
TODO: Funky active learning pun

Lukas Weber¹

Abstract

A lot of unlabelled data generally available but annotation is costly. A goal in fields where labelled data is not available is to maximize performance with minimum amount of labeled data instances. Active learning selects instances for labelling based on a selection strategy.

- Objective: Compare active learning query methods using modAL.
- Methods: Experiments with Logistic Regression (MNIST) and a CNN (CIFAR).
- Experiments: Vary initial dataset sizes (small, moderate, large).
- Key findings: Highlight main conclusions about performance in different scenarios.

1. Introduction

Complex tasks in supervised learning usually require large annotated datasets. However, the cost and time needed to label samples can be high, especially when expert knowledge in the according field is required.

Active learning aims to reduce the overall labeling cost by selecting the most informative samples for annotation by an oracle and thus enable using machine learning for tasks where annotated data is sparse.

The effectiveness of active learning relies heavily on the choice of query strategy, as it determines which samples are selected for annotation and thus directly influences model performance. In pool-based active learning, a model is trained on a small labeled dataset and used to evaluate a large pool of unlabeled data. Informative samples are selected based on a query strategy, labeled by an oracle, and added to the training set. The model is retrained, and the process is repeated until a stopping criterion is met.

In stream-based active learning, unlabeled data arrives sequentially. Thus, the model decides whether to query the

label of an incoming instance using an evaluation metric like the uncertainty of the prediction and checking whether a threshold is reached. If not, the instance is discarded.

This report focuses on comparing several pool-based as well as stream-based query strategies.

Related Work: Several previous works explore different query methods for several tasks. (Schröder et al., 2022)

- Previous works have explored strategies such as ..., ..., While they have been widely tested, most evaluations focus on single dataset scenarios or specific models
- Then the gaps: However, few studies analyze the impact of initial dataset size on query strategy performance, and even fewer investigate their behavior on simple versus complex models.

Previous papers:

- (Schröder et al., 2022) Compares various uncertainty based query strategies in context of fine-tuning transformer models for text classification. Does not compare across different sizes of the pool, and for simple models.
- (Ueno et al., 2023) TOWARDS FUTURE DEEP ACTIVE LEARNING”: Most papers concentrate on cifar or mnist, this paper evaluates six different datasets, including medical and visual inspection images. Also tackles the problem that experimental settings are not standardized, making evaluation of existing methods difficult. Also did a verification experiment where fully trained models are used to select most informative samples according to various query strategies and construct a labeled dataset. This is done to isolate and evaluate the effectiveness of query strategies indepently of training process/ underfitting or randomness due to early-stage training.
- (Zhan et al., 2022): compares different query methods on mnist and cifar with a cnn (resnet18). Does not compare impact of different initial train set size and strategies on shallow classifiers.
- (Werner et al., 2023) highlight critical challenges in active learning research, including the lack of repro-

^{*}Equal contribution ¹MSc Computer Science, Tübingen, GER. Correspondence to: Lukas Weber <lukas2.weber@student.uni-tuebingen.de>.

ducible experiments and fair comparisons across domains. They propose a benchmark framework evaluating active learning strategies on tabular, image, and text datasets using robust evaluation protocols to mitigate variance and ensure comparability. While their work focuses on a wide range of domains, this paper only compares methods in the image domain. Furthermore, this paper focuses on the image domain. Specifically, it compares strategies for a simple logistic regression model and a CNN across different starting sizes of the annotated dataset.

By examining different starting sizes, this work simulates various real-world scenarios: (1) starting with no labeled data, (2) having a small pool of labeled data, and (3) working with a well-trained model to identify data points for fine-tuning.

Modal: description and relevance for research maybe? The modAL framework(ref) has simplified the implementation of active learning pipelines, yet its comparative use across different strategies remains underexplored.

Prior studies

- existing studies that compare active learning methods
- do they have gaps in the literature? (maybe: comparison under different dataset sizes, use with cnns, comparison of simple or complicated models)

Contributions Comparative study of popular active learning query strategies using modAL framework. Contributions are as follows:

- systematic evaluation on two datasets: MNIST (Logistic Regression) and CIFAR (CNN)
- analysis across different scenarios: small, moderate and large initial datasets
- insight into performance impact of model complexity and initial dataset size on active learning performance

Structure of the paper

- outline paper

2. Methods

2.1. Datasets

- MNIST and CIFAR describe
- processing steps
MNIST: Vectorized, flattened (describe input size), normalized

2.2. Models

- Logistic Regression: Details on implementation and hyperparameters, how were they chosen?
- same for CNN

2.3. Active Learning

- How are the splits (initial pool, initial training data) generated
- number of iterations for active learning
- Used active learning query strategies

Metrics for evaluation (Accuracy, train error, confusion matrix, ...)

Tool: modAL, scikit

3. Experiments

Setup: details for reproducibility (random number generators)? Random seed = 42 (Werner et al., 2023) mention reproducibility problems in models even when setting random number generators before each run. By training models multiple times in succession and checking whether the same configurations lead to different outcomes, it is shown that in this scenario, reproducibility is guaranteed. **TODO: Schnell ein skript aufsetzen und dann finale accuracy und so weiter vergleichen.**

Did I perform multiple runs?

Scenarios

- describe experiments and their settings and purposes

4. Results

Quantitative results: Tables or plots comparing the query strategies across scenarios. Also maybe key performance differences (accuracy vs number of labeled examples)

5. Discussion

Qualitative analysis: Why certain strategies perform better or worse in specific settings. Trends between simple models and complex models

Impact of initial dataset size

Most methods require knowledge about the balance of the dataset and might be disadvantageous to use in case of a disbalanced dataset. In real world scenarios, this is not known.

Only image-based datasets were tested. Not vector- or text-based.

6. Conclusion

Summary of findings and insights (maybe that a strategy works best with small datasets while another one is robust on bigger ones)

Implications: practical recommendations?

Future work? (More diverse datasets, exploring additional strategies, applying in real-world tasks)

7. Cool ideas

- Show the typical framework of active learning in a figure in the introduction

8. Good formulations

From the active learning survey: Entropy is an information-theoretic measure that represents the amount of information needed to “encode” a distribution. As such, it is often thought of as a measure of uncertainty or impurity in machine learning.

References

- Schröder, C., Niekler, A., and Potthast, M. Revisiting uncertainty-based query strategies for active learning with transformers, 2022. URL <https://arxiv.org/abs/2107.05687>.
- Ueno, S., Yamada, Y., Nakatsuka, S., and Kato, K. Benchmarking of query strategies: Towards future deep active learning, 2023. URL <https://arxiv.org/abs/2312.05751>.
- Werner, T., Burchert, J., and Schmidt-Thieme, L. Towards comparable active learning, 2023. URL <https://arxiv.org/abs/2311.18356>.
- Zhan, X., Wang, Q., hao Huang, K., Xiong, H., Dou, D., and Chan, A. B. A comparative survey of deep active learning, 2022. URL <https://arxiv.org/abs/2203.13450>.