

# **Homework1 - Paper Review**

Guodong Dong

[gvd5289@psu.edu](mailto:gvd5289@psu.edu)

Pennsylvania State University

College of Information Sciences and Technology

CSE 584: Machine Learning: Tools and Algorithms

Prof. Wenpeng Yin

September 15, 2024

# Content

<b>Paper 1: Active Learning Using Pre-clustering</b> .....	1
<b>Motivation</b> .....	1
<b>Solutions</b> .....	1
<b>novelties/contributions</b> .....	2
<b>Downsides</b> .....	3
<b>Paper 2: Active learning with support vector machines</b> .....	4
<b>Motivation</b> .....	4
<b>Solutions</b> .....	4
<b>novelties/contributions</b> .....	5
<b>Downsides</b> .....	6
<b>Paper 3: Active learning for logistic regression an evaluation</b> .....	7
<b>Motivation</b> .....	7
<b>Solutions</b> .....	7
<b>novelties/contributions</b> .....	8
<b>Downsides</b> .....	8
<b>References</b> .....	10

# **Paper 1: Active Learning Using Pre-clustering**

## **Motivation**

Nguyen and Smeulders (2004) discuss the challenge of selecting the most valuable data for labeling in two-class active learning. Labeling is time-consuming, and standard methods often focus only on samples close to the classification boundary. They try to improve the quality of active learning by considering the distribution of previous data instead of depending just on the distance from the categorization boundary. This makes the labeling process more efficient by picking representative samples and eliminating redundant labeling within the same cluster.

## **Solutions**

This paper explores the problem of improving the efficiency of active learning by incorporating clustering into the data selection process. Traditional methods focus on selecting samples near the classification boundary, where the model is uncertain, but these approaches often ignore the underlying data distribution and can result in redundant labeling. To overcome this, Nguyen and Smeulders propose a framework that clusters the data and selects representative samples from each cluster for labeling. This ensures that the labeled set is diverse and avoids repeatedly labeling similar samples within the

same cluster. A classifier is then trained on the cluster representatives, and the classification decision is propagated to other samples in the dataset through a local noise model.

The clustering process adapts over time using a coarse-to-fine strategy, where large clusters are used initially and then refined into smaller ones as the learning progresses. This allows the model to focus on broader data distinctions early on and gradually improve classification precision. The selection of samples for labeling considers both the uncertainty of the sample and its representativeness within the data. This approach ensures that the most informative and diverse samples are labeled, leading to more efficient learning and better classifier performance compared to traditional boundary-based methods.

## **novelties/contributions**

One of the primary innovations in this paper is the integration of clustering into the active learning process. This framework uses clustering to identify representative samples from different parts of the data space. By selecting central samples from each cluster, the algorithm ensures diversity in the labeled set and reduces the chances of redundant labeling within the same cluster. Another contribution is the use of a local noise model to propagate classification decisions from labeled cluster representatives to other unlabeled samples within the same cluster. This model accounts for potential inaccuracies near cluster boundaries by assigning soft memberships to samples. This

ensures that the classifier can generalize better, even in areas where clusters overlap or are less well-defined.

The algorithm combines two criteria for selecting samples for labeling: uncertainty (how close a sample is to the classification boundary) and representativeness (whether the sample is a central point in a cluster). This dual focus ensures that the most informative samples are chosen while maintaining a broad and diverse labeled set, which improves classifier performance with fewer labeled samples.

## **Downsides**

The clustering process is computationally expensive, especially when applied to large datasets. While Nguyen and Smeulders attempt to mitigate this by using a coarse-to-fine clustering strategy and breaking the data into smaller subsets, the overall computational cost remains high. This can limit the scalability of the approach, especially for very large datasets. Moreover, the method is restricted to using logistic regression as the classifier, which may not perform as well in more complex or non-linear datasets. As a result, the approach might not generalize well to more challenging classification tasks that require non-linear decision boundaries.

## **Paper 2: Active learning with support vector machines**

### **Motivation**

Kremer et, al. (2014) aim at reducing the labeling costs in active learning, where large amounts of unlabeled data are available, but obtaining labels is expensive. They explore how active learning strategies can be used with Support Vector Machines (SVMs) to autonomously select the data points that, if labeled, would most enhance the model's accuracy.

### **Solutions**

Kremer et, al. address the problem of reducing labeling costs in machine learning by integrating active learning with SVMs. SVMs are particularly suitable for this because they provide a clear mathematical framework for identifying which samples are most informative. They leverage the concept of uncertainty sampling, where the algorithm selects samples that are closest to the SVM's decision boundary. These samples represent the points where the classifier is least confident, so labeling them would provide the most information to refine the model. By focusing on these uncertain samples, the algorithm can maximize the improvement of the model with fewer labeled data points, reducing the overall labeling cost.

In addition to uncertainty sampling, Kremer et, al. introduce strategies to ensure that the selected samples are not only uncertain but also representative of the broader data distribution. This is important to avoid overfitting to specific regions near the decision boundary. They propose a Max-Min Margin approach, which evaluates the impact of both possible labels on the decision boundary to ensure the most uncertain regions are addressed. Furthermore, it incorporates methods like batch-mode active learning and online learning, where models are updated incrementally as new labels are acquired, to ensure efficient computation even when working with large datasets. These approaches help to strike a balance between improving model performance and reducing the labeling and computational costs.

## **novelties/contributions**

This work provides an in-depth review of uncertainty sampling with SVMs. In active learning, samples near the decision boundary are the most informative. SVMs are particularly well-suited for this because the distance of a data point from the decision boundary can be easily calculated in the SVM framework. Kremer et, al. formalize this approach and show that the most uncertain samples which are closest to the decision boundary can significantly improve model accuracy with fewer labeled samples.

Another key novelty is the combination of informativeness and representativeness in the sample selection process. While uncertainty sampling focuses on samples near the decision boundary, this strategy could

lead to over-sampling outliers of the data that are not representative of the overall distribution. To counter this, Kremer et, al. select samples that are both uncertain and representative of the underlying data distribution, ensuring that the labeled set reflects a broader spectrum of the data. This method prevents the classifier from overfitting to specific regions and improves generalization.

## **Downsides**

Active learning may introduce a selection bias, as samples are not randomly selected but chosen based on specific criteria like uncertainty or representativeness. This selection bias can negatively impact the generalization performance of the model. Kremer et, al. discuss importance weighting as a way to address this bias, but the estimation of accurate weights is challenging. If the weights are not correctly determined, the correction for the selection bias may fail, leading to a model that performs well on the actively selected samples but poorly on unseen data. Uncertainty-based sampling can prioritize outliers or noisy data near the decision boundary, which may negatively affect model performance. While Kremer et, al. propose to be diversity in selected samples, there is still a risk of redundancy, particularly in densely populated data regions. This can waste labeling resources without significant model improvement.



## **Paper 3: Active learning for logistic regression an evaluation**

### **Motivation**

Schein and Ungar (2007) focus on the difficulty of identifying the most efficient active learning methods tailored specifically for logistic regression, particularly in pool-based scenarios where a large amount of unlabeled data is available, but the labeling process is both costly and time-consuming. In such settings, it is crucial to minimize the number of labeled data points required to build a reliable model. They try to enhance the training efficiency of logistic regression models by selecting data points that provide the most information about the underlying patterns in the data. The goal is to significantly reduce the overall cost of labeling while ensuring the model achieves high generalization accuracy. They explore and evaluate various active learning strategies to achieve this balance between cost and performance.

### **Solutions**

Schein and Ungar discuss several active learning methods, focusing on both theoretically grounded approaches like A-optimality (variance reduction) and heuristic methods such as uncertainty sampling and query by committee. For experimental design-based methods like A-optimality, the goal is to reduce

model prediction variance. In contrast, heuristic methods like margin sampling and query by bagging (QBB) prioritize samples based on classifier uncertainty or disagreement among committee members. Schein and Ungar run evaluations across various datasets to compare these methods, analyzing their effectiveness in reducing labeling costs while improving the performance of logistic regression classifiers.

### **novelties/contributions**

The key contributions include the application of A-optimality for active learning in logistic regression, adapting a variance reduction technique to minimize prediction uncertainty and improve model performance with fewer labeled data. Additionally, Schein and Ungar provide a comprehensive evaluation of various heuristic methods, such as uncertainty sampling, query by committee, and margin sampling, comparing their performance and computational efficiency. While A-optimality is computationally intensive, the heuristic methods offer faster alternatives, though some, like entropy-based sampling, perform poorly on noisy data. The detailed comparison highlights when and why certain methods are more effective, offering insights into optimizing labeling strategies in active learning for logistic regression.

### **Downsides**

While this approach provides strong theoretical benefits by minimizing the variance in the model's predictions, its practical application is limited due to

the large computational overhead required to evaluate the variance for each unlabeled data point. As the dataset size increases, the cost of computing the Fisher information matrix and performing matrix inversions becomes prohibitively expensive, making A-optimality less feasible for large-scale or real-time applications. This restricts its use to smaller datasets or cases where computational resources are not a constraint, limiting its overall applicability in real-world scenarios. Another major weakness is the inconsistent performance of heuristic methods, particularly in noisy datasets. While methods like uncertainty sampling and query by committee are computationally more efficient than A-optimality, their performance can be unreliable. This inconsistency makes it difficult to apply these heuristic methods universally across different types of datasets.

## References

- Nguyen, H. T., & Smeulders, A. (2004, July). Active learning using pre-clustering. In Proceedings of the twenty-first international conference on Machine learning (p. 79).
- Kremer, J., Steenstrup Pedersen, K., & Igel, C. (2014). Active learning with support vector machines. *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, 4(4), 313-326.
- Schein, A. I., & Ungar, L. H. (2007). Active learning for logistic regression: an evaluation. *Machine Learning*, 68, 235-265.