

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](#)¹ 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 large pre-trained models consistently improving state-of-the-art performance on NLP tasks (Devlin et al., 2018; Lewis et al., 2019), the performances of many of them have reached a plateau (Zhong et al., 2020; Fu et al., 2020). In other words, it is

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

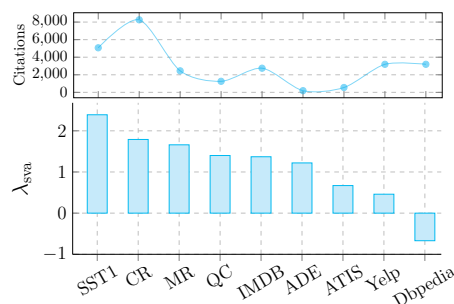


Figure 1: Illustrate different datasets’ distinguishing ability w.r.t top-scoring systems characterized by our measure $\log(\lambda_{sva})$ on text classification and their corresponding citations.

challenging to discriminate a better model using existing datasets (Wang 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 distin-

066 guishing top-scoring systems, suggesting that some
067 of the relatively neglected datasets are actually valu-
068 able in distinguishing models. This phenomenon
069 shows the significance of quantifying the discrim-
070 inative ability of datasets: it can not only help us
071 to **eliminate** those with lower discrimination from
072 *commonly-used datasets* (e.g., DBpedia), but also
073 help us to **recognize** the missing pearl in *seldom*
074 *used datasets* (e.g., ADE and ATIS (Hemphill et al.,
075 1990)).

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

095 We brief **takeaways** in this work based on our
096 observations:

097 (1) Not all datasets in *benchmark* are necessary
098 in terms of model selection²: empirical results
099 show that following datasets struggle at discrim-
100 inating current top-scoring systems: STS-B and
101 SST-2 from GLUE (Wang et al., 2019b); BUCC
102 and PAWX-X from XTREME, which is consis-
103 tent with the concurrent work (Ruder et al., 2021)
104 (§4.3.2).

105 (2) In regard to single-task benchmark datasets,
106 for Chinese Word Segmentation task, there are
107 multiple datasets (MSR, CityU, CTB) (Tseng
108 et al., 2005; Jin and Chen, 2008) that exhibit much
109 worse discriminative ability, suggesting that: fu-
110 ture works on this task are encouraged to either
111 (i) adopt other datasets to evaluate their systems
112 or (ii) at least make significant test³ if using these

²Caveat: Annotated datasets are always valuable, because the supervision signals provided there can not only help us directly train a system for specific use case, but also provide good supervised transfer for related tasks (Sanh et al., 2021).

³We randomly select 10 recently published papers (from

113 datasets. Similar observations happen in the dataset
114 CoNLL-2003 (Sang and De Meulder, 2003) from
115 Named Entity Recognition task and MultiNLI
116 (Williams et al., 2017) from natural language infer-
117 ence task (§4.3.2).

118 (3) Some seldom used datasets such as ADE from
119 text classification are actually better at distinguish-
120 ing top-performing systems, which highlights an
121 interesting and necessary future direction: *how to*
122 *identify infrequently-used but valuable (better dis-*
123 *crimination) datasets for NLP tasks, especially in*
124 *the age of dataset’s proliferation?*⁴ (§4.2)

125 (4) Quantifying a dataset’s discrimination (w.r.t
126 top-scoring systems) by calculating the statistical
127 measures (defined in §4.1.2) from leaderboard’s
128 results is a straightforward and effective way. But
129 for those datasets without rich leaderboard results,⁵
130 predicting the discrimination based on datasets’
131 characteristics would be an promising direction
132 (§4.3.1).

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

134 (1) We try to quantify the discrimination abil-
135 ity for datasets by designing two variance-based
136 measures. (2) We systematically investigate 4 text
137 classification models on 9 datasets, providing the
138 newest baseline performance for those seldom used
139 datasets. We released the code and all the uni-
140 formly formatted datasets at <https://github.com/annonnlp-demo/acl-v2> (3) We study
141 several popular NLP benchmarks, including GLUE,
142 XTREME, NLI, and so on. Some valuable sugges-
143 tions and observations will make research easier.
144

145 2 Related Work

146 **Benchmarks for NLP** In order to conveniently
147 keep themselves updated with the research
148 progress, researchers recently are actively build-
149 ing evaluation benchmarks for diverse tasks so
150 that they could make a comprehensive compari-
151 son of systems, and use a leaderboard to record the
152 evolving process of the systems of different NLP
153 tasks, such as SQuAD (Rajpurkar et al., 2016),
154 GLUE (Wang et al., 2018), XTREME (Hu et al.,
155 2020a), GEM (Gehrmann et al., 2021) and GE-
156 NIE (Khashabi et al., 2021). Despite their utility,
157 more recently, Bowman and Dahl (2021) highlight

ACL/EMNLP) that utilized these datasets and found only 2 of them perform significant test.

⁴<https://paperswithcode.com/datasets>

⁵The measure can keeps updated as the top-scoring systems of the leaderboard evolves, which can broaden its practical applicability

that unreliable and biased systems score so highly on standard benchmarks that there is little room for researchers who develop better systems to demonstrate their improvements. In this paper, we make a pilot study on meta-evaluating benchmark evaluation datasets and quantitatively characterize their discrimination in different top-scoring systems.

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.
- **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 collected from DBpedia.
- **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.⁶ The brief introduction of the four models is as follows.

- **LSTM** (Hochreiter and Schmidhuber, 1997) is a widely used sentence encoder. 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) was utilized to fine-tuning on our text classification datasets.
- **CNN** is a CNN-based text classification model (Kim, 2014) was expolred in our work.

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 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

We design several measures to judge dataset’s distinguishing ability based on the performances that top-performing systems have achieved on it.⁷ 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 λ_{var} suggests that the discrimination of the given dataset is more significant. λ_{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 $\lambda_{\text{var}} = 2.65$.

⁶We mainly focus on neural network-based models, since most top-scoring systems in the leaderboard are based on deep learning.

⁷A 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 λ_{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 ($t = 1, \dots, T$)⁸. 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

⁸For example, given a test set with 1000 samples, we sample 80% subset from it and repeat this process T times.

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., λ_{sva}), which are easier to obtain, by calculating its correlation with significant test-based measure λ_{hit} , which is costly to get. Since the implementation of λ_{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. λ_{var} and λ_{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 λ_{hit} .

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

where the q can be λ_{var} or λ_{sva} .

Result (1) λ_{var} and λ_{sva} are strong correlative ($S_\lambda > 0.6$) with λ_{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 λ_{sva} , dataset discrimination characterized by λ_{var} is more acceptable for λ_{hit} . The reason can be attributed to that the designing of the measure λ_{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 λ_{var} and λ_{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 and two multitask benchmarks. For the single-task

Method	BERT	LSTMAAttr	LSTM	CNN	λ_{hit}	λ_{var}	λ_{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 λ_{hit} , λ_{var} , and λ_{sva}) on 9 text classification datasets. The two correlation coefficients pass the significance test ($p < 0.05$). λ_{var} and λ_{sva} measures are designed based on performance variance.

benchmarks, we collect the top-performing models in a specific period for each dataset, provided by Paperswithcode⁹. For the multitask benchmarks, here, the GLUE¹⁰ and XTREME¹¹ 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 λ_{var} and λ_{sva} . We detail several main observations:

- λ_{var} and λ_{sva} have consistent evaluation results

⁹<https://paperswithcode.com/>

¹⁰<https://gluebenchmark.com/>

¹¹<https://sites.research.google/xtreme>

for both single-task (CWS, NER, NLI) and multitask (GLUE, XTREME) benchmarks.

- For the XTREME benchmark, BUCC and PAWSX have lowest λ_{var} and λ_{sva} , which suggest that they are hardly to discriminate the top-performing systems. Moreover, these two datasets 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 have the excellent ability to distinguish different top-scoring models (with higher λ_{var} and λ_{sva}), while the SST-2 and STS-B perform worse.
- For CWS benchmarks, there is a larger gap between the value of λ_{var} and λ_{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 λ_{var} , λ_{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 alleviate this problem, in this section, we focus on text classification task and investigate the possibility of

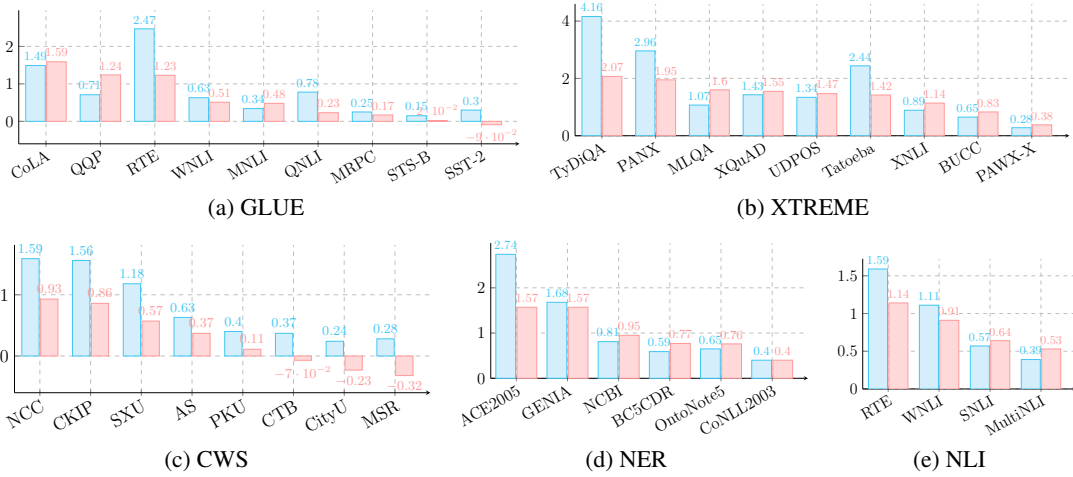


Figure 2: The dataset discrimination characterized by $\log(\lambda_{var})$ (the logarithm for better visualization) (blue) and $\log(\lambda_{sva})$ (pink) on five popular NLP benchmarks.

estimating datasets’ discrimination solely based on their characteristics without actual training systems on them.

5.1 Task Formulation

5.1.1 Regression-based Task Formulation

We formulate it as a performance prediction problem (Birch et al., 2008; Xia et al., 2020; Ye et al., 2021). Formally, we refer to \mathcal{M} , D^{tr} , D^{te} , \mathcal{S} as the machine learning system, training data, test data and training strategy respectively. The goal of performance prediction is to estimate actual performance y without actual training by using features of \mathcal{M} , D^{tr} , D^{te} , and \mathcal{S} .

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

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

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

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

5.1.2 Ranking-based Task Formulation

Instead of only regressing one dataset’s quality, we also care about the quality ranking of different datasets w.r.t discriminating systems in a task. Therefore, we also formulate it as a listwise LTR (learning to rank) task where a model takes individual lists as instances, to predict the rank of element among the list (Liu, 2011). Given a set of n datasets $d = \{d_1, d_2, \dots, d_n\}$ ($d \in D =$

$\{D^{tr}, D^{te}\}$), different d construct the dataset of LTR task, the target of the ranker is to predict the dataset quality ranking for each dataset in d according to the datasets’ features. The estimated rankings $\bar{\lambda} = \{\lambda_1, \lambda_2, \dots, \lambda_n\} \in [1, n]$ for set d can be defined as:

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

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

5.2 Characterization of Datasets

In this section, we will introduce three aspects that characterize datasets: Inherent Feature, Lexical Feature, and Semantic Feature. Due to space limitations, we move a more detailed feature introduction to the Appendix C.

5.2.1 Inherent Feature

Average length (ϕ_{len}): The average sentence length on a dataset, where the number of tokens on a sentence is considered as the sentence length. **Label number (ϕ_{lab}):** The number of labeled classes in a dataset. **Label balance (ϕ_{bal}):** The label balance metric measures the variance between the ideal and the true label distribution.

5.2.2 Lexical Feature

Basic English Words Ratio (ϕ_{basic}): The proportion of words belonging to the 1000 basic English ¹² words in the whole dataset. **Type-Token Ratio (ϕ_{ttr}):** We measure the text lexical richness by

¹²https://simple.wikipedia.org/wiki/Wikipedia:List_of_1000_basic_words

the type-token ratio (Richards, 1987) based on the lexical richness tool ¹³. **Language Mixedness Ratio** (ϕ_{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** (ϕ_{pmi}): PMI¹⁴ is a measurement to calculate the correlation between variables.

5.2.3 Semantic Feature

Perplexity (ϕ_{pp}): We calculate the perplexity ¹⁵ based on GPT2 (Radford et al., 2019) to evaluate the quality of the text. **Grammar Errors Ratio** (ϕ_{gerr}): We adopt the detection tool ¹⁶ to recognize words with grammatical errors, and then calculate the ratio of grammatical errors. **Flesch Reading Ease** ¹⁷ (ϕ_{fre}): To describe the readability of a text, we introduce the ϕ_{fre} achieving by textstat ¹⁸.

For feature $\phi_{\text{len}}, \phi_{\text{tr}}, \phi_{\text{lmix}}, \phi_{\text{gerr}}, \phi_{\text{pmi}}, \phi_{\text{fre}}$, and ϕ_{rfe} , we individually compute $\phi(\cdot)$ on the training, test set, as well as their interaction. Take average length (ϕ_{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: 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

¹³<https://github.com/LSYS/lexicalrichness>

¹⁴https://en.wikipedia.org/wiki/Pointwise_mutual_information

¹⁵<https://en.wikipedia.org/wiki/Perplexity>

¹⁶https://github.com/jxmorris12/language_tool_python

¹⁷https://en.wikipedia.org/wiki/Flesch%E2%80%93Kincaid_readability_tests

¹⁸<https://github.com/shivam5992/textstat>

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 gbtrees (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 ϕ_{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 ϕ (defined in Sec. 5.2) and dataset discrimination ability λ_{sva} and λ_{var} on these sub-datasets.

Regression Task Settings ϕ and λ_{sva} (λ_{var}) will be the input and target of the regression models, as defined by Eq. 6. For the experiment setting, we randomly select 287 ($\phi, \lambda_{\text{sva}}$ (λ_{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^{tr} (or D^{te}) to construct the datasets for the ranking task, and kept 4200, 600, 1200 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 coefficient (Zar, 1972) to evaluate how well the regression model predicts the discriminative ability for datasets. The Spearman rank correlation coefficient is used for the correlation between the output of a regression model and the ground truth.

Ranking Task NDCG (Järvelin and Kekäläinen, 2000) and MAP (Yue et al., 2007) are the evaluation metric of our ranking task. For NDCG, it considers the rank of a set of discriminative abilities. In our setting, every dataset has its own real discriminative ability. Here, We transfer the predicted discriminative ability to the rank of the dataset in the NDCG metric, so we can use NDCG to evaluate the model’s predicted effect. For MAP, it likes

how NDCG works, but it considers a set of binary values. Here, we set a threshold value of $\lambda_{\text{var}} = 3$ ($\lambda_{\text{sva}} = 28$) for λ_{var} (λ_{sva}) to distinguish the dataset discrimination ability from good (relevant) to bad (irrelevant).

Method	RMSE		Spearman			
	λ_{var}	λ_{sva}	λ_{var}		λ_{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	
		λ_{var}	λ_{sva}	λ_{var}	λ_{sva}
	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: 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 discrimination 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.

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

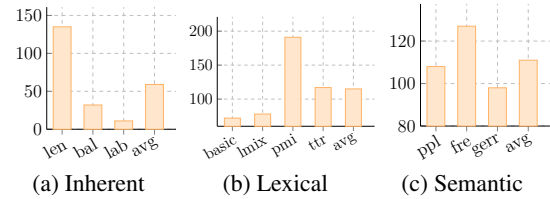


Figure 3: Feature importance for the text classification measured by LGBBoost with the target of λ_{sva} .

observe that: (1) The most influential features are ϕ_{pmi} , ϕ_{len} , and ϕ_{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 Discussion & Implications

Discussion Given a leaderboard of a dataset, metrics explored in this paper can be easily used to calculate its discrimination, while some limitations still exist. We make some discussion below to encourage more explorations on new measures: (a) **Interpretability**: current metrics can only identify which datasets are of lower indiscriminability while don’t present more explanation why it is the case. (b) **Functionality**: a dataset with lower discrimination doesn’t mean it’s useless since the supervision signals provided there can not only help us directly train a system for the specific use case but also provide good supervised transfer for related tasks. Metrics designed in this work focus on the role of discriminating models.

Calls Based on observations obtained from this paper, we make the following calls for future research: (1) Datasets’ discrimination ability w.r.t top-scoring systems could be included in the dataset schema (such as dataset statement (Bender and Friedman, 2018)), which would allow researchers to gain a saturated understanding of the dataset. (2) Leaderboard constructors could also report the discriminative ability of the datasets they aim to include. (3) Seldom used datasets are also valuable for model selection, and a more fair dataset searching system should be investigated, for example, relevance- and scientifically meaningful first, instead of other biases, like popularity.

626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680

References

Emily M. Bender and Batya Friedman. 2018. [Data statements for natural language processing: Toward mitigating system bias and enabling better science](#). *Transactions of the Association for Computational Linguistics*, 6:587–604.

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 Aggarwal, 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. 687 688

Trevor Hastie, Robert Tibshirani, and Jerome Friedman. 2009. *Boosting and Additive Trees*, pages 337–387. Springer New York, New York, NY. 689 690 691

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*. 692 693 694 695 696

Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory. *Neural computation*, 9(8):1735–1780. 697 698 699

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. 700 701 702 703 704 705

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. 706 707 708 709 710

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. 711 712 713 714 715 716

Kalervo Järvelin and Jaana Kekäläinen. 2000. [IR evaluation methods for retrieving highly relevant documents](#). In *SIGIR 2000: Proceedings of the 23rd Annual International ACM SIGIR Conference on Research and Development in Information Retrieval, July 24-28, 2000, Athens, Greece*, pages 41–48. ACM. 717 718 719 720 721 722 723

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*. 724 725 726 727 728 729

Armand Joulin, Edouard Grave, Piotr Bojanowski, Matthijs Douze, Herve Jégou, and Tomas Mikolov. 2016a. [Fasttext.zip: Compressing text classification models](#). *arXiv preprint arXiv:1612.03651*. 730 731 732 733

Armand Joulin, Edouard Grave, Piotr Bojanowski, and Tomas Mikolov. 2016b. [Bag of tricks for efficient text classification](#). *arXiv preprint arXiv:1607.01759*. 734 735 736 737

738	Guolin Ke, Qi Meng, Thomas Finley, Taifeng Wang,	<i>Neural Information Processing Systems 2013. Pro-</i>	794
739	Wei Chen, Weidong Ma, Qiwei Ye, and Tie-Yan Liu.	<i>ceedings of a meeting held December 5-8, 2013,</i>	795
740	2017. Lightgbm: A highly efficient gradient boost-	<i>Lake Tahoe, Nevada, United States</i> , pages 3111–	796
741	ing decision tree . In <i>Advances in Neural Informa-</i>	3119.	797
742	<i>tion Processing Systems 30: Annual Conference on</i>		
743	<i>Neural Information Processing Systems 2017, De-</i>	Bo Pang and Lillian Lee. 2005. Seeing stars: Exploit-	798
744	<i>cember 4-9, 2017, Long Beach, CA, USA</i> , pages	ing class relationships for sentiment categorization	799
745	3146–3154.	with respect to rating scales . <i>CoRR</i> , abs/cs/0506075.	800
746	Daniel Khashabi, Gabriel Stanovsky, Jonathan Bragg,	Jeffrey Pennington, Richard Socher, and Christopher D.	801
747	Nicholas Lourie, Jungo Kasai, Yejin Choi, Noah A	Manning. 2014. Glove: Global vectors for word	802
748	Smith, and Daniel S Weld. 2021. Genie: A leader-	representation . In <i>Proceedings of the 2014 Confer-</i>	803
749	board for human-in-the-loop evaluation of text gen-	<i>ence on Empirical Methods in Natural Language</i>	804
750	eration. <i>arXiv preprint arXiv:2101.06561</i> .	<i>Processing, EMNLP 2014, October 25-29, 2014,</i>	805
751	Yoon Kim. 2014. Convolutional neural networks for	<i>Doha, Qatar, A meeting of SIGDAT, a Special Inter-</i>	806
752	sentence classification . <i>CoRR</i> , abs/1408.5882.	<i>est Group of the ACL</i> , pages 1532–1543. ACL.	807
753	Philipp Koehn. 2004. Statistical significance tests	Leif E Peterson. 2009. K-nearest neighbor. <i>Scholarpe-</i>	808
754	for machine translation evaluation . In <i>Proceeed-</i>	<i>dia</i> , 4(2):1883.	809
755	<i>ings of the 2004 Conference on Empirical Meth-</i>	John Ross Quinlan. 1990. Probabilistic decision trees.	810
756	<i>ods in Natural Language Processing</i> , pages 388–	In <i>Machine Learning</i> , pages 140–152. Elsevier.	811
757	395, Barcelona, Spain. Association for Computa-		
758	tional Linguistics.	Alec Radford, Jeffrey Wu, Rewon Child, David Luan,	812
759	Prasanth Kolachina, Nicola Cancedda, Marc Dymet-	Dario Amodei, and Ilya Sutskever. 2019. Language	813
760	man, and Sriram Venkatapathy. 2012. Prediction of	models are unsupervised multitask learners. <i>OpenAI</i>	814
761	learning curves in machine translation . In <i>Proceeed-</i>	<i>blog</i> , 1(8):9.	815
762	<i>ings of the 50th Annual Meeting of the Association</i>	Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and	816
763	<i>for Computational Linguistics (Volume 1: Long Pa-</i>	Percy Liang. 2016. SQuAD: 100,000+ questions for	817
764	<i>pers)</i> , pages 22–30, Jeju Island, Korea. Association	machine comprehension of text . In <i>Proceedings of</i>	818
765	for Computational Linguistics.	<i>the 2016 Conference on Empirical Methods in Natu-</i>	819
766	Mike Lewis, Yinhan Liu, Naman Goyal, Mar-	<i>ral Language Processing</i> , pages 2383–2392, Austin,	820
767	jan Ghazvininejad, Abdelrahman Mohamed, Omer	Texas. Association for Computational Linguistics.	821
768	Levy, Ves Stoyanov, and Luke Zettlemoyer. 2019.	Brian Richards. 1987. Type/token ratios: what do	822
769	Bart: Denoising sequence-to-sequence pre-training	they really tell us? <i>Journal of Child Language</i> ,	823
770	for natural language generation, translation, and	14(2):201–209.	824
771	comprehension. <i>ArXiv</i> , abs/1910.13461.		
772	Xin Li and Dan Roth. 2002. Learning question clas-	Sebastian Ruder, Noah Constant, Jan Botha, Aditya	825
773	sifiers . In <i>COLING 2002: The 19th International</i>	Siddhant, Orhan Firat, Jinlan Fu, Pengfei Liu,	826
774	<i>Conference on Computational Linguistics</i> .	Junjie Hu, Graham Neubig, and Melvin John-	827
775	Zhouhan Lin, Minwei Feng, Cícero Nogueira dos San-	son. 2021. XTREME-R: towards more challeng-	828
776	tos, Mo Yu, Bing Xiang, Bowen Zhou, and Yoshua	ing and nuanced multilingual evaluation . <i>CoRR</i> ,	829
777	Bengio. 2017. A structured self-attentive sentence	abs/2104.07412.	830
778	embedding . <i>CoRR</i> , abs/1703.03130.	Erik F Sang and Fien De Meulder. 2003. Intro-	831
779	Tie-Yan Liu. 2011. Learning to rank for information	duction to the conll-2003 shared task: Language-	832
780	retrieval.	independent named entity recognition. <i>arXiv</i>	833
781	Andrew L. Maas, Raymond E. Daly, Peter T. Pham,	<i>preprint cs/0306050</i> .	834
782	Dan Huang, Andrew Y. Ng, and Christopher Potts.	Victor Sanh, Albert Webson, Colin Raffel, Stephen H	835
783	2011. Learning word vectors for sentiment analy-	Bach, Lintang Sutawika, Zaid Alyafeai, Antoine	836
784	sis . In <i>Proceedings of the 49th Annual Meeting of</i>	Chaffin, Arnaud Stiegler, Teven Le Scao, Arun	837
785	<i>the Association for Computational Linguistics: Hu-</i>	Raja, et al. 2021. Multitask prompted training en-	838
786	<i>man Language Technologies</i> , pages 142–150, Port-	ables zero-shot task generalization. <i>arXiv preprint</i>	839
787	land, Oregon, USA. Association for Computational	<i>arXiv:2110.08207</i> .	840
788	Linguistics.	Claude E Shannon. 1948. A mathematical theory of	841
789	Tomás Mikolov, Ilya Sutskever, Kai Chen, Gregory S.	communication. <i>The Bell system technical journal</i> ,	842
790	Corrado, and Jeffrey Dean. 2013. Distributed rep-	27(3):379–423.	843
791	resentations of words and phrases and their com-	Richard Socher, Alex Perelygin, Jean Wu, Jason	844
792	positionality . In <i>Advances in Neural Information</i>	Chuang, Christopher D. Manning, Andrew Ng, and	845
793	<i>Processing Systems 26: 27th Annual Conference on</i>	Christopher Potts. 2013. Recursive deep models	846

847	for semantic compositionality over a sentiment tree-bank . In <i>Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing</i> , pages 1631–1642, Seattle, Washington, USA. Association for Computational Linguistics.	903
848		904
849		
850		
851		
852	Johan AK Suykens and Joos Vandewalle. 1999. Least squares support vector machine classifiers. <i>Neural processing letters</i> , 9(3):293–300.	
853		
854		
855	Huihsin Tseng, Pichuan Chang, Galen Andrew, Daniel Jurafsky, and Christopher Manning. 2005. A conditional random field word segmenter for sighan bake-off 2005. In <i>Proceedings of the fourth SIGHAN workshop on Chinese language Processing</i> , volume 171.	
856		
857		
858		
859		
860		
861	Marco Turchi, Tjil De Bie, and Nello Cristianini. 2008. Learning performance of a machine translation system: a statistical and computational analysis. In <i>Proceedings of the Third Workshop on Statistical Machine Translation</i> , pages 35–43.	
862		
863		
864		
865		
866	Alex Wang, Yada Pruksachatkun, Nikita Nangia, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R Bowman. 2019a. Super-glue: A stickier benchmark for general-purpose language understanding systems. <i>arXiv preprint arXiv:1905.00537</i> .	
867		
868		
869		
870		
871		
872	Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel Bowman. 2018. GLUE: A multi-task benchmark and analysis platform for natural language understanding . In <i>Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP</i> , pages 353–355, Brussels, Belgium. Association for Computational Linguistics.	
873		
874		
875		
876		
877		
878		
879		
880	Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R. Bowman. 2019b. GLUE: A multi-task benchmark and analysis platform for natural language understanding . In <i>7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019</i> . OpenReview.net.	
881		
882		
883		
884		
885		
886		
887	Adina Williams, Nikita Nangia, and Samuel R Bowman. 2017. A broad-coverage challenge corpus for sentence understanding through inference. <i>arXiv preprint arXiv:1704.05426</i> .	
888		
889		
890		
891	Mengzhou Xia, Antonios Anastasopoulos, Ruochen Xu, Yiming Yang, and Graham Neubig. 2020. Predicting performance for natural language processing tasks . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 8625–8646, Online. Association for Computational Linguistics.	
892		
893		
894		
895		
896		
897		
898	Zihuiwen Ye, Pengfei Liu, Jinlan Fu, and Graham Neubig. 2021. Towards more fine-grained and reliable NLP performance prediction . In <i>Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main</i>	
899		
900		
901		
902		
	<i>Volume</i> , pages 3703–3714, Online. Association for Computational Linguistics.	905
		906
	Yisong Yue, Thomas Finley, Filip Radlinski, and Thorsten Joachims. 2007. A support vector method for optimizing average precision . In <i>SIGIR 2007: Proceedings of the 30th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval, Amsterdam, The Netherlands, July 23-27, 2007</i> , pages 271–278. ACM.	907
		908
		909
		910
		911
	Jerrold H Zar. 1972. Significance testing of the spearman rank correlation coefficient. <i>Journal of the American Statistical Association</i> , 67(339):578–580.	912
		913
		914
	Xiang Zhang, Junbo Jake Zhao, and Yann LeCun. 2015. Character-level convolutional networks for text classification . <i>CoRR</i> , abs/1509.01626.	915
		916
		917
	Ming Zhong, Pengfei Liu, Yiran Chen, Danqing Wang, Xipeng Qiu, and Xuanjing Huang. 2020. Extractive summarization as text matching . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 6197–6208, Online. Association for Computational Linguistics.	918
		919
		920
		921
		922
		923

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 (ϕ_{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 (ϕ_{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_{type} is the number

of unique words, and n_{token} is the number of tokens. We use lexical richness¹⁹ to calculate the TTR for each sentence and then average them.

Language Mixedness Ratio (ϕ_{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 (ϕ_{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, ϕ_{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 (ϕ_{gerr}): The proportion of words with grammatical errors in the whole dataset. We adopt the detection tool²⁰ 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 (ϕ_{fre}): Flesch Reading Ease²¹ calculated by textstat²² is a way to describe the simplicity of a reader who can read a text. First, we calculate the ϕ_{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.

¹⁹<https://github.com/LSYS/lexicalrichness>

²⁰https://github.com/jxmorris12/language_tool_python

²¹https://en.wikipedia.org/wiki/Flesch%E2%80%93Kincaid_readability_tests

²²<https://github.com/shivam5992/textstat>