data Engineering*, 29(12):2724–2743, 2017.
- [58] Sida I Wang and Christopher D Manning. Baselines and bigrams: Simple, good sentiment and topic classification. In *Proceedings of the 50th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 90–94, 2012.
- [59] Xinyu Wang, Yong Jiang, Nguyen Bach, Tao Wang, Zhongqiang Huang, Fei Huang, and Kewei Tu. Automated concatenation of embeddings for structured prediction. *arXiv preprint arXiv:2010.05006*, 2020.
- [60] Christopher John Cornish Hellaby Watkins. Learning from delayed rewards. 1989.
- [61] Max Welling and Yee W Teh. Bayesian learning via stochastic gradient langevin dynamics. In *Proceedings of the 28th International Conference on Machine Learning (ICML-11)*, pages 681–688, 2011.
- [62] P. Xia, J. Sedoc, and B. Van Durme. Incremental neural coreference resolution in constant memory, 2020.