
Towards Comparable Active Learning

Thorben Werner *
Institute of Computer Science
University of Hildesheim
Universitätsplatz 1, 31141 Hildesheim
werner@ismll.de

Johannes Burchert*
Institute of Computer Science
University of Hildesheim
Universitätsplatz 1, 31141 Hildesheim
burchert@ismll.de

Prof. Lars Schmidt-Thieme*
Institute of Computer Science
University of Hildesheim
Universitätsplatz 1, 31141 Hildesheim
schmidt-thieme@ismll.uni-hildesheim.de

Abstract

Active Learning (AL) is a technique used in machine learning to improve data efficiency by selecting the most informative samples for annotation. In this paper we address the issue of inconsistent results in active learning literature. Authors of previous papers are constantly reporting significant performance improvements in their original work, while subsequent literature fails to reproduce those results. This inconsistency leads to a chaotic landscape of AL algorithms. We propose the first AL benchmark that tests algorithms in all three major domains of tabular, image and text data. Furthermore, we discuss a central problem for reproducing AL experiments with the default seeding setup that depends on a single seed per experiment. ...

1 Introduction

Version: 1.0

Active Learning (AL) plays an important role in our society that applies machine learning to more and more areas and therefore has a high demand for labeled data in more and more areas. A problem that concerns academic researchers and practitioners in businesses alike and even could be extended to education in schools and hobbyists around the world. On top of providing a principled way to labeled unlabeled datasets, active learning is one of the two major approaches besides semi-supervised learning to make deep learning models more data efficient by requiring only a limited set of manually labeled data. Both approaches are at their core orthogonal and can freely be combined and therefore we should continue our research efforts for both approaches.

Among others, the authors of [15] have pointed out severe inconsistencies in results of AL papers in recent years. In their supplementary materials they conducted a meta analysis of reported results of several different AL algorithms and found that all considered algorithms only provided significant lifts in their own original papers, while all following literature reported performances no better than uncertainty sampling, or in some cases no better than random sampling for the same algorithm. The result of these inconsistencies is a chaotic landscape of AL algorithms where every paper claims to

*Information Systems and Machine Learning Lab (ISMILL)

archive state-of-the-art results by significantly outperforming everyone else, while the vast majority of results proves to be non-reproducible.
To allow comparison between

1.1 Contributions

1. Evaluation of Active Learning algorithms on datasets from 4 different domains, including synthetic data that highlights principled shortcomings of existing approaches.
2. Novel experimental protocol for seeding the experiment with 3 different seeds to allow full control and reproducibility and analysis of how many restarts are required to converge to the true median performance reliably.
3. Simple algorithm for an Oracle-Curve that can be constructed greedily and does not rely on search.

2 Overview

We constrain our work on pool-based active learning where a pool of unlabeled samples is fixed at the start of each experiment and samples are chosen sequentially. Specifically, we are not experimenting on so-called batch active learning, where at each iteration multiple unlabeled samples are chosen at the same time. Even though batch AL is the more active research domain, it does not have a principled advantage over single-sample AL except speed of computation. Not only is the problem of optimizing a portfolio of unlabeled samples more complicated to solve, the algorithms also have systematically less information per sample to work with. For this reason we propose to focus more research effort on single-sample AL to find better algorithms in an environment that is simpler to solve and easier to control. Table 1 shows a feature comparison between our proposed benchmark and several existing benchmarks in the literature, as well as methodological AL papers with experiments on at least two data domains.

Paper	Sampling	# Datasets	Domains	Algorithms	Oracle
Beck et al. [2]	batch	4	1	7	-
Hu et al. [5]	batch	5	2	13	-
Li et al. [9]	batch	5	1	13	-
Zhou et al. [15]	batch	3	2	2	✓
Ours	single	9	4	6	✓

Table 1: Comparison of our benchmark with the existing literature

49

2.1 Problem Description

Version: 1.0

Given two spaces $\mathcal{X} := \mathcal{R}^M$ and $\mathcal{Y} := \mathcal{R}^C$, a sample $\mathcal{D}_1, \dots, \mathcal{D}_N \subseteq (\mathcal{X} \times \mathcal{Y})^*$ of sequences of pairs (x, y) from an unknown distribution p called datasets and a number $B \in \mathcal{N}$ with $B < |\mathcal{D}|$.
Given two functions $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathcal{R}$ called loss, and $A : (\mathcal{X} \times \mathcal{Y})^* \times \mathcal{X}^* \rightarrow \mathcal{Y}^{\mathcal{X}}$ called learning algorithm, find a function

$$a : (\mathcal{X} \times \mathcal{Y})^* \times \mathcal{X}^* \rightarrow \{0, 1\}^* \quad (\text{equivariant in the second argument})$$

called acquisition function, s.t. the expected loss of a model learned on all predictors plus B acquired targets is minimal:

$$\begin{aligned} & \mathbb{E}_{\mathcal{D}_{\text{train}}, \mathcal{D}_{\text{test}} \sim p} \text{avg}_{(x, y) \in \mathcal{D}_{\text{test}}} \ell(y, \hat{y}(x)) \\ & \text{with } \hat{y} := A((\mathcal{D}_{\text{train}_{n_1}}, \dots, \mathcal{D}_{\text{train}_{n_B}}), \mathcal{D}_{\text{train}}|_{\mathcal{X}}) \\ & n_b := \text{index}(a((\mathcal{D}_{\text{train}_{n_1}}, \dots, \mathcal{D}_{\text{train}_{n_{b-1}}}), \mathcal{D}_{\text{train}}|_{\mathcal{X}})), \quad b \in 1:B \end{aligned}$$

(We would need to switch from lowest expected loss to highest AUC)

Since combinatorial problem of finding the optimal subset $\mathcal{D}_{\text{train}}$ is computational not feasible, we

allow sequential construction of $\mathcal{D}_{\text{train}}$.
 To construct the active learning setting from existing labeled datasets, we first split it into $\mathcal{D}_{\text{train}}$, \mathcal{D}_{val} and $\mathcal{D}_{\text{test}}$ and then suppress the labels \mathcal{Y} of $\mathcal{D}_{\text{train}}$ to form the unlabeled pool $\mathcal{U} := \mathcal{X}^*$ and form and initial labeled pool by uniformly sampling k number of instances per class from \mathcal{U} and recovering their label $\mathcal{L} := (\mathcal{X}, \mathcal{Y})^{k \times C}$.
 Following [15], the quality of an active learning algorithm is evaluated by an "anytime" protocol that incorporates classification performance at every iteration, not just the final performance after the budget is exhausted. We employ the normalized area under the accuracy curve (AUC):

$$\text{auc}(\mathcal{U}, \mathcal{L}, \hat{y}, B) := \frac{1}{B} \sum_{i=1}^B \text{Acc}(y_{\text{test}}, \hat{y}_i(x_{\text{test}})) \quad (1)$$

Where \hat{y}_i is the retrained classification model after the i -th instance was selected.

Algorithm 1 Oracle

Require: \mathcal{U} ▷ Unlabeled Pool
Require: τ ▷ Unlabeled Sample Size

- 1: $\mathcal{L}^{(1)} \leftarrow \text{seed}(\mathcal{U})$
- 2: $\mathcal{U}^{(1)} \leftarrow \mathcal{U} \setminus \mathcal{L}^{(1)}$
- 3: **for** $i := 1 \dots B$ **do**
- 4: $\text{acc}^{(i)} \leftarrow \text{Retrain}(\mathcal{L}^{(i)})$
- 5: $u^{(i)} \sim \text{unif}(1 : |\mathcal{U}^{(i)}|)$
- 6: $r^* \leftarrow 0$
- 7: **for** $j := 1 \dots \tau$ **do**
- 8: $y^{(j)} \leftarrow \text{label}(u_j^{(i)})$
- 9: $\mathcal{L}^{(j)} \leftarrow \mathcal{L}^{(i)} \cup \{(u_j^{(i)}, y^{(j)})\}$
- 10: $\text{acc}^{(j)} \leftarrow \text{Retrain}(\mathcal{L}^{(j)})$
- 11: $r^{(j)} \leftarrow \text{acc}^{(j)} - \text{acc}^{(i)}$
- 12: **if** $r^{(j)} > r^*$ **then**
- 13: $r^* \leftarrow r^{(j)}$
- 14: $j^* \leftarrow j$
- 15: **end if**
- 16: **end for**
- 17: **if** $r^* = 0$ **then**
- 18: $j^* \leftarrow \text{margin}(u^{(i)}, \hat{y}_\theta)$
- 19: **end if**
- 20: $y^{(i)} \leftarrow \text{label}(u_{j^*}^{(i)})$
- 21: $\mathcal{L}^{(i+1)} \leftarrow \mathcal{L}^{(i)} \cup \{(u_{j^*}^{(i)}, y^{(i)})\}$
- 22: $\mathcal{U}^{(i+1)} \leftarrow \mathcal{U}^{(i)} \setminus \{u_{j^*}^{(i)}\}$
- 23: **end for**
- 24: **return** $\frac{1}{B} \sum_{i=1}^B \text{acc}^{(i)}$

Algorithm 2 Retrain

Require: \mathcal{L} ▷ Labeled Pool
Require: \mathcal{D}_{val} ▷ Validation Data
Require: $\mathcal{D}_{\text{test}}$ ▷ Test Data
Require: \hat{y}_θ ▷ Class. Model
Require: e^{max} ▷ Maximum Epochs

- 1: $\text{loss}^* \leftarrow \infty$
- 2: **for** $i := 1 \dots e^{\text{max}}$ **do**
- 3: $\theta_{i+1} \leftarrow \theta_i - \eta \nabla_{\theta} \ell(\mathcal{L}, \hat{y}_\theta)$
- 4: $\text{loss}_i \leftarrow \ell(\mathcal{D}_{\text{val}}, \hat{y}_\theta)$
- 5: **if** $\text{loss}_i < \text{loss}^*$ **then**
- 6: $\text{loss}^* \leftarrow \text{loss}_i$
- 7: **else**
- 8: Break
- 9: **end if**
- 10: **end for**
- 11: **return** $\text{Acc}(\mathcal{D}_{\text{test}}, \hat{y}_\theta)$

In Alg. 1 $\text{Retrain}(\mathcal{L}^{(i)})$ is shorthand for $\text{Retrain}(\mathcal{L}^{(i)}, \mathcal{D}_{\text{test}}, \hat{y}_\theta, e^{\text{max}})$, $\text{label}(u_j^{(i)})$ recovers the label of instance $u_j^{(i)}$ and $\text{margin}(u^{(i)}, \hat{y}_\theta)$ applies margin sampling to select as a fallback, when the oracle has not found another option that induces a positive change in classification performance.

3 Related Work

Version: Braindump

Many different algorithms have been proposed for active learning. In this work we focus on those approaches that have shown consistent results over the years as well as newer approaches that have demonstrated significant lifts in their initial experiments. AL algorithms can be categorized into two classes: Geometric approaches and uncertainty-based approaches. Geometric approaches include CoreSet [13] and TypiClust [4], which use clustering techniques to partition the data and then sample their unlabeled points based on the clusters. Uncertainty-based approaches include classic uncertainty sampling (based on Shannon-Entropy and the margin-score), BALD [7] and BADGE [1], which use metrics to measure the classifiers state.

Some previous work also aimed to provide a benchmark suite for active learning: The authors of [2] and [9] both focus on active learning in the image domain. While [2] discuss a new metric to measure AL performance, which they call "Label Efficiency" and provide experiments on many common configurations for data preparation, model training and other hyperparameters, [9] focuses on combined approaches of AL and semi-supervised learning to aid model training. The authors of [5] study models that are learned with AL techniques in the image and text domain. They test for several different properties of the models including robustness, response to compression techniques and final performance.

4 Methodology

4.1 Evaluation

Version: 1.0

We compare different AL algorithms based on their median AUC score (Eq. 1) across multiple restarts of the experiment. Each restart will retain the train/test split (often given by the dataset itself), but introduces a new validation split to mimic the leave-one-out protocol for cross-validation. The AUC incorporates performance in early stages (low budget) as well as capabilities to push the classifier in later stages (high budget). A good AL algorithm should be able to perform well in both scenarios.

Since AUC is dependent on the chosen budget, we need a general rule on how to set this hyperparameter that does not inherently benefit a subset of algorithms. In this work, we choose the budget per dataset to be the first point at which any algorithm (except oracle) manages to reach a percentage of the upper bound performance measured on the full dataset. Even though we would like to propose a single percentage value for all datasets, we found that different data modalities and use cases need different percentages to produce sensible budgets. We propose the following values: **Tabular**: 99%, **Image**: 90% and **Text**: 95%.

Additionally, we provide evidence in Fig. 1 that previous works have not evaluated their experiments with a sufficient number of restarts. To create Fig. 1 we used all our 50 runs from the margin/random sampling algorithm on the splice dataset. From these 50 runs we uniformly sampled subsets of runs and calculated the median AUC for this subset. One AUC value like this corresponds to one cross-validated experiment sampled from the distribution of experiments that are restarted exactly this many times. To create one slice in Fig. 1, we drew 50 samples from this distribution. Each box-plot represents the variance of an evaluation if conducted with the respective number of restarts. We can clearly observe that low repetitions (< 10) provide an uncertain evaluation where lucky and unlucky draws of the same experiment give drastically different median AUC values. To combat this uncertain evaluation, we propose to repeat every experiment 50 times, to arrive at the true median AUC for each algorithm.

4.2 Reproducibility

Version: 1.0

A big focus in this work is to provide an experimental setup that is fully reproducible independent of the dataset, classification model or AL algorithm used. In our opinion, an evaluation on one dataset with a given seed should always be done on the same validation split as well as using the same random state for all included systems like the mini batch sampler for the classifier model training. Even though different AL algorithms will pick different samples, making them unavailable for sampling in earlier or later batches, the theoretical decision tree for every possible choice in every iteration i should stay the same. Since every possible trajectory cannot be precomputed and stored to disk, we need to resort to seeding the evaluation. The default choice of setting a global seed at the start of the experiment is not sufficient here, since a single additional random draw from the random number generator completely changes the behavior of all other systems. This additional random number might be drawn during the initialization of the classification model, the AL algorithm, or even during every AL iteration if Ω is stochastic. This problem also applies to the initialization

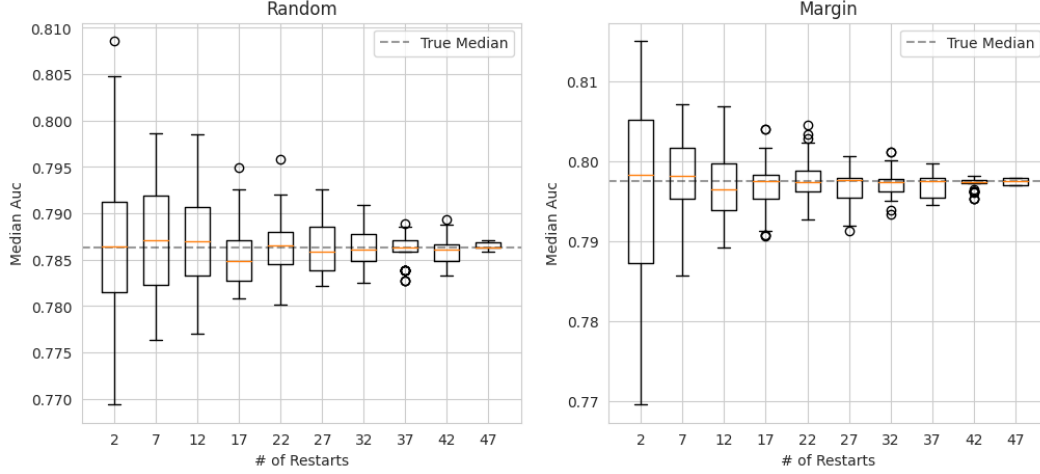


Figure 1: Random draws from an experimental distribution on the Splice datasets with different numbers of repetitions. Each point on the Y-axis represents a cross-validated result that could have been reported in a paper. This analysis shows the drastic differences in performance one could observe even when repeating an experiment 2-10 times.

of the classification model θ and the initialization and querying of Ω . The desired control over all these processes can only be archived by assigning a separate random number generator to all these processes. To the best of our knowledge, we are the first work that discusses this issue and proposes a solution for it. We hypothesize that the insufficient setup with global seeds contributes to the ongoing problem of inconsistent results of AL algorithms in different papers.

In summary, we introduce three different seeds: s_Ω for the acquisition function, $s_{\mathcal{D}}$ for dataset splitting and mini batch sampling and s_θ for model initialization and sampling of dropout masks. Unless stated otherwise, we will keep s_Ω fixed for restarts of the same experiment, while $s_{\mathcal{D}}$ and s_θ are incremented by 1 between restarts to introduce stochasticity into our framework.

4.3 Oracle

Version: 1.0

Posing active learning as a sequence ordering problem, the oracle sequence for a given combination of dataset, model and training procedure would be the sequence that induces the highest AUC score for a given budget. However, since this combinatorial problem is not solvable for realistic datasets, previous works have proposed approximations to this oracle sequence. [15] has used simulated annealing to search for the optimal sequence and used the best solution found after a fixed time budget. Even though their reported performance curves display a significant lift over all other algorithms, we found the computational cost of reproducing this oracle for all our datasets to be prohibitive (The authors reported the search to take several days per dataset on 8 V100 GPUs).

In this paper we propose a greedy oracle algorithm that constructs an approximation of the optimal sequence in an iterative fashion. Our oracle simply tests every data point in the provided sample of unlabeled points by fitting the classifier and directly measuring the resulting validation performance. The point with the best validation performance is selected and added to the labeled pool for that iteration. We noticed that this oracle is overfitting on the validation set, resulting in stagnating or even decreasing performance curves in later AL iterations. To circumvent this problem, we introduced margin sampling as a fallback option for the oracle. Whenever the oracle does not find an unlabeled point that results in an increase in performance (indicating an overfitting position), it defaults to margin sampling in that iteration. The pseudocode for our oracle can be found in Alg. 1.

5 Implementation Details

5.1 Available Information

At each iteration i the AL algorithm needs to pick an unlabeled datapoint based on a fixed set of information $\{\mathcal{L}^{(i)}, \mathcal{U}^{(i)}, B, |\mathcal{L}^{(i)}| - |\mathcal{L}^{(1)}|, \text{acc}^{(i)}, \text{acc}^{(1)}, \theta^{(i)}, \text{opt}_\theta\}$, where $\theta^{(i)}$ is the current classifier and opt_θ is the optimizer used to fit $\theta^{(i)}$. We allow algorithms to derive additional information of this set like predictions of the classifier, K-Means clustering or even training new classifiers. However, the algorithm may not incorporate external information like other datasets, queries to recover additional labels, or the test/validation set.

5.2 Sampling Strategies

We selected AL algorithms that have good performances reported by multiple different sources. To ensure a fair comparison we fixed the training process of our classification model as well as the set of available information for the algorithm and selected only those that can work under these restrictions:

Uncertainty Sampling Tries to find the sample that the classifier is most uncertain about. For our benchmark we use entropy and margin (a.k.a. best-vs-second-best) sampling.

BALD [7] Applies the query-by-committee strategy of model ensembles to a single model by interpreting the model’s parameters as distributions and then sample multiple outputs from them via Monte-Carlo dropout.

BADGE [1] Uses gradient embeddings of unlabeled points to select samples where the classifier is expected to change a lot. The higher the magnitude of the gradient the higher is the expected improvement in model performance.

Coreset [13] Employs K-Means clustering to try to cover all modalities (is this the right word?) that are observed in the data. Selects the unlabeled sample that is the furthest away from all cluster centers. In this work we use the greedy variant of Coreset and we also encode the data with the current classifier θ_i before clustering.

TypiClust Relies on clustering similar to Coreset but proposes a new measure called ”Typicality” to select unlabeled centers. Tries to select points that are in the densest regions of clusters that do not contain labeled samples yet. For this work we also use the current classifier θ_i to encode data if necessary (raw images or text).

5.2.1 Honorable Mentions

Learning Loss for AL Introduces an updated training of the classification model with an auxiliary loss and therefore cannot be compared fairly against classification models without this boosted training regime.

5.3 Choosing the Classifier

Traditionally, the classifier is chosen per dataset so that it is capable of solving the dataset close to the SOTA performance reported in the literature. Similar to our hypothesis in section 4.1 we hypothesize that AL algorithms will perform similarly on small classifiers and more complex ones, so that the overall ranking of algorithms stays the same. **TODO: Transform this into full hypothesis incl. formatting?**

On the basis of this hypothesis we opt to use smaller classifiers that still solve the dataset to a reasonable degree. Smaller classifiers also require fewer labeled datapoints to acquire performance close to the upper bound. For every dataset the chosen architecture’s hyperparameters are optimized by to

199 archive maximum upper bound performance. For an overview of architectures and hyperparameters
200 please refer to Appendix C.

201 5.4 Training the Classifier

202 Version: 0.1

203 Generally, the classification model can be trained in two ways. Either you reset the parameters
204 after each AL iteration and train the classifier from scratch with the updated labeled set $\mathcal{L}^{(i)}$, or you
205 retain the previous state and only fine-tune the classifier on $\mathcal{L}^{(i)}$ for a reduced number of epochs.
206 In this work we use the fine-tuning method for raw datasets to save computation, while we use the
207 from-scratch training for embedded dataset, since they have very small classifiers and this method
208 generally produces better results. Our fine-tuning scheme always trains for at least one epoch and
209 employs an aggressive early stopping after that. The early stopping has patience 0, so it will stop as
210 soon as the validation loss does no longer decrease.

211 6 Experiments

212 6.1 Datasets

213 Version: 1.0

214 For all our datasets we use the pre-defined train / test splits, if given. In the remaining cases, we
215 define test sets upfront and store them into separate files to keep them fixed across all experiments.
216 The validation set is split during experiment-time and depends on the dataset-seed.

217 **Tabular** We use **Splice**, **DNA** and **USPS** from LibSVMTools [11]. All three datasets are normal-
218 ized between [0, 1].

219 **Image** We use **FashionMNIST** [14] and **Cifar10** [8]. Both datasets are normalized according to
220 their standard protocols.

221 **Text** We use **News Category** [10] and **TopV2** [3]. For News Category we use the 15 most common
222 categories as indicated by its Kaggle site. We additionally drop sentences above 80 words to reduce
223 the needed padding (retaining 99,86% of the data). For TopV2, we are only using the "alarm"
224 domain. Both datasets are encoded with pre-trained GloVe embeddings [12]. Since neither dataset
225 provided a fixed test set, we randomly split 7000 datapoints into a test set.

226 We would like to point out that these datasets can be considered "toy-datasets" and therefore not
227 relevant for practical purposes. This might be true if we aimed to develop novel classification models
228 on these datasets, however we are solely focused on comparing different AL algorithms in this paper.
229 Our core assumption is that a well-performing algorithm in our benchmark will transfer well into
230 more practical use-cases.

231 Adapting the experimental setting from [4] we offer all our datasets in the raw setting as well as pre-
232 encoded by a fixed embedding model that was trained by unsupervised contrastive learning. The text
233 datasets are an exception, as they are only offered in their encoded form. The pre-encoded datasets
234 offer multiple advantages in that they require smaller classification models (i.e. small MLPs) and
235 in general much less sampled datapoints to archive close to upper bound performance (measured
236 with the full dataset). They unify experimental setups for different datasets and therefore make them
237 much more comparable.

238 6.2 Results

239 From Fig. 2 we notice drastically different qualities for the same AL algorithm for different datasets.
240 We would like to highlight that both datasets are tabular from the medical domain with similar
241 number of features and classes, yet we see that i.e. BALD is the best algorithm for Splice and

	Splice	DNA	USPS
Oracle	0.830 +- 0.01	0.836 +- 0.02	0.823 +- 0.01
SAL	0.799 +- 0.01	0.797 +- 0.03	0.809 +- 0.01
Coreset	0.800 +- 0.01	0.795 +- 0.03	0.787 +- 0.02
TypiClust	0.790 +- 0.01	0.771 +- 0.04	0.761 +- 0.02
MarginScore	0.797 +- 0.02	0.795 +- 0.04	0.808 +- 0.01
ShannonEntropy	0.799 +- 0.02	0.794 +- 0.04	0.807 +- 0.01
RandomAgent	0.788 +- 0.01	0.765 +- 0.03	0.772 +- 0.01
Badge	0.807 +- 0.01	0.769 +- 0.06	0.797 +- 0.02
BALD	0.811 +- 0.01	0.743 +- 0.04	0.717 +- 0.05

Table 2: AUC values for all algorithms on the tabular datasets. Higher is better. **Is the STD even helping us here? Since we repeated 50 times, we can use the values as "true medians"**

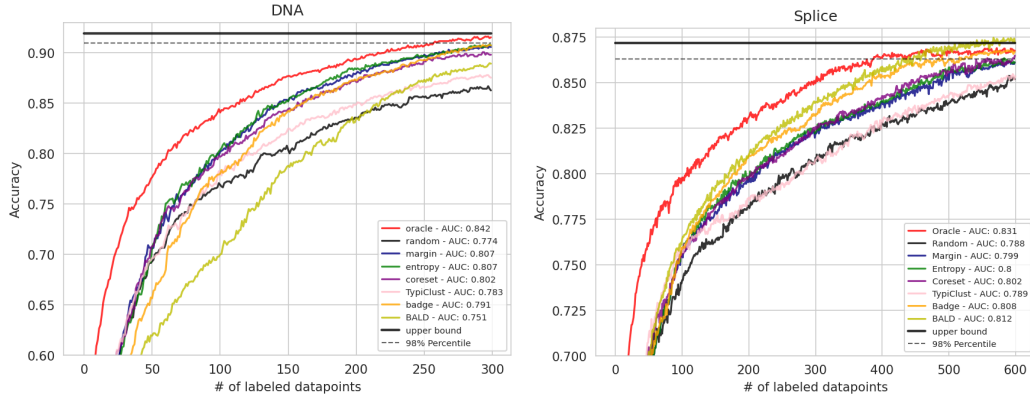


Figure 2: Results for all algorithms on Splice and DNA, both from the tabular domain. Even within one domain, the performance of the same algorithm can vary drastically.

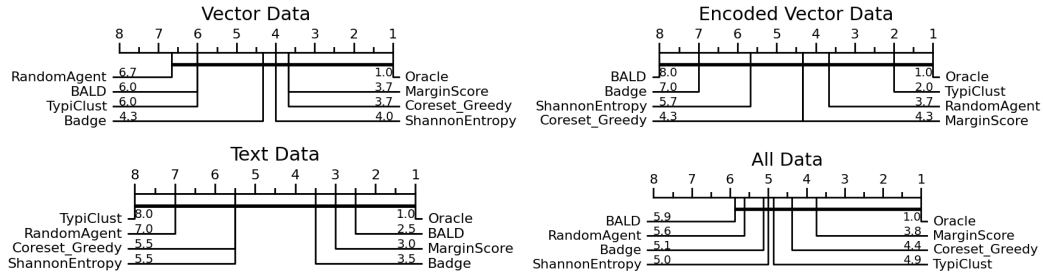


Figure 3: Critical Difference Diagram for all algorithms grouped by domain and all domains combined. Ranks are computed based on median AUC for each algorithm and dataset combination. Lower ranks are better.

242 the worst algorithm for DNA. These inconsistencies are present between the datasets of all our
243 tested domains, further highlighting the difficulties for comparing AL algorithms in terms of average
244 performance. In order to provide a meaningful analysis of which algorithm can be expected to
245 perform best on average we ranked the algorithm for each dataset based on their median AUC and
246 displayed these rankings in critical difference diagrams [6]. In Fig. 3 we report the rankings split by
247 domain as well as across all domains (excluding the toydata).

248 **7 Conclusion**

249 **8 Limitations and Future Work**

- 250 No batch AL
- 251 No learned algorithms
- 252 No improved classifier training (semi-supervised, etc.)

References

- [1] Jordan T Ash, Chicheng Zhang, Akshay Krishnamurthy, John Langford, and Alekh Agarwal. Deep batch active learning by diverse, uncertain gradient lower bounds. In *International Conference on Learning Representations*, 2020.
- [2] Nathan Beck, Durga Sivasubramanian, Apurva Dani, Ganesh Ramakrishnan, and Rishabh Iyer. Effective evaluation of deep active learning on image classification tasks. *arXiv preprint arXiv:2106.15324*, 2021.
- [3] Xilun Chen, Asish Ghoshal, Yashar Mehdad, Luke Zettlemoyer, and Sonal Gupta. Low-resource domain adaptation for compositional task-oriented semantic parsing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. Association for Computational Linguistics, 2020.
- [4] Guy Hacohen, Avihu Dekel, and Daphna Weinshall. Active learning on a budget: Opposite strategies suit high and low budgets. *arXiv preprint arXiv:2202.02794*, 2022.
- [5] Qiang Hu, Yuejun Guo, Maxime Cordy, Xiaofei Xie, Wei Ma, Mike Papadakis, and Yves Le Traon. Towards exploring the limitations of active learning: An empirical study. In *2021 36th IEEE/ACM International Conference on Automated Software Engineering (ASE)*, pages 917–929. IEEE, 2021.
- [6] Hassan Ismail Fawaz, Germain Forestier, Jonathan Weber, Lhassane Idoumghar, and Pierre-Alain Muller. Deep learning for time series classification: a review. *Data Mining and Knowledge Discovery*, 33(4):917–963, 2019.
- [7] Andreas Kirsch, Joost Van Amersfoort, and Yarin Gal. Batchbald: Efficient and diverse batch acquisition for deep bayesian active learning. *Advances in neural information processing systems*, 32, 2019.
- [8] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- [9] Yu Li, Muxi Chen, Yannan Liu, Daojing He, and Qiang Xu. An empirical study on the efficacy of deep active learning for image classification. *arXiv preprint arXiv:2212.03088*, 2022.
- [10] Rishabh Misra. News category dataset. *arXiv preprint arXiv:2209.11429*.
- [11] Information Engineering Graduate Institute of Taiwan University. Libsvmtools.
- [12] Jeffrey Pennington, Richard Socher, and Christopher D. Manning. Glove: Global vectors for word representation. In *Empirical Methods in Natural Language Processing (EMNLP)*, pages 1532–1543, 2014.
- [13] Ozan Sener and Silvio Savarese. Active learning for convolutional neural networks: A core-set approach. *arXiv preprint arXiv:1708.00489*, 2017.
- [14] Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms. *arXiv preprint arXiv:1708.07747*, 2017.
- [15] Yilun Zhou, Adithya Renduchintala, Xian Li, Sida Wang, Yashar Mehdad, and Asish Ghoshal. Towards understanding the behaviors of optimal deep active learning algorithms. In *International Conference on Artificial Intelligence and Statistics*, pages 1486–1494. PMLR, 2021.

	SpliceEncoded	DNAEncoded	USPSEncoded
Oracle	0.754 +- 0.02	0.726 +- 0.02	0.674 +- 0.01
SAL	0.675 +- 0.03	0.640 +- 0.04	0.634 +- 0.01
Coreset	0.690 +- 0.02	0.644 +- 0.05	0.607 +- 0.02
TypiClust	0.695 +- 0.02	0.660 +- 0.03	0.643 +- 0.01
MarginScore	0.675 +- 0.03	0.643 +- 0.05	0.632 +- 0.02
ShannonEntropy	0.673 +- 0.03	0.638 +- 0.05	0.626 +- 0.02
RandomAgent	0.680 +- 0.03	0.633 +- 0.04	0.594 +- 0.02
Badge	0.670 +- 0.04	0.600 +- 0.07	0.597 +- 0.02
BALD	0.660 +- 0.04	0.597 +- 0.06	0.652 +- 0.01

Table 3: AUC values for all algorithms on the encoded tabular datasets. Higher is better.

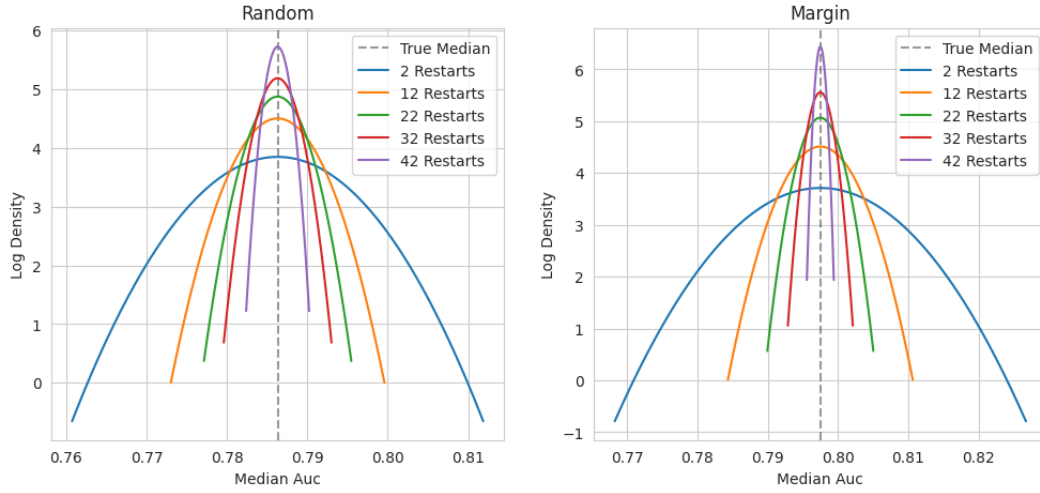
	Cifar10Encoded	FashionMnistEncoded
Oracle	0.699 +- 0.01	0.721 +- 0.01
SAL	0.623 +- 0.01	0.667 +- 0.01
Coreset	0.641 +- 0.01	0.674 +- 0.01
TypiClust	0.627 +- 0.01	0.638 +- 0.01
MarginScore	0.641 +- 0.01	0.676 +- 0.01
ShannonEntropy	0.637 +- 0.02	0.680 +- 0.01
RandomAgent	0.598 +- 0.02	0.633 +- 0.01
Badge	0.645 +- 0.02	0.681 +- 0.01
BALD	0.633 +- 0.01	0.666 +- 0.01

Table 4: AUC values for all algorithms on the encoded image datasets. Higher is better.

	TopV2	News
Oracle	0.860 +- 0.01	0.448 +- 0.01
SAL	0.831 +- 0.01	0.312 +- 0.01
Coreset	0.816 +- 0.02	0.355 +- 0.02
TypiClust	0.718 +- 0.02	0.314 +- 0.01
MarginScore	0.821 +- 0.02	0.357 +- 0.01
ShannonEntropy	0.794 +- 0.02	0.342 +- 0.01
RandomAgent	0.776 +- 0.02	0.349 +- 0.01
Badge	0.824 +- 0.01	0.343 +- 0.01
BALD	0.825 +- 0.02	0.347 +- 0.01

Table 5: AUC values for all algorithms on the text datasets. Higher is better.

293 **B Alternative Plot for Restarts Ablation**

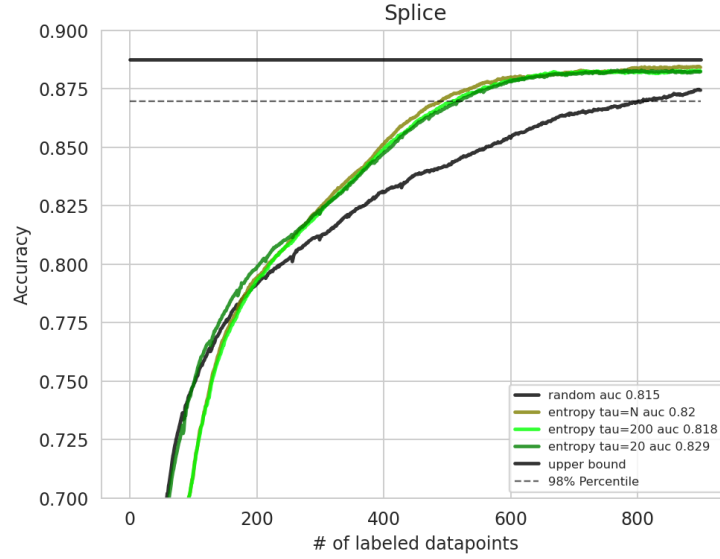


294 **C Hyperparameters per Dataset**

Dataset	Classifier	Optimizer	LR	Weight Decay	Dropout
Splice	[24, 12]	NAdam	1.2e-3	5.9e-5	0
SpliceEnc.	linear	NAdam	6.2e-4	5.9e-6	0
DNA	[24, 12]	NAdam	3.9e-2	3.6e-5	0
DNAEnc	linear	NAdam	1.6e-3	4e-4	0
USPS	[24, 12]	Adam	8.1e-3	1.5e-6	0
USPS	linear	NAdam	7.8e-3	1.9e-6	0
FashionMnist	ResNet18	NAdam	1e-3	0	0
FashionMnistEnc	linear	Adam	1.6e-3	1e-5	5e-2
Cifar10	ResNet18	NAdam	1e-3	0	0
Cifar10Enc	linear	NAdam	1.7e-3	2.3e-5	0
TopV2	BiLSTM	NAdam	1.5e-3	1.7e-7	5e-2
News	BiLSTM	NAdam	1.5e-3	1.7e-7	5e-2

Table 6: Classifier architectures and optimized hyperparameters per dataset. Numbers in brackets signify a MLP with corresponding hidden layers.

295 **D Comparison of different sample sizes**



296 **E AL Pseudocode**

Algorithm 3 Active Learning

Require: \mathcal{U}	▷ Unlabeled Pool
Require: τ	▷ Unlabeled Sample Size
Require: Ω	▷ AL Agent
Require: ω	▷ Environment State function
	▷ Create the initial labeled set

```

1:  $\mathcal{L}^{(1)} \leftarrow \text{seed}(\mathcal{U})$ 
2:  $\mathcal{U}^{(1)} \leftarrow \mathcal{U}$ 
3: for  $i := 1 \dots B$  do
4:    $\text{acc}^{(i)} \leftarrow \text{Retrain}(\mathcal{L}^{(i)})$ 
5:    $a^{(i)} \leftarrow \Omega(\mathcal{U}^{(i)})$ ;  $a \in 1 : |\mathcal{U}^{(i)}|$ 
6:    $y^{(i)} \leftarrow \text{label}(\mathcal{U}_a^{(i)})$ 
7:    $\mathcal{L}^{(i+1)} \leftarrow \mathcal{L}^{(i)} \cup \{(\mathcal{U}_a^{(i)}, y^{(i)})\}$ 
8:    $\mathcal{U}^{(i+1)} \leftarrow \mathcal{U}^{(i)} \setminus \{\mathcal{U}_a^{(i)}\}$ 
9: end for
10: return  $\frac{1}{B} \sum_{i=1}^B \text{acc}^{(i)}$ 

```
