

# Are All the Datasets in Benchmark Necessary? A Pilot Study of Dataset Evaluation for Text Classification

Anonymous ACL submission

## Abstract

In this paper, we ask the research question if all the datasets in the benchmark are necessary. We approach this by first characterizing the distinguishability of datasets when comparing different systems. Experiments on 9 datasets and 36 systems show that several existing benchmark datasets contribute little to discriminating top-scoring systems, while those less used datasets exhibit impressive discriminative power. We further, taking the text classification task as a case study, investigate the possibility of predicting dataset discrimination based on its properties (e.g., average sentence length). Our preliminary experiments promisingly show that given a sufficient number of training experimental records, a meaningful predictor can be learned to estimate dataset discrimination over unseen datasets.

We released all related code at [Github](#)<sup>1</sup> and a new benchmark dataset for text classification based on our observations.

## 1 Introduction

In natural language processing (NLP) tasks, there are often datasets that we use as benchmarks against which to evaluate machine learning models, either explicitly defined such as GLUE (Wang et al., 2018) and XTREME (Hu et al., 2020a) or implicitly bound to the task (e.g., DPedia (Zhang et al., 2015) has become a default dataset for the evaluation of text classification systems). Given this mission, one important feature of a good benchmark dataset is the ability to statistically differentiate diverse systems (Bowman and Dahl, 2021). With the large pre-trained model (Devlin et al., 2018; Lewis et al., 2019) constantly updating the best performance of NLP tasks, the performances of many of them have reached a plateau (Zhong et al., 2020; Fu et al., 2020). In other words, it is challenging to discriminate a better model using existing datasets (Wang

<sup>1</sup><https://github.com/annonnlp-demo/acl-v2>

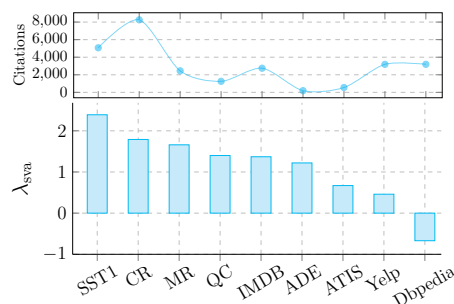


Figure 1: Illustrate different datasets’ distinguishing ability w.r.t top-scoring systems characterized by our measure  $\log(\lambda_{sva})$  on text classification and their corresponding citations.

et al., 2019a). In this context, we ask the question: *are all benchmark’s datasets necessary?* We use the text classification task as a case study and try to answer the following two sub-questions:

**RQ1:** *How can we quantify the distinguishing ability of benchmark datasets?* To answer this question, we first design measures with varying calculation difficulties (§4) to judge datasets’ discrimination ability based on top-scoring systems’ performances. By exploring correlations among different measures, we then evaluate how reliable a dataset’s discrimination is when discrimination is calculated solely based on overall results that top-scoring systems have achieved, and generalize this measure to other NLP tasks. Fig. 1 illustrates how different text classification datasets are ranked (the bottom one) based on measures devised in this work (a smaller value suggests lower discrimination) and the corresponding citations of these datasets (the upper one). One can observe that: (i) The highly-cited dataset DBpedia (Zhang et al., 2015) (more than 3,000 times since 2015) shows the worst discriminative power. (ii) By contrast, dataset like ADE (Gurulingappa et al., 2012) (less than 200 times since 2012) does better in distinguishing top-scoring systems. This phenomenon shows the significance of quantifying the discrim-

067 inative ability of datasets: it can not only help us  
068 to **eliminate** those with lower discrimination from  
069 *commonly-used datasets* (e.g., DBpedia), but also  
070 help us to **recognize** the missing pearl in *seldom*  
071 *used datasets* (e.g., ADE and ATIS (Hemphill et al.,  
072 1990)).

073 **RQ2:** *Can we try to predict the discriminative*  
074 *power of the model?* Given a dataset, we investi-  
075 gate if we can judge its ability to distinguish models  
076 based on its characteristics (e.g., average sentence  
077 length), which is motivated by the scenario where  
078 a new dataset has just been constructed without  
079 sufficient top-scoring systems to calculate discrim-  
080 ination defined in RQ1. To answer this question,  
081 inspired by recent literature on performance pre-  
082 diction (Domhan et al., 2015; Turchi et al., 2008;  
083 Birch et al., 2008; Xia et al., 2020; Ye et al., 2021),  
084 we conceptualize this problem as a *discrimination*  
085 *regression task*. We define 11 diverse features to  
086 characterize a text classification dataset and regress  
087 its discrimination scores using different parame-  
088 terized models. Preliminary experiments (§5.4)  
089 indicate that a meaningful regressor can be learned  
090 to estimate the discrimination of unseen datasets  
091 without actual training using top-scoring systems.

092 We brief **takeaways** in this work based on our  
093 observations:

094 (1) In regard to multitask benchmark datasets,  
095 empirical results show that following datasets  
096 struggle at discriminating current top-scoring sys-  
097 tems: STS-B and SST-2 from GLUE (Wang  
098 et al., 2019b); BUCC and PAWX-X from XTREME,  
099 which is consistent with the concurrent work  
100 (Ruder et al., 2021) (§4.3.2).

101 (2) In regard to single-task benchmark datasets,  
102 for Chinese Word Segmentation task, there are  
103 multiple datasets (MSR, CityU, CTB) (Tseng  
104 et al., 2005; Jin and Chen, 2008) that exhibit much  
105 worse discriminative ability, suggesting that: fu-  
106 ture works on this task are encouraged to either  
107 (i) adopt other datasets to evaluate their systems  
108 or (ii) at least make significant test <sup>2</sup> if using these  
109 datasets. Similar observations happen in the dataset  
110 CoNLL-2003 (Sang and De Meulder, 2003) from  
111 Named Entity Recognition task and MultiNLI  
112 (Williams et al., 2017) from natural language infer-  
113 ence task (§4.3.2).

114 (3) Some seldom used datasets such as ADE from  
115 text classification are actually better at distinguish-

116 ing top-performing systems, which highlights an  
117 interesting and necessary future direction: *how to*  
118 *identify infrequently-used but valuable (better dis-*  
119 *crimination) datasets for NLP tasks, especially in*  
120 *the age of dataset’s proliferation?*<sup>3</sup> (§4.2)

121 (4) Quantifying a dataset’s discrimination (w.r.t  
122 top-scoring systems) by calculating the statistical  
123 measures (defined in §4.1.2) from leaderboard’s  
124 results is a straightforward and effective way. But  
125 for those datasets without rich leaderboard results,<sup>4</sup>  
126 predicting the discrimination based on datasets’  
127 characteristics would be an promising direction  
128 (§4.3.1).

129 Our **contributions** can be summarized as:

130 (1) We try to quantify the discrimination abil-  
131 ity for datasets by designing two variance-based  
132 measures. (2) We systematically investigate 4 text  
133 classification models on 9 datasets, providing the  
134 newest baseline performance for those seldom used  
135 datasets. We released the code and all the uni-  
136 formly formatted datasets at <https://github.com/annonnlp-demo/acl-V2> (3) We study  
137 several popular NLP benchmarks, including GLUE,  
138 XTREME, NLI, and so on. Some valuable sugges-  
139 tions and observations will make research easier.  
140

## 141 2 Related Work

142 **Benchmarks for NLP** In order to conveniently  
143 keep themselves updated with the research  
144 progress, researchers recently are actively build-  
145 ing evaluation benchmarks for diverse tasks so  
146 that they could make a comprehensive compari-  
147 son of systems, and use a leaderboard to record the  
148 evolving process of the systems of different NLP  
149 tasks, such as SQuAD (Rajpurkar et al., 2016),  
150 GLUE (Wang et al., 2018), XTREME (Hu et al.,  
151 2020a), GEM (Gehrmann et al., 2021) and GE-  
152 NIE (Khashabi et al., 2021). Despite their utility,  
153 more recently, Bowman and Dahl (2021) highlight  
154 that unreliable and biased systems score so highly  
155 on standard benchmarks that there is little room for  
156 researchers who develop better systems to demon-  
157 strate their improvements. In this paper, we make  
158 a pilot study on meta-evaluating benchmark evalua-  
159 tion datasets and quantitatively characterize their  
160 discrimination in different top-scoring systems.

<sup>3</sup><https://paperswithcode.com/datasets>

<sup>4</sup>The measure can keep updated as the top-scoring sys-  
tems of the leaderboard evolves, which can broaden its practi-  
cal applicability

<sup>2</sup>We randomly select 10 recently published papers (from  
ACL/EMNLP) that utilized these datasets and found only 2 of  
them perform significant test.

**Performance Prediction** Performance prediction is the task of estimating a system’s performance without the actual training process. With the recent booming of the number of machine learning models (Goodfellow et al., 2016) and datasets, the technique of performance prediction become rather important when applied to different scenarios ranging from early stopping training iteration (Kolachina et al., 2012), architecture searching (Domhan et al., 2015), and attribution analysis (Birch et al., 2008; Turchi et al., 2008). In this work, we aim to calculate a dataset’s discrimination without actual training top-scoring systems on it, which can be formulated as a performance prediction problem.

### 3 Preliminaries

#### 3.1 Task and Dataset

Text classification aims to assign a label defined beforehand to a given input document. In the experiment, we choose nine datasets, and their statistics can be found in the Appendix A.

- **IMDB** (Maas et al., 2011) consists of movie reviews with binary classes.
- **Yelp** (Zhang et al., 2015) is a part of the Yelp Dataset Challenge 2015 data, which is collected from Yelp.
- **CR** (Hu and Liu, 2004) is a product review dataset with binary classes.
- **MR** (Pang and Lee, 2005) is a movie review dataset collected from Rotten Tomatoes.
- **SST1** (Socher et al., 2013) is collected from HTML files of Rotten Tomatoes reviews with fully labeled parse trees.
- **DBpedia14** (Zhang et al., 2015) is a dataset for ontology classification that is collected from DBpedia 2014.
- **ATIS** (Hemphill et al., 1990) is an intent detection dataset that contains audio recordings of flight reservations.
- **QC** (Li and Roth, 2002) is a question classification dataset.
- **ADE** (Gurulingappa et al., 2012) is a subset of “Adverse Drug Reaction Data”.

#### 3.2 Model

We re-implement 4 top-scoring systems with typical neural architectures for each dataset.<sup>5</sup> The

<sup>5</sup>We mainly focus on neural network-based models, since most top-scoring systems in the leaderboard are based on deep learning.

brief introduction of the four models is as follows.

- **LSTM** (Hochreiter and Schmidhuber, 1997) is a widely used sentence encoder. To get left-to-right and right-to-left features, here, we adopt the bidirectional LSTM.
- **LSTMAtt** is proposed by Lin et al. (2017) that designed the self-attention mechanism to extract different aspects of features for a sentence.
- **BERT** (Devlin et al., 2018) utilizes the LSTM as the sentence encoder and gets word representation by BERT.
- **CNN** (LeCun and Bengio, 1995) extracts the sentence representation on the sequence of word representations.

Except for BERT, the other three models (e.g. LSTM) are initialized by GloVe (Pennington et al., 2014) or Word2Vec (Mikolov et al., 2013) pre-trained word embeddings. When the performance of the model on the dev set doesn’t improve within 20 epochs, the training will be stopped, and the best performing model will be kept. More detailed model parameter settings can be found in the Appendix B.

### 4 How to Characterize Discrimination?

To achieve this goal, we design measures based on the performance of different models for a dataset.

#### 4.1 Measures

The general idea of the measure designing is to judge dataset’s distinguishing ability based on the performances that top-performing systems have achieved on it.<sup>6</sup> Specifically, given a dataset  $D$  together with  $k$  top-scoring model *performance list*  $\mathbf{v} = [v_1, \dots, v_k]$ , we define the following measures.

##### 4.1.1 Performance Variance

We use the standard deviation to quantify the degree of variation or dispersion of a set of performance values. A larger value of  $\lambda_{\text{var}}$  suggests that the discrimination of the given dataset is more significant.  $\lambda_{\text{var}}$  can be defined as:

$$\lambda_{\text{var}} = \text{Std}(\mathbf{v}), \quad (1)$$

where  $\text{Std}(\cdot)$  is the function to compute the standard deviation. Assume that the performance list ( $k = 3$ ) on dataset  $D$  is  $\mathbf{v} = [88, 92, 93]$ , we can get  $\lambda_{\text{var}} = 2.65$ .

<sup>6</sup>A dataset’s discrimination is defined w.r.t top-scoring models from a leaderboard, keeping itself updated with systems’ evolution.

#### 4.1.2 Scaled Performance Variance

For the above measure, it can only reflect the variances of the performance of different models, without considering whether the model’s performance is close to the upper limit (e.g., 100% accuracy) on a given data set. To address this problem, we defined a modified variance by scaling  $\lambda_{\text{var}}$  with the difference between the upper limit performance  $u$  and average performance  $\text{Avg}(\mathbf{v})$  of  $\mathbf{v}$ .

$$\lambda_{\text{sva}} = \lambda_{\text{var}}(u - \text{Avg}(\mathbf{v})). \quad (2)$$

In practice,  $u$  can be defined flexibly based on tasks’ metrics. For example, in text classification task,  $u$  could be 100% (w.r.t F1 or accuracy), while in summarization task,  $u$  could be the results of oracle sentences (w.r.t ROUGE). Intuitively, given a performance list on text classification dataset:  $\mathbf{v} = [88, 92, 93]$ , we can obtain the  $\lambda_{\text{sva}} = 23.81$ .

#### 4.1.3 Hit Rate

The previous two measures quantify dataset’s discriminative ability w.r.t  $k$  top-performing systems in an *indirect* way (i.g, solely based on the overall results of different models). However, sometimes, small variance does not necessarily mean that the dataset fail to distinguish models, as long as the difference between models is statistically significant. To overcome this problem, we borrow the idea of bootstrap-based significant test (Koehn, 2004) and define the measure *hit rate*, which quantify the degree to which a given dataset could successfully differentiate  $k$  top-scoring systems.

Specifically, we take all  $\binom{k}{2}$  pairs of systems ( $m_i$  and  $m_j$ ) and compare their performances on a subset of test samples  $D_t$  that is generated using paired bootstrap re-sampling. Let  $v_i(D) > v_j(D)$  be the performance of  $m_1$  and  $m_2$  on the full test set, we define  $P(m_i, m_j)$  as the frequency of  $v_i(D_t) > v_j(D_t)$  over all  $T$  times of re-sampling ( $t = 1, \dots, T$ ) Then we have

$$\lambda_{\text{hit}} = \frac{1}{\binom{k}{2}} \sum P(m_i, m_j) \quad (3)$$

**Metric Comparison** The first two metrics, performance variance and scaled performance variance, are relative easily to obtain since they only require holistic performances of different top-scoring models on a given dataset, which can be conveniently collected from existing leaderboards. By contrast, although the metric *hit rate* can directly reflect dataset’s ability in discriminating diverse

systems, its calculation not only require more fine-grained information of system prediction but also complicated bootstrap re-sampling process.

#### 4.2 Exp-I: Exploring Correlation Between Variance and Hit Rate

The goal of this experiment is to investigate the reliability of the variance-based discrimination measures (e.g.,  $\lambda_{\text{sva}}$ ), which are easier to obtain, by calculating its correlation with significant test-based measure  $\lambda_{\text{hit}}$ , which is costly to get. Since the implementation of  $\lambda_{\text{hit}}$  relies on the bootstrap-based significant test, we choose text classification as the tested and re-implement 4 classification models (defined in Sec. 3.2) on 9 datasets. The performance and the distinction degree on the 9 text classification dataset are shown in Tab. 1.  $\lambda_{\text{var}}$  and  $\lambda_{\text{sva}}$  measures are designed based on performance variance, even if BERT always achieves the best performance on the same dataset, it will not affect the observed results from our experiments.

**Correlation measure** Here, we adopt the Spearman rank correlation coefficient (Zar, 1972) to describe the correlation between our variance-based measures and the hit rate measure  $\lambda_{\text{hit}}$ .

$$S_\lambda = \text{Spearman}(q, \lambda_{\text{hit}}), \quad (4)$$

where the  $q$  can be  $\lambda_{\text{var}}$  or  $\lambda_{\text{sva}}$ .

**Result** (1)  $\lambda_{\text{var}}$  and  $\lambda_{\text{sva}}$  are strong correlative ( $S_\lambda > 0.6$ ) with  $\lambda_{\text{hit}}$  respectively, which suggests that variance-based metrics could be a considerably reliable alternatives of significant test-based metric. (2)  $\text{Spearman}(\lambda_{\text{var}}, \lambda_{\text{hit}}) > \text{Spearman}(\lambda_{\text{sva}}, \lambda_{\text{hit}})$ , which indicate that comparing with  $\lambda_{\text{sva}}$ , dataset discrimination characterized by  $\lambda_{\text{var}}$  is more acceptable for  $\lambda_{\text{hit}}$ . The reason can be attributed to that the designing of the measure  $\lambda_{\text{hit}}$  does not consider the upper limit of the model’s performance. (3) DPdedia and Yelp are commonly used text classification datasets, while they have the worst ability to discriminate the top-scoring models since they get the lowest value of  $\lambda_{\text{var}}$  and  $\lambda_{\text{sva}}$ . By contrast, these two seldom used datasets ADE and ATIS show the better discriminative ability.

#### 4.3 Exp-II: Evaluation of Other Benchmarks

##### 4.3.1 Popular Benchmark Datasets

We also investigate how benchmark datasets from other NLP task perform using two devised measures. Specifically, we collected three single-task

Method	BERT	LSTMAAttr	LSTM	CNN	$\lambda_{hit}$	$\lambda_{var}$	$\lambda_{sva}$
SST1	54.12	43.80	47.60	44.80	0.88	4.65	243.56
CR	91.75	83.25	82.50	84.25	0.91	4.27	62.17
MR	85.55	79.92	79.80	82.00	0.86	2.69	48.83
QC	97.19	90.36	89.96	92.17	0.92	3.32	25.18
IMDB	93.34	89.45	89.65	87.81	0.87	2.33	23.18
ADE	93.48	92.90	92.65	89.54	0.78	1.77	13.90
ATIS	97.64	97.42	97.31	94.62	0.78	1.42	4.63
Yelp	97.52	96.60	96.60	95.46	0.81	0.84	2.91
DPedia	99.27	99.01	99.05	98.75	0.68	0.22	0.21
Spearman						0.83	0.73

Table 1: Illustration the 4 models’ performance and discrimination degree (characterized by  $\lambda_{hit}$ ,  $\lambda_{var}$ , and  $\lambda_{sva}$ ) on 9 text classification datasets. The two correlation coefficients pass the significance test ( $p < 0.05$ ).  $\lambda_{var}$  and  $\lambda_{sva}$  measures are designed based on performance variance.

and two multitask benchmarks. For the single-task benchmarks, we collect the top-performing models in a specific period for each dataset, provided by Paperswithcode<sup>7</sup>. For the multitask benchmarks, here, the GLUE<sup>8</sup> and XTREME<sup>9</sup> are considered in this work. Since Paperswithcode provided 5 models for each dataset in most case, for fairness and uniformity, we keep top-5 models for both single-task and multitask benchmark datasets.

**Named Entity Recognition (NER)** aims to identify named entities of an input text, for which we choose 5 top-scoring systems on 6 datasets and collect results from Paperswithcode.

**Chinese Word Segmentation (CWS)** aims to detect the boundaries of Chinese words in a sentence. We select 5 top-scoring systems on 8 datasets and collect results from Paperswithcode.

**Natural Language Inference (NLI)** targets at predicting whether a premise sentence can infer the hypothesis sentence. We select 5 top-performing models on 4 datasets from Paperswithcode.

**GLUE** (Wang et al., 2019b) covers 9 sentence- or sentence-pair tasks with different dataset sizes, text genres, and degrees of difficulty. Fig. 2-(a) shows the tasks/datasets that are considered in GLUE.

**XTREME** (Hu et al., 2020b) is the first benchmark that evaluates models across a wide variety of languages and tasks. The tasks/datasets that are covered by XTREME are shown in Fig. 2-(b).

### 4.3.2 Results and Analysis

Fig. 2 shows the results of dataset quality measure by  $\lambda_{var}$  and  $\lambda_{sva}$ . We detail several main observations:

<sup>7</sup><https://paperswithcode.com/>

<sup>8</sup><https://gluebenchmark.com/>

<sup>9</sup><https://sites.research.google/xtreme>

- $\lambda_{var}$  and  $\lambda_{sva}$  have consistent evaluation results for both single-task (CWS, NER, NLI) and multitask (GLUE, XTREME) benchmarks.
- For the XTREME benchmark, BUCC and PAWSX have lowest  $\lambda_{var}$  and  $\lambda_{sva}$ , which suggest that they are hardly to discriminate the top-performing systems. Moreover, these two data sets will be removed from the new version of the XTREME leaderboard called XTREME-R (Ruder et al., 2021). This consistent observation also shows the effectiveness of our measure.
- For GLUE benchmark, CoLA, QQP, and RTE datasets have the excellent ability to distinguish different top-scoring models (with higher  $\lambda_{var}$  and  $\lambda_{sva}$ ), while the SST-2 and STS-B datasets have the opposite conclusions.
- For CWS benchmarks, there is a larger gap between the value of  $\lambda_{var}$  and  $\lambda_{sva}$ , which indicate that the performance of top-scoring models considered are close to 100%. Furthermore, MSR, CityU and CTB are not suitable as benchmarks since they have poor discrimination ability with  $\lambda_{sva} < 0$ . So as MultiNLI for NLI task.
- CoNLL 2003 is a widely used NER dataset, but it is the lowest quality dataset under our dataset quality measure. The reason can be attributed to contain much annotation errors (Fu et al., 2020) in the CoNLL 2003 dataset, which makes its performance reach the bottleneck.

## 5 Can we Predict Discrimination?

Although metrics  $\lambda_{var}$ ,  $\lambda_{sva}$  ease the burden for us to calculate the datasets’ discrimination, one major limitation is: given a new dataset without results from leaderboards, we need to train multiple top-scoring systems and calculate corresponding results on it, which is computationally expensive. To alle-

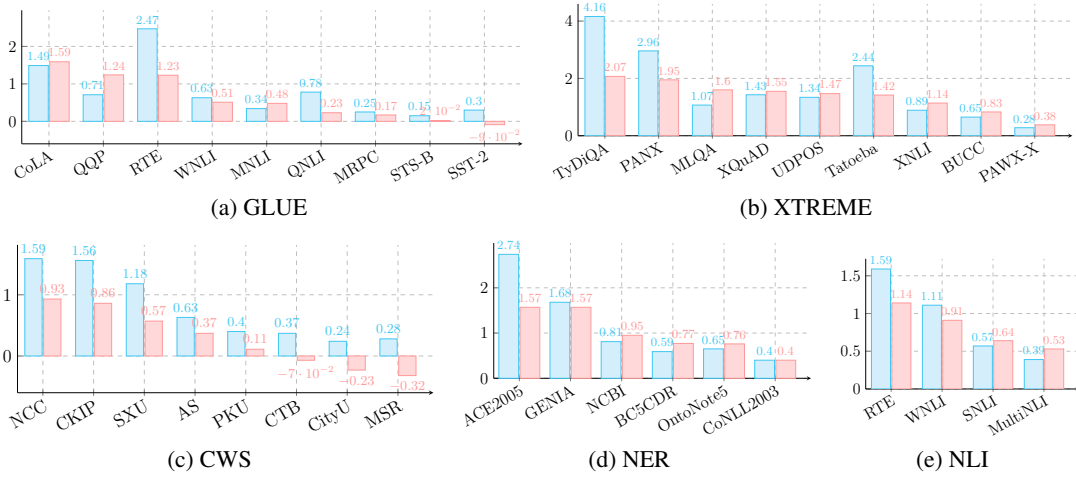


Figure 2: The dataset discrimination characterized by  $\lambda_{var}$  (blue) and  $\lambda_{sva}$  (pink) on five popular NLP benchmarks.

415 viate this problem, in this section, we focus on text  
 416 classification task and investigate the possibility of  
 417 estimating datasets’ discrimination solely based on  
 418 their characteristics without actual training systems  
 419 on them.

## 420 5.1 Task Formulation

### 421 5.1.1 Regression-based Task Formulation

422 We formulate it as a performance prediction prob-  
 423 lem (Birch et al., 2008; Xia et al., 2020; Ye et al.,  
 424 2021). Formally, we refer to  $\mathcal{M}$ ,  $D^{tr}$ ,  $D^{te}$ ,  $\mathcal{S}$   
 425 as the machine learning system, training data, test  
 426 data and training strategy respectively. The goal of  
 427 performance prediction is to estimate actual perfor-  
 428 mance  $y$  without actual training by using features  
 429 of  $\mathcal{M}$ ,  $D^{tr}$ ,  $D^{te}$ , and  $\mathcal{S}$ .

$$430 \hat{y} = \hat{f}(\Phi_{\mathcal{M}}, \Phi_{D^{tr}}, \Phi_{D^{te}}, \Phi_{\mathcal{S}}; \hat{\Theta}) \quad (5)$$

431 where  $\hat{y}$  denotes estimated prediction and  $\Phi(\cdot)$  is  
 432 a feature extractor. Following Xia et al. 2020, we  
 433 only use the features of the datasets as variables and  
 434 adapt it to our discriminative prediction scenario,  
 435 we can obtain:

$$436 \hat{\lambda} = \hat{f}(\Phi_{D^{tr}}, \Phi_{D^{te}}; \hat{\Theta}) \quad (6)$$

437 where  $\hat{\lambda}$  denotes predicted variance defined in  
 438 §4.1.2 such as  $\lambda_{var}$  or  $\lambda_{sva}$ .

### 439 5.1.2 Ranking-based Task Formulation

440 Instead of only regressing one dataset’s quality,  
 441 we also care about the quality ranking of dif-  
 442 ferent datasets w.r.t discriminating systems in a  
 443 task. Therefore, we also formulate it as a listwise  
 444 LTR(learning to rank) task where a model takes

445 individual lists as instances, to predict the rank of  
 446 element among the list (Liu, 2011). Given a set  
 447 of  $n$  datasets  $d = \{d_1, d_2, \dots, d_n\}$  ( $d \in D =$   
 448  $\{D^{tr}, D^{te}\}$ ), different  $d$  construct the dataset of  
 449 LTR task, the target of the ranker is to predict the  
 450 dataset quality ranking for each dataset in  $d$  ac-  
 451 cording to the datasets’ features. The estimated  
 452 rankings  $\bar{\lambda} = \{\lambda_1, \lambda_2, \dots, \lambda_n\} \in [1, n]$  for set  $d$   
 453 can be defined as:

$$454 \bar{\lambda} = \bar{f}(\Phi_{(d)}; \bar{\Theta}) \quad (7)$$

455 where  $\Phi(\cdot)$  is the dataset feature extractor,  $\bar{f}$  is the  
 456 ranking model.  $\bar{\lambda} \in [1, n]$  is the estimated rankings  
 457 of the variance (e.g.  $\lambda_{var}$  or  $\lambda_{sva}$ ) for datasets in  
 458 set  $d$ .

## 459 5.2 Characterization of Datasets

460 In this section, we will introduce three aspects that  
 461 characterize datasets: Inherent Feature, Lexical  
 462 Feature, and Semantic Feature. Due to space limita-  
 463 tions, we move a more detailed feature introduction  
 464 to the Appendix C.

### 465 5.2.1 Inherent Feature

466 **Average length ( $\phi_{len}$ ):** The average sentence  
 467 length on a dataset, where the number of tokens on  
 468 a sentence is considered as the sentence length.

469 **Label number ( $\phi_{lab}$ ):** The number of labeled  
 470 classes in a dataset.

471 **Label balance ( $\phi_{bal}$ ):** The label balance metric  
 472 measures the variance between the ideal and the  
 473 true label distribution.

### 474 5.2.2 Lexical Feature

475 **Basic English Words Ratio ( $\phi_{basic}$ ):** The propor-  
 476 tion of words belonging to the 1000 basic English

<sup>10</sup> words in the whole dataset.

**Type-Token Ratio** ( $\phi_{\text{trr}}$ ): We measure the text lexical richness by the type-token ratio (Richards, 1987) based on the lexical richness tool <sup>11</sup>.

**Language Mixedness Ratio** ( $\phi_{\text{lmix}}$ ): To detect the ratio of other languages mixed in the text, we utilize the models proposed by Joulin et al. (2016b) for language identification from fastText (Joulin et al., 2016a) which can recognize 176 languages.

**Pointwise Mutual Information** ( $\phi_{\text{pmi}}$ ): PMI<sup>12</sup> is a measurement to calculate the correlation between variables.

### 5.2.3 Semantic Feature

**Perplexity** ( $\phi_{\text{ppl}}$ ): We calculate the perplexity <sup>13</sup> based on GPT2 (Radford et al., 2019) to evaluate the quality of the text.

**Grammar Errors Ratio** ( $\phi_{\text{gerr}}$ ): We adopt the detection tool <sup>14</sup> to recognize words with grammatical errors, and then calculate the ratio of grammatical errors.

**Flesch Reading Ease** <sup>15</sup> ( $\phi_{\text{fre}}$ ): To describe the readability of a text, we introduce the  $\phi_{\text{fre}}$  achieving by textstat <sup>16</sup>.

For feature  $\phi_{\text{len}}$ ,  $\phi_{\text{trr}}$ ,  $\phi_{\text{lmix}}$ ,  $\phi_{\text{gerr}}$ ,  $\phi_{\text{pmi}}$ ,  $\phi_{\text{fre}}$ , and  $\phi_{\text{rfre}}$ , we individually compute  $\phi(\cdot)$  on the training, test set, as well as their interaction. Take average length ( $\phi_{\text{len}}$ ) as an example, we compute the average length on training set  $\phi_{\text{tr, len}}$ , test set  $\phi_{\text{te, len}}$ , and their interaction  $((\phi_{\text{tr, len}} - \phi_{\text{te, len}}) / \phi_{\text{tr, len}})^2$ .

### 5.3 Parameterized Models

The dataset discrimination prediction (ranking) model takes a series of dataset features as the input and then predicts discrimination(rank) based on  $\hat{f}(\cdot)$  ( $\bar{f}(\cdot)$ ) defined in Eq. 6 (Eq. 7). We explore the effectiveness of four variations of regression methods and two ranking frameworks.

#### Regression Models

<sup>10</sup>[https://simple.wikipedia.org/wiki/Wikipedia:List\\_of\\_1000\\_basic\\_words](https://simple.wikipedia.org/wiki/Wikipedia:List_of_1000_basic_words)

<sup>11</sup><https://github.com/LSYS/lexicalrichness>

<sup>12</sup>[https://en.wikipedia.org/wiki/Pointwise\\_mutual\\_information](https://en.wikipedia.org/wiki/Pointwise_mutual_information)

<sup>13</sup><https://en.wikipedia.org/wiki/Perplexity>

<sup>14</sup>[https://github.com/jxmorris12/language\\_tool\\_python](https://github.com/jxmorris12/language_tool_python)

<sup>15</sup>[https://en.wikipedia.org/wiki/Flesch%E2%80%93Kincaid\\_readability\\_tests](https://en.wikipedia.org/wiki/Flesch%E2%80%93Kincaid_readability_tests)

<sup>16</sup><https://github.com/shivam5992/textstat>

• **LightGBM** (Ke et al., 2017) is a gradient boosting framework with faster training and better performance than XGBoost.

• **K-nearest Neighbor (KNN)** (Peterson, 2009) is a non-parametric model that makes the prediction by exploring the k neighbors.

• **Support Vector Machine (SVM)** (Suykens and Vandewalle, 1999) uses kernel trick to solve both linear and non-linear problems.

• **Decision Tree (DT)** (Quinlan, 1990) is a tree-based algorithm that gives an understandable interpretation of predictions.

#### Ranking Frameworks

• **LightGBM** with Gradient Boosting Decision Tree (Friedman, 2001) boosting strategy was selected as our ranking model.

• **XGBoost** (Chen and Guestrin, 2016) with gbtree(Hastie et al., 2009) boosting strategy was another ranking model.

### 5.4 Experiments

#### 5.4.1 Data Construction

To construct a collection with large amount of discriminative datasets, we randomly select three dataset features (e.g. average sentence length  $\phi_{\text{len}}$ ) to divide the original dataset into several non-overlapping sub-datasets. As a result, we collect 987 sub-datasets. Then, we train four text classification models (CNN, LSTM, LSTMAtt, BERT) on these sub-datasets. Next, we calculate the dataset features  $\phi$  (defined in Sec. 5.2) and dataset discrimination ability  $\lambda_{\text{sva}}$  and  $\lambda_{\text{var}}$  on these sub-datasets.

**Regression Task Settings**  $\phi$  and  $\lambda_{\text{sva}}$  ( $\lambda_{\text{var}}$ ) will be the input and target of the regression models, as defined by Eq. 6. For the experiment setting, we randomly select 287 ( $\phi$ ,  $\lambda_{\text{sva}}$  ( $\lambda_{\text{var}}$ )) pairs as the test set and the rest as the training set (700).

**Ranking Task Settings** We construct datasets for ranking task from the dataset used in regression task. Here, we explored the value of  $n$  (defined in §5.1.2) to be 5, 7 and 9 to randomly choose samples from  $D^{\text{tr}}$  (or  $D^{\text{te}}$ ) to construct the datasets for the ranking task, and kept 4, 200, 600, 1, 200 samples for training, development and testing set respectively.

#### 5.4.2 Evaluation Metric

**Regression Task** We use RMSE (Chai and Draxler, 2014) and Spearman rank correlation co-

efficient (Zar, 1972) to evaluate how well the regression model predicts the discriminative ability for datasets. The Spearman rank correlation coefficient is used to calculate the correlation between the output of a regression model and the ground truth.

**Ranking Task** NDCG and MAP are the evaluation metric of our ranking task. MAP is a binary preference metric, which focuses on whether the relevant document has a higher ranking than the irrelevant document. Here, we set a threshold value of  $\lambda_{\text{var}} = 3$  ( $\lambda_{\text{sva}} = 28$ ) for  $\lambda_{\text{var}}$  ( $\lambda_{\text{sva}}$ ) to distinguish the dataset discrimination ability from good (relevant) to bad (irrelevant).

Method	RMSE		Spearman			
	$\lambda_{\text{var}}$	$\lambda_{\text{sva}}$	$\lambda_{\text{var}}$		$\lambda_{\text{sva}}$	
			corr	p	corr	p
KNN	2.42	51.21	0.77	9.75E-40	0.87	1.62E-63
LightGBM	1.53	32.74	0.72	2.23E-33	0.87	7.01E-61
DT	1.73	43.33	0.64	9.25E-25	0.84	1.33E-53
SVM	2.83	62.44	0.68	1.14E-28	0.77	7.26E-40

Table 2: The performance of regressing dataset discrimination for the text classification. “corr” denotes the “correlation”.

Model	n	NDCG		MAP	
		$\lambda_{\text{var}}$	$\lambda_{\text{svar}}$	$\lambda_{\text{var}}$	$\lambda_{\text{svar}}$
LightGBM	9	98.20	98.85	97.50	98.27
	7	97.76	98.73	97.01	99.05
	5	96.73	97.08	96.56	98.15
XGBoost	9	96.66	97.13	92.91	93.62
	7	96.74	97.65	94.77	96.11
	5	95.93	97.10	95.49	98.25

Table 3: The performance of ranking dataset discrimination for the text classification task.  $n$  is the number of datasets in  $d$  defined in §5.1.2

### 5.4.3 Results and Analysis

Tab. 2 and Tab. 3 show the results of four regression models and two ranking models that characterize the dataset discrimination ability, respectively. We can observe that:

(1) **Both the regression models and the ranking models can well describe the discrimination ability of different datasets.** For these four regression models, the prediction is highly correlated with the ground truth (with a correlation value larger than 0.6), passing the significance testing ( $p < 0.05$ ). This suggests that the dataset discrimi-

nation can be successfully predicted. For these two ranking models, their performance on NDCG and MAP is greater than 95%, which indicates that the discriminative ability of the data set can be easily ranked.

(2)  **$\lambda_{\text{sva}}$  measure is better to characterize the discrimination ability of different datasets compared with  $\lambda_{\text{var}}$ .** For a regression model (e.g. KNN), the performance of  $\lambda_{\text{sva}}$  is better than  $\lambda_{\text{var}}$  significantly (higher correlation on  $\lambda_{\text{sva}}$ ), indicating that the dataset properties designed are more suitable for characterizing  $\lambda_{\text{sva}}$ . This conclusion can also be observed in the ranking models.

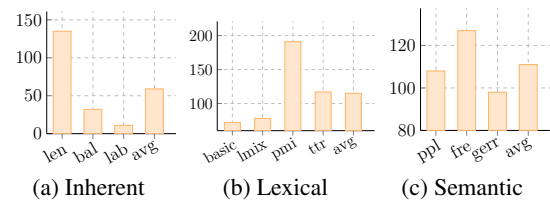


Figure 3: Feature importance for the text classification measured by LGBost with the target of  $\lambda_{\text{sva}}$ .

**Feature Importance Analysis** Fig. 3 illustrates the feature importance characterized by LightGBM. For a given feature, the number of times that is chosen as the splitting feature in the node of the decision trees is defined as its importance degree. We observe that: (1) The most influential features are  $\phi_{\text{pmi}}$ ,  $\phi_{\text{len}}$ , and  $\phi_{\text{fre}}$ , which come from the lexical, inherent, and semantic features, respectively. This indicated that the LightGBM can extract features from different aspects to make predictions. (2) In the perspective of feature groups, the semantic features are more influential than the inherent features and lexical features.

## 6 Implications and Future Directions

This paper has attempted to provide a methodology to characterize the discrimination ability (w.r.t top-scoring models) of the dataset, which allows to re-rank the value of the datasets. Some seldom used datasets may distinguish top-scoring models better than those frequently-used datasets. For the dataset with lower discrimination ability, we suggest proposing a more challenging dataset or make significance testing on these datasets. The idea can be applied to other NLP tasks. The suggestions and observations provided by this paper will inspire the future research.



## References

- Alexandra Birch, Miles Osborne, and Philipp Koehn. 2008. [Predicting success in machine translation](#). In *Proceedings of the 2008 Conference on Empirical Methods in Natural Language Processing*, pages 745–754, Honolulu, Hawaii. Association for Computational Linguistics.
- Samuel R. Bowman and George E. Dahl. 2021. [What will it take to fix benchmarking in natural language understanding?](#) *CoRR*, abs/2104.02145.
- Tianfeng Chai and Roland R Draxler. 2014. Root mean square error (rmse) or mean absolute error (mae)?—arguments against avoiding rmse in the literature. *Geoscientific model development*, 7(3):1247–1250.
- Tianqi Chen and Carlos Guestrin. 2016. [Xgboost](#). *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. [BERT: pre-training of deep bidirectional transformers for language understanding](#). *CoRR*, abs/1810.04805.
- Tobias Domhan, Jost Tobias Springenberg, and Frank Hutter. 2015. Speeding up automatic hyperparameter optimization of deep neural networks by extrapolation of learning curves. In *Twenty-Fourth International Joint Conference on Artificial Intelligence*.
- Jerome H Friedman. 2001. Greedy function approximation: a gradient boosting machine. *Annals of statistics*, pages 1189–1232.
- Jinlan Fu, Pengfei Liu, and Qi Zhang. 2020. [Rethinking generalization of neural models: A named entity recognition case study](#). In *The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, The Thirty-Second Innovative Applications of Artificial Intelligence Conference, IAAI 2020, The Tenth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2020, New York, NY, USA, February 7-12, 2020*, pages 7732–7739. AAAI Press.
- Sebastian Gehrmann, Tosin Adewumi, Karmanya Agarwal, Pawan Sasanka Ammanamanchi, Aremu Anuoluwapo, Antoine Bosselut, Khyathi Raghavi Chandu, Miruna Clinciu, Dipanjan Das, Kaustubh D Dhole, et al. 2021. The gem benchmark: Natural language generation, its evaluation and metrics. *arXiv preprint arXiv:2102.01672*.
- Ian Goodfellow, Yoshua Bengio, and Aaron Courville. 2016. *Deep Learning*. MIT Press. <http://www.deeplearningbook.org>.
- Harsha Gurulingappa, Abdul Mateen Rajput, Angus Roberts, Juliane Fluck, Martin Hofmann-Apitius, and Luca Toldo. 2012. [Development of a benchmark corpus to support the automatic extraction of drug-related adverse effects from medical case reports](#). *Journal of Biomedical Informatics*, 45(5):885–892.
- Text Mining and Natural Language Processing in Pharmacogenomics.
- Trevor Hastie, Robert Tibshirani, and Jerome Friedman. 2009. *Boosting and Additive Trees*, pages 337–387. Springer New York, New York, NY.
- Charles T Hemphill, John J Godfrey, and George R Doddington. 1990. The atis spoken language systems pilot corpus. In *Speech and Natural Language: Proceedings of a Workshop Held at Hidden Valley, Pennsylvania, June 24-27, 1990*.
- Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory. *Neural computation*, 9(8):1735–1780.
- Junjie Hu, Sebastian Ruder, Aditya Siddhant, Graham Neubig, Orhan Firat, and Melvin Johnson. 2020a. Xtreme: A massively multilingual multi-task benchmark for evaluating cross-lingual generalisation. In *International Conference on Machine Learning*, pages 4411–4421. PMLR.
- Junjie Hu, Sebastian Ruder, Aditya Siddhant, Graham Neubig, Orhan Firat, and Melvin Johnson. 2020b. [XTREME: A massively multilingual multi-task benchmark for evaluating cross-lingual generalization](#). *CoRR*, abs/2003.11080.
- Minqing Hu and Bing Liu. 2004. [Mining and summarizing customer reviews](#). In *Proceedings of the Tenth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD ’04*, page 168–177, New York, NY, USA. Association for Computing Machinery.
- Guangjin Jin and Xiao Chen. 2008. The fourth international chinese language processing bakeoff: Chinese word segmentation, named entity recognition and chinese pos tagging. In *Proceedings of the sixth SIGHAN workshop on Chinese language processing*.
- Armand Joulin, Edouard Grave, Piotr Bojanowski, Matthijs Douze, H erve J egou, and Tomas Mikolov. 2016a. Fasttext.zip: Compressing text classification models. *arXiv preprint arXiv:1612.03651*.
- Armand Joulin, Edouard Grave, Piotr Bojanowski, and Tomas Mikolov. 2016b. Bag of tricks for efficient text classification. *arXiv preprint arXiv:1607.01759*.
- Guolin Ke, Qi Meng, Thomas Finley, Taifeng Wang, Wei Chen, Weidong Ma, Qiwei Ye, and Tie-Yan Liu. 2017. [Lightgbm: A highly efficient gradient boosting decision tree](#). In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, pages 3146–3154.

734	Daniel Khashabi, Gabriel Stanovsky, Jonathan Bragg,	Jeffrey Pennington, Richard Socher, and Christopher D.	791
735	Nicholas Lourie, Jungo Kasai, Yejin Choi, Noah A	Manning. 2014. <a href="#">Glove: Global vectors for word</a>	792
736	Smith, and Daniel S Weld. 2021. Genie: A leader-	<a href="#">representation</a> . In <i>Proceedings of the 2014 Confer-</i>	793
737	board for human-in-the-loop evaluation of text gen-	<i>ence on Empirical Methods in Natural Language</i>	794
738	eration. <i>arXiv preprint arXiv:2101.06561</i> .	<i>Processing, EMNLP 2014, October 25-29, 2014,</i>	795
		<i>Doha, Qatar; A meeting of SIGDAT, a Special Inter-</i>	796
739	Philipp Koehn. 2004. <a href="#">Statistical significance tests</a>	<i>est Group of the ACL</i> , pages 1532–1543. ACL.	797
740	<a href="#">for machine translation evaluation</a> . In <i>Proceed-</i>		
741	<i>ings of the 2004 Conference on Empirical Meth-</i>	Leif E Peterson. 2009. K-nearest neighbor. <i>Scholarpe-</i>	798
742	<i>ods in Natural Language Processing</i> , pages 388–	<i>dia</i> , 4(2):1883.	799
743	395, Barcelona, Spain. Association for Computa-		
744	tional Linguistics.	John Ross Quinlan. 1990. Probabilistic decision trees.	800
		In <i>Machine Learning</i> , pages 140–152. Elsevier.	801
745	Prasanth Kolachina, Nicola Cancedda, Marc Dymet-		
746	man, and Sriram Venkatapathy. 2012. <a href="#">Prediction of</a>	Alec Radford, Jeffrey Wu, Rewon Child, David Luan,	802
747	<a href="#">learning curves in machine translation</a> . In <i>Proceed-</i>	Dario Amodei, and Ilya Sutskever. 2019. Language	803
748	<i>ings of the 50th Annual Meeting of the Association</i>	models are unsupervised multitask learners. <i>OpenAI</i>	804
749	<i>for Computational Linguistics (Volume 1: Long Pa-</i>	<i>blog</i> , 1(8):9.	805
750	<i>pers)</i> , pages 22–30, Jeju Island, Korea. Association		
751	for Computational Linguistics.	Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and	806
		Percy Liang. 2016. <a href="#">SQuAD: 100,000+ questions for</a>	807
752	Yann LeCun and Yoshua Bengio. 1995. Convolutional	<a href="#">machine comprehension of text</a> . In <i>Proceedings of</i>	808
753	networks for images, speech, and time series.	<i>the 2016 Conference on Empirical Methods in Natu-</i>	809
754	<i>The handbook of brain theory and neural networks</i> ,	<i>ral Language Processing</i> , pages 2383–2392, Austin,	810
755	3361(10).	Texas. Association for Computational Linguistics.	811
756	Mike Lewis, Yinhan Liu, Naman Goyal, Mar-	Brian Richards. 1987. <a href="#">Type/token ratios: what do</a>	812
757	jan Ghazvininejad, Abdelrahman Mohamed, Omer	<a href="#">they really tell us?</a> <i>Journal of Child Language</i> ,	813
758	Levy, Ves Stoyanov, and Luke Zettlemoyer. 2019.	14(2):201–209.	814
759	Bart: Denoising sequence-to-sequence pre-training		
760	for natural language generation, translation, and	Sebastian Ruder, Noah Constant, Jan Botha, Aditya	815
761	comprehension. <i>ArXiv</i> , abs/1910.13461.	Siddhant, Orhan Firat, Jinlan Fu, Pengfei Liu,	816
		Junjie Hu, Graham Neubig, and Melvin John-	817
762	Xin Li and Dan Roth. 2002. <a href="#">Learning question clas-</a>	son. 2021. <a href="#">XTREME-R: towards more challeng-</a>	818
763	<a href="#">sifiers</a> . In <i>COLING 2002: The 19th International</i>	<a href="#">ing and nuanced multilingual evaluation</a> . <i>CoRR</i> ,	819
764	<i>Conference on Computational Linguistics</i> .	abs/2104.07412.	820
765	Zhouhan Lin, Minwei Feng, Cícero Nogueira dos San-	Erik F Sang and Fien De Meulder. 2003. Intro-	821
766	tos, Mo Yu, Bing Xiang, Bowen Zhou, and Yoshua	duction to the conll-2003 shared task: Language-	822
767	Bengio. 2017. <a href="#">A structured self-attentive sentence</a>	independent named entity recognition. <i>arXiv</i>	823
768	<a href="#">embedding</a> . <i>CoRR</i> , abs/1703.03130.	<i>preprint cs/0306050</i> .	824
769	Tie-Yan Liu. 2011. Learning to rank for information	Claude E Shannon. 1948. A mathematical theory of	825
770	retrieval.	communication. <i>The Bell system technical journal</i> ,	826
		27(3):379–423.	827
771	Andrew L. Maas, Raymond E. Daly, Peter T. Pham,		
772	Dan Huang, Andrew Y. Ng, and Christopher Potts.	Richard Socher, Alex Perelygin, Jean Wu, Jason	828
773	2011. <a href="#">Learning word vectors for sentiment analy-</a>	Chuang, Christopher D. Manning, Andrew Ng, and	829
774	<a href="#">sis</a> . In <i>Proceedings of the 49th Annual Meeting of</i>	Christopher Potts. 2013. <a href="#">Recursive deep models</a>	830
775	<i>the Association for Computational Linguistics: Hu-</i>	<a href="#">for semantic compositionality over a sentiment tree-</a>	831
776	<i>man Language Technologies</i> , pages 142–150, Port-	<a href="#">bank</a> . In <i>Proceedings of the 2013 Conference on</i>	832
777	land, Oregon, USA. Association for Computational	<i>Empirical Methods in Natural Language Processing</i> ,	833
778	Linguistics.	pages 1631–1642, Seattle, Washington, USA. Asso-	834
		ciation for Computational Linguistics.	835
779	Tomás Mikolov, Ilya Sutskever, Kai Chen, Gregory S.	Johan AK Suykens and Joos Vandewalle. 1999. Least	836
780	Corrado, and Jeffrey Dean. 2013. <a href="#">Distributed rep-</a>	squares support vector machine classifiers. <i>Neural</i>	837
781	<a href="#">resentations of words and phrases and their com-</a>	<i>processing letters</i> , 9(3):293–300.	838
782	<a href="#">positionality</a> . In <i>Advances in Neural Information</i>		
783	<i>Processing Systems 26: 27th Annual Conference on</i>	Huihsin Tseng, Pichuan Chang, Galen Andrew, Daniel	839
784	<i>Neural Information Processing Systems 2013. Pro-</i>	Jurafsky, and Christopher Manning. 2005. A condi-	840
785	<i>ceedings of a meeting held December 5-8, 2013,</i>	tional random field word segmenter for sighthan bake-	841
786	<i>Lake Tahoe, Nevada, United States</i> , pages 3111–	off 2005. In <i>Proceedings of the fourth SIGHAN</i>	842
787	3119.	<i>workshop on Chinese language Processing</i> , volume	843
		171.	844
788	Bo Pang and Lillian Lee. 2005. <a href="#">Seeing stars: Exploit-</a>		
789	<a href="#">ing class relationships for sentiment categorization</a>		
790	<a href="#">with respect to rating scales</a> . <i>CoRR</i> , abs/cs/0506075.		

- 845 Marco Turchi, Tijl De Bie, and Nello Cristianini. 2008.  
846 Learning performance of a machine translation sys-  
847 tem: a statistical and computational analysis. In *Pro-*  
848 *ceedings of the Third Workshop on Statistical Ma-*  
849 *chine Translation*, pages 35–43.
- 850 Alex Wang, Yada Pruksachatkun, Nikita Nangia,  
851 Amanpreet Singh, Julian Michael, Felix Hill, Omer  
852 Levy, and Samuel R Bowman. 2019a. Super-  
853 glue: A stickier benchmark for general-purpose  
854 language understanding systems. *arXiv preprint*  
855 *arXiv:1905.00537*.
- 856 Alex Wang, Amanpreet Singh, Julian Michael, Fe-  
857 lix Hill, Omer Levy, and Samuel Bowman. 2018.  
858 [GLUE: A multi-task benchmark and analysis plat-](#)  
859 [form for natural language understanding](#). In *Pro-*  
860 *ceedings of the 2018 EMNLP Workshop Black-*  
861 *boxNLP: Analyzing and Interpreting Neural Net-*  
862 *works for NLP*, pages 353–355, Brussels, Belgium.  
863 Association for Computational Linguistics.
- 864 Alex Wang, Amanpreet Singh, Julian Michael, Felix  
865 Hill, Omer Levy, and Samuel R. Bowman. 2019b.  
866 [GLUE: A multi-task benchmark and analysis plat-](#)  
867 [form for natural language understanding](#). In *7th*  
868 *International Conference on Learning Representa-*  
869 *tions, ICLR 2019, New Orleans, LA, USA, May 6-9,*  
870 *2019*. OpenReview.net.
- 871 Adina Williams, Nikita Nangia, and Samuel R Bow-  
872 man. 2017. A broad-coverage challenge corpus for  
873 sentence understanding through inference. *arXiv*  
874 *preprint arXiv:1704.05426*.
- 875 Mengzhou Xia, Antonios Anastasopoulos, Ruochen  
876 Xu, Yiming Yang, and Graham Neubig. 2020. [Pre-](#)  
877 [dicting performance for natural language process-](#)  
878 [ing tasks](#). In *Proceedings of the 58th Annual Meet-*  
879 *ing of the Association for Computational Linguistics*,  
880 pages 8625–8646, Online. Association for Computa-  
881 tional Linguistics.
- 882 Zihuiwen Ye, Pengfei Liu, Jinlan Fu, and Graham Neu-  
883 big. 2021. [Towards more fine-grained and reliable](#)  
884 [NLP performance prediction](#). In *Proceedings of the*  
885 *16th Conference of the European Chapter of the*  
886 *Association for Computational Linguistics: Main*  
887 *Volume*, pages 3703–3714, Online. Association for  
888 Computational Linguistics.
- 889 Jerrold H Zar. 1972. Significance testing of the spear-  
890 man rank correlation coefficient. *Journal of the*  
891 *American Statistical Association*, 67(339):578–580.
- 892 Xiang Zhang, Junbo Jake Zhao, and Yann LeCun. 2015.  
893 [Character-level convolutional networks for text clas-](#)  
894 [sification](#). *CoRR*, abs/1509.01626.
- 895 Ming Zhong, Pengfei Liu, Yiran Chen, Danqing Wang,  
896 Xipeng Qiu, and Xuanjing Huang. 2020. [Extrac-](#)  
897 [tive summarization as text matching](#). In *Proceedings*  
898 *of the 58th Annual Meeting of the Association for*  
899 *Computational Linguistics*, pages 6197–6208, On-  
900 line. Association for Computational Linguistics.

## A Statistics of Datasets

Tab. 4 shows the statistical information of the nine datasets of text classification task used in our work. For those datasets without explicit the development set, we randomly selected 12.5% samples from the training set as the development set.

Dataset	Train	Test	Development
IMDB	25,000	25,000	-
Yelp	560,000	38,000	-
QC	5,452	500	-
DPedia	560,000	70,000	-
CR	3,594	400	-
ATIS	4,978	893	-
SST1	8,544	2,210	1,101
MR	9,596	1,066	-
ADE	23,516	-	-

Table 4: Statistics of datasets.

## B Parameter Settings for Text Classification Model

In this section, we will introduce the parameter settings of the neural network-based models explored in Section 3.2. The optimizer is AdamW for the four models. The settings of other parameters are shown in Tab. 5.

Parameter	BERT	CNN	LSTM	LSTMAtt
learning rate	2*e-5	1*e-4	1*e-3	1*e-3
batch size	4	4	32	32
word emb	-	Word2vec	GloVe	GloVe
word emb size	-	300	300	300
hidden size	768	120	256	256
max sent len	512	-	-	-
filter size	-	1,3,5	-	-

Table 5: the parameters of four models.

## C Characterization of Datasets

### C.1 Inherent Feature

**Label balance** ( $\phi_{\text{bal}}$ ): The label balance metric measures the variance between the ideal and the true label distribution:  $\phi_{\text{bal}} = (c_t - c_s)/c_s$ , where the  $c_t$  and  $c_s$  are the true and ideal label information entropy (Shannon, 1948), respectively.

### C.2 Lexical Feature

**Type-Token Ratio** ( $\phi_{\text{ttr}}$ ): TTR (Richards, 1987) is a way to measure the documents lexical richness:  $\phi_{\text{ttr}} = n_{\text{type}}/n_{\text{token}}$ , where the  $n_{\text{type}}$  is the number

of unique words, and  $n_{\text{token}}$  is the number of tokens. We use lexical richness<sup>17</sup> to calculate the TTR for each sentence and then average them.

**Language Mixedness Ratio** ( $\phi_{\text{lmix}}$ ): The proportion of sentence that contains other languages in the whole dataset. To detect the mixed other languages, we utilize the models proposed by Joulin et al. (2016b) for language identification from fast-Text (Joulin et al., 2016a) which can recognize 176 languages.

**Pointwise Mutual Information** ( $\phi_{\text{pmi}}$ ): is a measurement to calculate the correlation between variables. Specifically, for a word in one class  $\phi_{\text{pmi}(c,w)} = \log(\frac{p(c,w)}{p(c)p(w)})$ , where  $p(c)$  is the proportion of the tokens belonging to label  $c$ ,  $p(w)$  is the proportion of the word  $w$ , and  $p(c,w)$  is the proportion of the word  $w$  which belongs to class  $c$ . For every class, all the  $\phi_{\text{pmi}(c,w)}$ , larger than zero, are added to get the sum, which serve as the dataset’s pmi. Finally,  $\phi_{\text{pmi}}$  is calculated by dividing the sum by the numbers of pairs(c,w) of the train dataset. We pick up the top-ten words sorted by  $\phi_{\text{pmi}(c,w)}$  in all classes, then the ration related to the class-related word( $\phi_{\text{r_pmi}}$ ) is calculated by dividing the number of samples who contain the top-ten words by the total samples in the train set.

### C.3 Semantic Feature

**Grammar errors ratio** ( $\phi_{\text{gerr}}$ ): The proportion of words with grammatical errors in the whole dataset. We adopt the detection tool<sup>18</sup> to recognize words with grammatical errors. We first compute the grammar errors ratio for each sentence:  $n/m$ , where the n and m denote the number of words with grammatical errors and the number of the token for a sentence, averaging them.

**Flesch Reading Ease** ( $\phi_{\text{fre}}$ ): Flesch Reading Ease<sup>19</sup> calculated by textstat<sup>20</sup> is a way to describe the simplicity of a reader who can read a text. First, we calculate the  $\phi_{\text{fre}}$  for each sample, and then average them as the dataset’s feature. Then we pick out the samples whose score below 60, then the ration related to the low score samples( $\phi_{\text{r_fre}}$ ) is calculated by dividing the number of the picked samples by the total samples in the train set.

<sup>17</sup><https://github.com/LSYS/lexicalrichness>

<sup>18</sup>[https://github.com/jxmorris12/language\\_tool\\_python](https://github.com/jxmorris12/language_tool_python)

<sup>19</sup>[https://en.wikipedia.org/wiki/Flesch%E2%80%93Kincaid\\_readability\\_tests](https://en.wikipedia.org/wiki/Flesch%E2%80%93Kincaid_readability_tests)

<sup>20</sup><https://github.com/shivam5992/textstat>