



**FACULTY
OF MATHEMATICS
AND PHYSICS**
Charles University

MASTER THESIS

Mitchell Borchers

**Active learning in E-Commerce
Merchant Classification using Website
Information**

Department of Theoretical Computer Science and Mathematical Logic

Supervisor of the master thesis: Mgr. Marta Vomlelová, Ph.D.

Study programme: Artificial Intelligence

Study branch: IUIPA

Prague 2023

I declare that I carried out this master thesis independently, and only with the cited sources, literature and other professional sources. It has not been used to obtain another or the same degree.

I understand that my work relates to the rights and obligations under the Act No. 121/2000 Sb., the Copyright Act, as amended, in particular the fact that the Charles University has the right to conclude a license agreement on the use of this work as a school work pursuant to Section 60 subsection 1 of the Copyright Act.

In date
Author's signature

Dedication.

Title: Active learning in E-Commerce Merchant Classification using Website Information

Author: Mitchell Borchers

Department: Department of Theoretical Computer Science and Mathematical Logic

Supervisor: Mgr. Marta Vomlelová, Ph.D., Department of Theoretical Computer Science and Mathematical Logic

Abstract: Data and the collection and analysis of data has become an important part of everyday life. For example, navigation, e-commerce, and social media all make use of immense amounts of data to provide users with suggestions on the best routes to take, which new items they might be most interested in, and which content might fit best with their interests. A variety of algorithms and methods exist to process the data and use it to make predictions. One such algorithm is xPAL, which is a decision-theoretic approach to measure the usefulness of a labeling candidate in terms of its expected performance gain. With xPAL and other active machine learning methods an optimal strategy can be explored to classify new data points.

Keywords: probabilistic active learning xPAL machine learning multi-class classification active learning

Contents

Introduction	2
1 Machine Learning Review	3
1.1 Definitions	3
1.2 Active Learning	3
1.2.1 Random Sampling	4
1.2.2 Diversity Sampling	4
1.2.3 Uncertainty Sampling	4
1.2.4 xPAL	4
1.2.5 EER	5
1.2.6 Summary	5
2 Related Works	6
2.1 Probabilistic Active Learning	6
2.2 Multiclass Probabalistic Active Learning	6
2.3 xPAL	6
3 Analysis	7
3.1 Problem Recap	7
3.2 Experiments	7
3.3 Results and Analysis	7
Conclusion	8
Bibliography	9
List of Figures	10
List of Tables	11
List of Abbreviations	12
A Attachments	13
A.1 First Attachment	13

Introduction

One of the main challenges of creating a successful machine learning model is obtaining labeled data. With easy access to a variety of modern tools, devices, and sensors, we are able to rapidly collect unlabeled data. But, in supervised learning, prediction models are trained using labeled data. The problem is that acquiring labels for the collected data can be expensive, time-consuming, or even impossible in some cases.

However, methods have been developed to help reduce the time required to label this data. Active learning is a semi-supervised machine learning method where the model is trained with a smaller set of labeled data but also aims to exploit trends within the unlabeled data. Active learning has been heavily researched in the past but typically with binary data. Multi-class active learning research is still in its infancy.

Active learning is different from other machine learning methods because it uses the unlabeled data and some evaluation criteria to determine which candidate could be the most beneficial to the model if it was given a label. In summary, the model requests the label from some oracle that provides the label then it takes this new labeled data point and rebuilds the model. We describe it as semi supervised active learning because of the oracle (typically a human) involved in the process that provides the label for the requested candidate data.

In our case we have some data (website urls) for some company or business that are given to us from our partner. From this data our partner currently utilizes human labor to browse the website and then label the url with a category (23 labels) and a sub-category (234+ tags) that branch from the main category but still have some relation. This is a repetitive and expensive task that could be automated using active learning.

To reduce the burden of human labeling we propose creating a workflow using Scrapy, Postgres, translation services, and semi supervised Active Learning (bayesian, expected error, etc.) that only require occasional interaction where a human can label a candidate that is most beneficial to the model such as using a decision-theoretic approach to measure the usefulness of a labeling candidate in terms of expected performance gain. The workflow takes the website as input, navigates to the webpage, collects and translates the text, and adds it to a database. We then run the model using the data from the database and return the model.

1. Machine Learning Review

1.1 Definitions

Russel and Norvig succinctly define an agent and differnt types of learning in their book "Artificial Intelligence: A Modern Approach" (Russell and Norvig [2009]), their definition is paraphrased here. They define an agent as something that acts and a rational agent as one that acts so as to achieve the best outcome. If there there is uncertainty, then the agent tires to achieve the best expected outcome. Any component of an agent can be improved by learning from data. The improvements and techniques used to make them depend on four major factors:

- Which component is to be improved.
- What prior knowledge the agent already has.
- What representation is used for the data and the component.
- What feedback is available to learn from.

Here we will mostly be focused on the final point, "What feedback is available to learn from" but also slightly on the second point because we will incorporate Bayesian learning. There are three main types of feedback that determine the three main types of learning which are unsupervised, reinforcement, and supervised learning.

In unsupervised learning an agent learns patterns even though no feedback is provided. In reinforcement learning, the agent learns from a series of rewards or punishments. In supervised learning, an agent learns from input-output pairs, which can be discrete or continuous, to find a function that maps the pairs.

The goal of supervised learning is given a training set of N example input-output pairs:

$$(x_1, y_1), (x_2, y_2), \dots (x_N, y_N),$$

where each y_j was generated by some unknown function $y = f(x)$, find a function h that approximates the true function (Russell and Norvig [2009]).

In reality, the lines separating the types of learning aren't so clear. Semi-supervised learning is also an important and widely used method. In semi-supervised learning we are given a few labeled examples that were labled by some orcale (labeler, data annotator, etc.) and we must then make the most of a large collection of unlabeled examples. But what can we do with all the unlabeled data? This is where active learning comes in.

1.2 Active Learning

Supervised learning models almost always get more accurate with more labeled data. Active learning is the process of deciding which data to sample for annotation (Munro [2021]). In other words, the central component of an active learning

algorithm is the selection strategy, or deciding which of the unlabeled data could be the most useful to the model if it was labeled. Active learning uses a selection strategy that augments the existing classifier, it is not itself a classifier but rather a tool paired with a classifier.

Many strategies for choosing the next points to label exist. Here we will briefly define three basic approaches: uncertainty, diversity, and random sampling to get an idea of sampling. As well as two more advanced sampling approaches: xPAL and EER. When sampling the unlabeled data an ordered list is returned and the top candidate is the candidate that is expected to be most valuable for the model, but we are not strictly limited to taking just one candidate.

1.2.1 Random Sampling

Random sampling is rather self explanatory as we randomly select an unlabeled data point from the pool and have an oracle provide a label.

1.2.2 Diversity Sampling

Diversity sampling is the set of strategies for identifying unlabeled items that are underrepresented or unknown to the machine learning model in its current state. The items may have features that are unique or obscure in the training data, or they might represent data that are currently under-represented in the model.

Either way this can result in poor or uneven performance when the model is applied or the data is changing over time. The goal of diversity sampling is to target new, unusual, or underrepresented items for annotation to give the algorithm a more complete picture of the problem space (Munro [2021]).

1.2.3 Uncertainty Sampling

Uncertainty sampling is the set of strategies for identifying unlabeled items that are near a decision boundary in your current machine learning model. If you have a binary classification task, these items will have close to a 50% probability of belonging to either label; therefore, the model is called uncertain or confused.

These items are most likely to be wrongly classified, so they are the most likely to result in a label that differs from the predicted label, moving the decision boundary after they have been added to the training data and the model has been retrained (Munro [2021]).

1.2.4 xPAL

Extended probabilistic gain for active learning (xPAL) is a decision-theoretic selection strategy that directly optimizes the gain and misclassification error, and uses a Bayesian approach by introducing a conjugate prior distribution to determine the class posterior to deal with uncertainties. Although the data distribution can be estimated, there is still uncertainty about the true class posterior probabilities.

These class posterior probabilities can be modeled as a random variable based on the current observations in the dataset. For this model, a Bayesian approach

is used by incorporating a conjugate prior to the observations. This produces more robust usefulness estimates for the candidates Kottke et al. [2021].

1.2.5 EER

Monte Carlo estimation of error reduction (EER) estimates future error rate by log-loss, using the entropy of the posterior class distribution on a sample fo the unlablled examples, or by 0-1 loss, using the posterior probabilities of the most probable class for the sampled unlabelled examples. Roy and McCallum [2001]

1.2.6 Summary

After we select and label the candidate data we retrain the model with the updated data and we continue this process as new data is obtained or until we are satisfied with the performance of the model

2. Related Works

Ideas listed here... is related works a good topic to discuss next?

2.1 Probabilistic Active Learning

2.2 Multiclass Probabalistic Active Learning

2.3 xPAL

3. Analysis

Outline the problem, the experiments, the results, and analysis.

3.1 Problem Recap

3.2 Experiments

3.3 Results and Analysis

Conclusion

Bibliography

Daniel Kottke, Marek Herde, Christoph Sandrock, Denis Huseljic, Georg Kreml, and Bernhard Sick. Toward optimal probabilistic active learning using a bayesian approach. *Machine Learning*, 110(6):1199–1231, May 2021. doi: 10.1007/s10994-021-05986-9. URL <https://doi.org/10.1007/s10994-021-05986-9>.

Robert Munro. *Human-in-the-loop machine learning*. Manning Publications, New York, NY, July 2021.

Nicholas Roy and Andrew McCallum. Toward optimal active learning through monte carlo estimation of error reduction. *ICML, Williamstown*, 2:441–448, 2001.

Stuart Russell and Peter Norvig. *Artificial Intelligence: A Modern Approach*. Prentice Hall, Upper Saddle River, NJ, 3rd edition, 2009. ISBN 9780136042594.

List of Figures

List of Tables

List of Abbreviations

A. Attachments

A.1 First Attachment