

# CSE 584 Homework1: Active Learning Paper Reviews

Yu Tianrun

September 16, 2024

## **Paper 1: Support Vector Machine Active Learning with Applications to Text Classification**

### **1. What problem is this paper trying to solve, i.e., its motivation?**

This paper aims to address the problem of high labeling costs in supervised learning tasks, particularly for support vector machines (SVMs). Most machine learning algorithms, including SVMs, rely on randomly selected, pre-labeled training sets, which can be costly and time-consuming to create. The paper proposes using pool-based active learning to minimize the number of labeled training instances required by allowing the learner to select which data points to query for labels, thus reducing the overall labeling burden.

### **2. How does it solve the problem?**

The paper introduces a new algorithm for active learning with SVMs. Instead of relying on random sampling for labeling, the algorithm leverages pool-based active learning, where the model selects the most informative data points from a pool of unlabeled instances to query for labels. The authors provide a theoretical motivation using the concept of a "version space" and design strategies for selecting queries that maximize learning efficiency. They conduct experiments to demonstrate that this method can significantly reduce the number of labeled instances needed in both inductive and transductive learning settings.

### **3. List of novelties/contributions:**

- The paper presents a novel active learning algorithm specifically designed for support vector machines (SVMs), enabling more efficient selection of instances for labeling.
- It introduces the notion of using version space as a theoretical foundation for query selection in active learning.
- The authors demonstrate that their proposed active learning approach can significantly reduce the need for labeled data in both inductive and transductive tasks.
- The algorithm is empirically validated through experiments on text classification, showing its effectiveness in real-world applications.

#### **4. What do you think are the limitations of this work?**

- The abstract does not specify how scalable the approach is for very large datasets, especially considering the computational cost of SVMs and version space estimation in high-dimensional feature spaces.
- While the method is proven effective for text classification, the paper may not thoroughly explore its applicability to other domains or more complex, noisy datasets.
- The approach is specific to SVMs, so its applicability to other types of models or tasks outside the scope of SVM-based learning could be limited.

## **Paper 2: Deep Bayesian Active Learning with Image Data**

#### **1. What problem is this paper trying to solve, i.e., its motivation?**

This paper addresses the challenges of applying active learning in deep learning frameworks, specifically for high-dimensional image data. Deep learning models typically require large amounts of labeled data, which makes labeling an expensive and time-consuming process. Additionally, many active learning techniques rely on model uncertainty, which is not well represented in traditional deep learning models. The motivation of the paper is to create an active learning framework that works effectively with deep learning, particularly for image data, reducing the number of labeled instances required.

#### **2. How does it solve the problem?**

The paper combines recent advances in Bayesian deep learning with active learning. The authors propose using Bayesian convolutional neural networks (BCNNs) to model uncertainty, which is crucial for effective active learning. The BCNNs allow the system to identify which data points would be most beneficial to label, based on the uncertainty of the model. By applying this framework to image data, including the MNIST dataset and a real-world skin cancer diagnosis task (ISIC2016), the authors demonstrate that their approach significantly reduces the need for labeled data while maintaining high accuracy.

#### **3. List of novelties/contributions:**

- The paper integrates Bayesian deep learning into active learning, specifically using Bayesian convolutional neural networks (BCNNs) to address model uncertainty.
- It introduces an active learning framework that scales to high-dimensional data, an area where existing active learning methods struggle.
- Experimental results show that the proposed method achieves better performance with significantly fewer labeled instances compared to traditional random sampling and other active learning methods.
- The method is empirically validated on both synthetic datasets (MNIST) and real-world medical datasets (ISIC2016 for skin cancer diagnosis).

#### **4. What do you think are the limitations of this work?**

- While the paper demonstrates success with BCNNs, Bayesian deep learning approaches, including BCNNs, are computationally expensive and may not be practical for very large datasets or models with numerous parameters.
- The experiments are limited to image data, particularly MNIST and ISIC2016, which are specific in scope. The generalization of this method to other types of high-dimensional data or more complex real-world tasks is not explored in depth.
- The approach assumes that the uncertainty estimation provided by BCNNs is reliable, but the quality of uncertainty estimation in practice may vary depending on the data and model architecture.
- As with many active learning approaches, the method requires repeated retraining of the model as more data is labeled, which could become computationally costly in large-scale applications.

### **Paper 3: Active Learning for Convolutional Neural Networks: A Core-Set Approach**

#### **1. What problem is this paper trying to solve, i.e., its motivation?**

This paper aims to address the challenge of efficiently training convolutional neural networks (CNNs) with a limited labeling budget. CNNs require a large amount of labeled data to achieve high performance, but collecting and labeling large datasets is expensive and time-consuming. Traditional active learning methods, which select a subset of data points for labeling, are often ineffective when applied to CNNs in a batch setting. The paper’s motivation is to create an active learning framework tailored to CNNs, capable of selecting the most informative data points to label, reducing the labeling burden while maintaining high accuracy.

#### **2. How does it solve the problem?**

The paper proposes redefining the active learning problem for CNNs as a “core-set selection” problem. Core-set selection aims to choose a small subset of unlabeled data such that the model trained on this subset performs comparably to one trained on the entire dataset. To solve this, the authors derive a theoretical bound between the loss of the selected subset and the remaining data points, using the geometry of the data. This problem is reduced to the k-center problem, and an efficient approximate solution is applied. Their experiments show that this core-set selection method significantly improves performance in image classification tasks.

#### **3. List of novelties/contributions:**

- The paper redefines active learning for CNNs as a core-set selection problem, offering a novel approach for selecting subsets of data in a batch setting.

- It provides a theoretical analysis of the relationship between a selected subset’s loss and the geometry of the remaining dataset, introducing a new bound for core-set selection.
- The proposed method applies the k-center problem to active learning, presenting an efficient algorithm to solve it.
- Empirical results demonstrate that the method significantly outperforms existing active learning techniques on image classification tasks using CNNs.
- The method is applicable in both fully supervised and weakly supervised learning scenarios, offering flexibility across different learning settings.

#### **4. What do you think are the limitations of this work?**

- The core-set selection method, while theoretically sound, may still face scalability issues when applied to very large datasets due to the computational cost associated with solving the k-center problem.
- The approach is tested primarily on image classification tasks, so its generalizability to other types of data (such as text or audio) or more complex real-world applications is not thoroughly explored.
- The method does not incorporate model uncertainty into its selection process, which might limit its effectiveness in cases where uncertainty-based acquisition strategies are more appropriate.
- Although the proposed method is more efficient for batch-mode active learning, the paper does not deeply analyze the computational cost of iterative model retraining, which could be significant in large-scale applications.