# Review of Three Papers on Active Learning

**Debarshi Kundu**
dqk5620@psu.edu

## 1 Introduction

Active learning is a specialized machine learning technique aimed at reducing the amount of labeled data required to train models effectively. Unlike traditional supervised learning, where models are trained on large datasets, active learning selectively queries the most informative data points for labeling, minimizing the need for extensive annotation efforts. This is particularly valuable in contexts such as medical diagnosis, natural language processing, and image recognition, where labeling data can be costly and time-consuming.

In a typical active learning process, a model is trained on a small labeled dataset and iteratively queries an oracle (such as a human annotator) to label additional data points that are likely to enhance the model's performance. These data points are chosen based on strategies like uncertainty sampling or query-by-committee. This report reviews three recent papers that contribute to the advancement of active learning methods:

- **"Active Learning Through a Covering Lens" (NeurIPS 2022)**Yehuda et al. [2022]
- **"Active Learning with Neural Networks: Insights from Nonparametric Statistics" (NeurIPS 2022)**Zhu and Nowak [2022]
- **"Empowering Active Learning to Jointly Optimize System and User Demands" (ACL 2020)**Lee et al. [2020]

## 2 "Active Learning Through a Covering Lens" (NeurIPS 2022)

### 2.1 Problem the Paper Tries to Solve (Motivation)

This paper tackles the challenge of low-budget active learning, where the number of labeled data points that can be queried is extremely limited. Existing methods tend to fail under such conditions, leading to suboptimal performance and sometimes performing no better than random sampling.

### 2.2 How It Solves the Problem

The paper introduces **ProbCover**, a novel algorithm that approaches active learning from a covering perspective. ProbCover focuses on maximizing Probability Coverage by selecting data points that cover the most densely populated areas of the data distribution. The algorithm employs a greedy approximation to solve the NP-hard problem of maximizing coverage, making it especially useful in low-budget and semi-supervised learning scenarios.

### 2.3 List of Novelties/Contributions

- **Introduction of ProbCover**: A novel active learning algorithm designed for low-budget settings, significantly outperforming other methods on various benchmarks.
- **Connection to High-Budget Approaches**: The paper shows that high-budget methods, such as Coreset, are ineffective in low-budget scenarios and introduces a duality between these regimes.

- **Efficiency in Semi-Supervised Learning**: ProbCover allows models to achieve state-of-the-art performance using fewer labeled examples.
- **Greedy Approximation**: The algorithm guarantees a near-optimal solution with a 1 1/e approximation ratio.

## 2.4 Downsides of the Work

- **Dependence on Embedding Spaces**: The success of ProbCover relies heavily on the existence of high-quality embedding spaces, which may not always be available or easy to generate.
- **Approximation Limits**: The greedy approximation, while efficient, does not achieve the optimal solution.
- **Dataset Sensitivity**: Performance may vary significantly depending on the specific dataset and its structure.

# 3 "Active Learning with Neural Networks: Insights from Nonparametric Statistics" (NeurIPS 2022)

## 3.1 Problem the Paper Tries to Solve (Motivation)

This paper addresses the theoretical gap in deep active learning, which has shown promising empirical results but lacks rigorous label complexity guarantees. The authors aim to provide the first near-optimal label complexity guarantees for deep active learning from a nonparametric classification perspective.

## 3.2 How It Solves the Problem

The paper leverages nonparametric statistical tools to provide label complexity guarantees for deep active learning. It proves that deep learning models can achieve near-minimax label complexity under standard low-noise conditions. Furthermore, the paper introduces an oracle-efficient algorithm that can achieve polylogarithmic label complexity when an abstention option is allowed, even without low-noise assumptions.

## 3.3 List of Novelties/Contributions

- **First near-optimal label complexity guarantees**: This paper is the first to provide such guarantees in the context of deep active learning.
- **Active learning with abstention**: The introduction of an abstention option leads to exponential reductions in label complexity.
- **Extension to Radon BV2 spaces**: The paper extends its theoretical framework to these spaces, which are more naturally aligned with neural networks than traditional Sobolev spaces.
- **Generalization to broader function spaces**: This makes the framework applicable beyond typical Sobolev or Hölder spaces.

## 3.4 Downsides of the Work

- **Dependence on noise and smoothness parameters**: Some theoretical guarantees assume the availability of noise and smoothness parameters, which may not always be practical.
- **Limited analysis of the disagreement coefficient**: A more thorough examination of the role of the disagreement coefficient could improve the theoretical understanding of deep active learning.
- **Abstention strategy limitations**: Relying on the abstention option might not be suitable for all real-world applications, where making decisions, even with uncertainty, is required.

# 4 "Empowering Active Learning to Jointly Optimize System and User Demands" (ACL 2020)

## 4.1 Problem the Paper Tries to Solve (Motivation)

This paper addresses the challenge of balancing system efficiency with user satisfaction in active learning settings. Traditional active learning optimizes model performance but may overlook the user's needs, which can be particularly problematic in personalized education platforms.

## 4.2 How It Solves the Problem

The authors propose a joint optimization framework that considers both system uncertainty and user satisfaction. They introduce two sampling strategies:

- **Combined Sampling**: Selects data points that are both uncertain for the system and useful for the user.
- **Trade-off Sampling**: Balances system performance and user satisfaction using a weighted function.

The framework is tested in an educational context, optimizing the selection of exercises for users based on their proficiency levels.

## 4.3 List of Novelties/Contributions

- **Joint optimization framework**: This is one of the first approaches to simultaneously optimize system performance and user satisfaction.
- **Novel sampling strategies**: These strategies ensure that users receive personalized tasks while also improving the model's performance.
- **Real-world educational application**: The proposed methods were validated on a language-learning platform, enhancing user engagement and system efficiency.
- **Experimental validation**: The results demonstrated superior performance compared to random and traditional active learning methods.

## 4.4 Downsides of the Work

- **Implementation complexity**: Jointly optimizing two objectives increases the complexity of the system, making it harder to implement and scale.
- **Limited generalization**: While the approach is promising in educational settings, its applicability to other domains, such as recommendation systems, is not fully explored.
- **Data dependency**: The method's success hinges on the availability of reliable data for both system uncertainty and user satisfaction.

# 5 Conclusion

These three papers each contribute unique advancements to the field of active learning. **"Active Learning Through a Covering Lens"** offers a practical solution for low-budget active learning, while **"Active Learning with Neural Networks"** provides much-needed theoretical guarantees. Finally, **"Empowering Active Learning"** highlights the importance of aligning model objectives with user needs. Together, they represent significant strides toward making active learning more efficient, robust, and user-centric.

# References

Ji-Ung Lee, Christian M Meyer, and Iryna Gurevych. Empowering active learning to jointly optimize system and user demands. *arXiv preprint arXiv:2005.04470*, 2020.

Ofer Yehuda, Avihu Dekel, Guy Hacohen, and Daphna Weinshall. Active learning through a covering lens. *Advances in Neural Information Processing Systems*, 35:22354–22367, 2022.

Yinglun Zhu and Robert Nowak. Active learning with neural networks: Insights from nonparametric statistics. *Advances in Neural Information Processing Systems*, 35:142–155, 2022.