

Active Learning Methods for Annotating Training Sets

Akash R, Amit M Madiwalar, Bhoomica Basavaraju, G Tharun Kumar

Department of Computer Science and Engineering, RV Institute of Technology and Management, Bengaluru, India

Abstract: Active learning, a machine learning approach, identifies data requiring human annotations, thereby reducing the cost and time of data collection while maintaining high accuracy. This method involves training machine learning models on a small set of labeled data and then leveraging the model to predict labels for unlabeled objects. Selection of data points where the model is most uncertain for annotation iteratively refines the model until desired results are achieved. Active learning has proven beneficial across various machine learning tasks, including text classification, image classification, entity recognition, and natural language processing, particularly in scenarios where annotation is resource-intensive. In this study, we investigate active learning's application on CIFAR10, EuroSAT, and Fashion MNIST datasets, comparing different active learning methods such as minimum confidence, probability models, and entropy models. Our findings illustrate that both approaches enhance model performance compared to random sampling, underscoring the efficacy of active learning in improving image classification tasks across diverse datasets.

Keywords: Active learning , Human Labeling, Least Confidence, Margin Sampling, Entropy Sampling, CIFAR-10 , EuroSAT, CNN, Fashion MNIST.

I. INTRODUCTION

In machine learning, data labeling is the process of assigning labels to data points, which significantly influences a model's performance. However, data collection can be labor-intensive and costly, especially with large datasets. Active learning addresses this challenge by reducing labeling requirements, where a model learns from a minimal set of labeled data. Experts identify the model's minimum evidence, repeating this process until the model is sufficiently trained, which can improve performance by training on a more representative data subset. Minimum confidence involves selecting data points for which the model has the lowest confidence, assuming that uncertain data points are more informative. Marginal sampling chooses the most uncertain data by computing the margin for each data point, representing the difference between the highest and second-highest probabilities. Entropy sampling selects data points with high entropy, indicating greater uncertainty and potential usefulness. These techniques are applied to the CIFAR-10, EuroSAT, and FashionMNIST datasets.

II. LITERATURE SURVEY

Active Learning with Bayesian UNet for Efficient Semantic Image Segmentation. [1] The principal aim of their research was to present a sample-efficient methodology for image segmentation, known as AB-UNet (Active Bayesian UNet). This entails a convolutional neural network that incorporates batch normalization and max-pool dropout. The outcomes of their investigation exhibited notable enhancements through the utilization of AB-UNet across all datasets. Their methodology is straightforward to deploy, manageable, and facilitates quicker generalization in contrast to alternative methods such as graphical-based segmentation. An area requiring further exploration is the absence of investigation into the Game theoretic strategy for leveraging predictions derived from AB-UNet.

Active learning with point supervision for cost-effective panicle detection in cereal crops. [2] The primary aim was to suggest a point supervision-based active learning approach for detecting panicles in cereal crops. Within their methodology, the system engages in continual interaction with a human annotator, iteratively soliciting labels solely for the most informative images instead of the entire dataset. Consequently, they have introduced a cost-efficient technique for training dependable panicle detectors in cereal crops. A cost-effective panicle detection approach for cereal crops is exceedingly advantageous for both breeders and agronomists. Plant breeders possess the capacity to swiftly obtain crop yield estimations, empowering them to make crucial decisions regarding crop management. The deficiency lies in the fact that the Active learning framework has not been expanded to include Regression bounding box size information for uncertainty estimation.

Gradient and Log-based Active Learning for Semantic Segmentation of Crop and Weed for Agricultural Robots. [3] Their main objective was to introduce and compare active learning strategies that intelligently pick images taken under new conditions to re-train an existing network: first by picking samples based on a log-space ranking of their loss with respect to pseudo labels. and second by selecting training samples that are expected to have a maximum effect on the network weights. The effectiveness of their method produces higher semantic segmentation accuracies with a small number of training samples, compared to random sampling as well as entropy based sampling. As a result of that, the effort in human annotation is reduced without compromising performance.

Active learning for object detection in high resolution satellite images. [4]The primary aim was to employ two active learning methodologies in the segmentation of satellite images: Bayesian dropout and Core-Set. In order to assess the efficacy of the uncertainty and core-set strategies, they are juxtaposed with two distinct baseline methods. It appears that the Core-Set technique proves to be the most efficient for the less robust model. Conversely, for the more robust model, although the Core-Set approach shows promising results, the Bayesian dropout technique outperforms it significantly. DIAL: Deep Interactive and Active Learning for Semantic Segmentation in Remote Sensing, [5] Their primary goal was to conduct Semantic segmentation: The neural network has the capability to autonomously produce accurate semantic segmentation maps, without relying on external instructions. Interactive Learning: Additionally, the neural network can enhance these segmentation maps by incorporating annotations to effectively rectify any errors. Active Learning: This involves assessing the neural network's uncertainty in order to direct the user towards specific queries. Consequently, they have developed a framework aimed at iteratively improving segmentation maps initially proposed by a neural network.

III METHODOLOGY

The CIFAR-10 dataset consists of 60,000 color images, each sized 32x32 pixels and categorized into 10 classes such as airplanes, cars, birds, cats, deer, horses, ships, and trucks. Among these, 50,000 images are allocated for training and 10,000 for testing, serving as a prominent benchmark for evaluating image classification algorithms

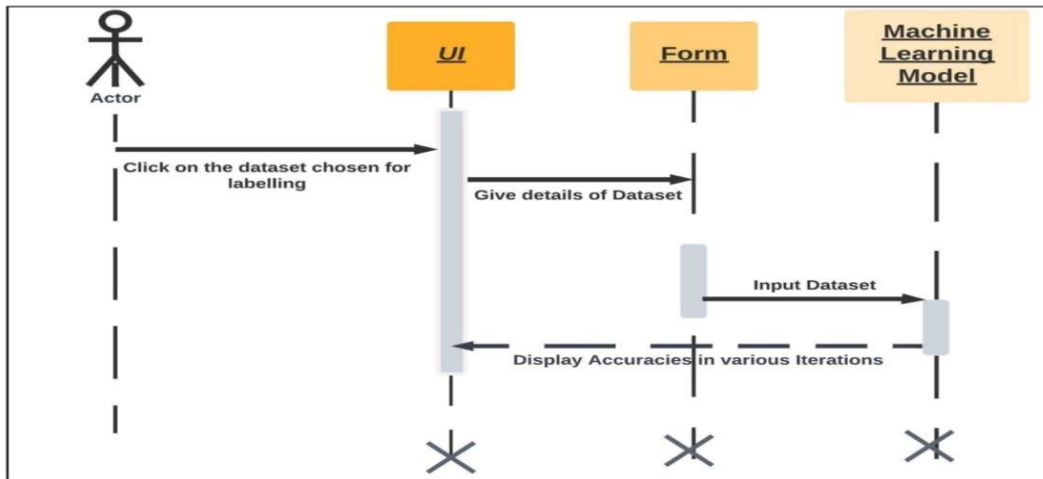
The EuroSAT RGB dataset is a publicly accessible repository containing 27,000 labeled and georeferenced satellite images representing 10 land use and land cover categories. These categories include annual crops, forests, herbaceous plants, roads, industries, pastures, permanent crops, residences, rivers, and oceans and ponds, each comprising 2,700 images. Each image has a resolution of 64x64 pixels and is stored in RGB format.

Fashion-MNIST dataset comprises 60,000 grayscale images, with dimensions of 28x28 pixels, representing 10 different fashion categories. Additionally, there is a test set consisting of 10,000 images. Each image is labeled with one of 10 categories: t-shirt/top, pants, sweater, dress, coat, sandals, jacket, sneakers, bag, and heels.

In the depicted architecture, a two-dimensional Convolutional Neural Network (CNN) is employed, utilizing the Rectified Linear Unit (ReLU) activation function and 3x3 kernel size. The model begins with an input of 32x32x3 image dimensions. Following three consecutive two-dimensional convolutional layers, a max-pooling layer with a pool size of 2x2 is applied. This layer selects the most dominant pixels within each block, reducing the overall matrix size while preserving information integrity. Subsequently, a Dropout layer with a dropout rate of 25% is introduced to mitigate overfitting by randomly deactivating connections of hidden units during training. Lastly, a flatten layer is utilized to convert the output from the convolutional layers into a single elongated feature vector. Dense layer is provided after the Dropout layer as each neuron in a dense layer is connected to every neuron in the previous layer and helps to learn complex relationships between the input and output data. Again a Dropout layer of 40% is provided to avoid further overfitting. Finally a Dense layer is provided with softmax as activation function.

System Visualization

The UML diagram shown in Figure 1 shows a sequence diagram where the user is the person who provides the data set to the active learning algorithm. Active Learning Algorithm is a computer program that asks the user for the labels of a subset of a data set. The user then labels the desired data points. The active learning algorithm updates its model uses the labeled data points and repeats this process until the model is sufficiently accurate. The user can then use the model to predict new data points.

**Fig.1.0: System Visualization****DATASET DESCRIPTION**

The CIFAR-10 dataset comprises 60,000 color images sized 32x32 pixels, categorized into 10 classes including airplane, automobile, bird, cat, deer, dog, frog, horse, ship, and truck. These images are divided into 50,000 training images and 10,000 test images, serving as a standard benchmark for image classification in machine learning algorithms. The dataset was sourced from the internet and resized for uniformity.

EuroSAT RGB dataset contains 27,000 geo-referenced satellite images, categorized into 10 land use and land cover classes. These categories include Annual Crop, Forest, Herbaceous Vegetation, Highway, Industrial, Pasture, Permanent Crop, Residential, River, and Sea and Lake, each comprising 2,700 images. The dataset, derived from Sentinel-2 satellite images, features a spatial resolution of 64x64 pixels and is stored in RGB format.

Fashion-MNIST dataset consists of 60,000 grayscale images sized 28x28 pixels, representing 10 fashion categories. Additionally, it includes a test set of 10,000 images. The categories span T-shirt/top, Trouser, Pullover, Dress, Coat, Sandal, Shirt, Sneaker, Bag, and Ankle boot. Designed as a substitute for the MNIST dataset, it facilitates benchmarking of machine learning algorithms with identical image size and splits for training and testing.

Figure 2 shows the structure of active learning. The model is based on labeling an unlabeled pool (U) which is first trained on a subset of the manually labeled dataset. Predictions are assigned as points to each associated image label. Based on these results, predictions are prioritized using a query strategy or method of minimum confidence, marginal sampling and entropic sampling. This is done to select samples from the subset most likely mislabeled. The samples selected are correct were reported against the database and then added directly to the Deep training set Learning model. The training set is used to update the unlabeled pool and deep train learning model.

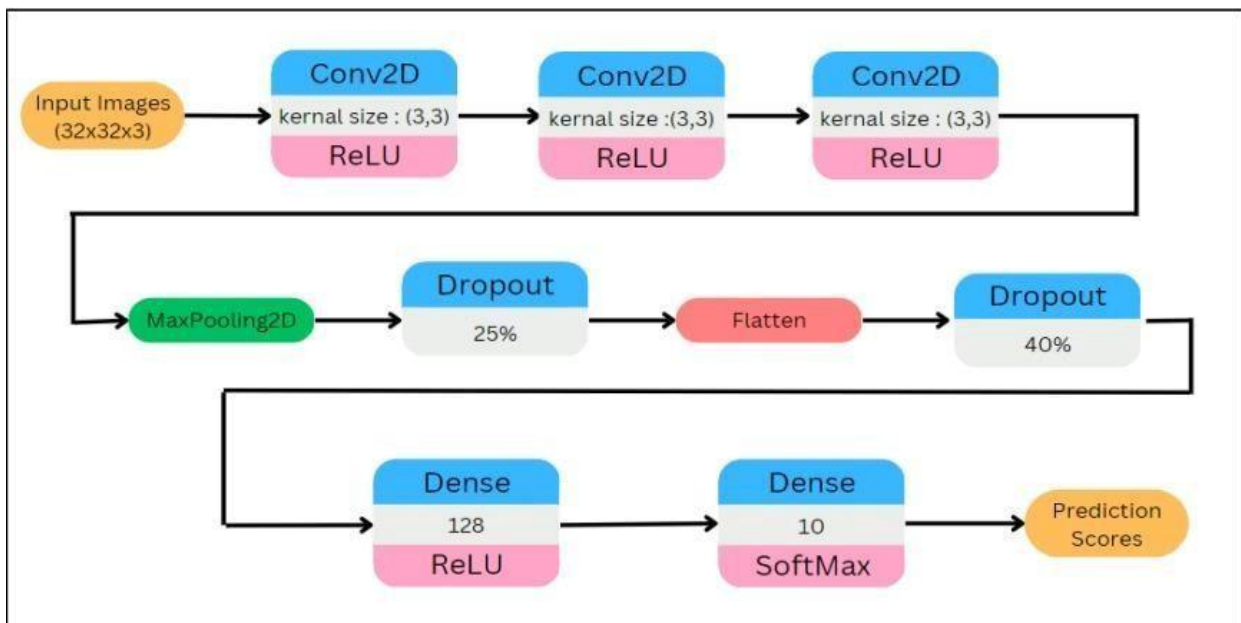


Fig. 2.0: Classifier Architecture: The various layers of the CNN model

Convolutional Neural Networks (CNN) are commonly employed in binary classifiers to discern the hyperplane, which serves as the boundary between two classes, maximizing the distinction between them. During the determination of this hyperplane, the support vectors, which are the structures nearest to the edges, are identified and utilized.

Region-based active learning offers a promising approach to minimize the labeling data needed for training semantic segmentation models, leveraging Bayesian principles and combining entropy and uncertainty metrics to identify data regions. This methodology has demonstrated efficacy across various scenarios and is anticipated to further enhance performance with the increasing prevalence of deep learning models. The AL-SSL scheme proposed for training aims to diminish the requisite training data volume, offering straightforward visualization tools compatible with diverse machine learning techniques. Additionally, the framework is scalable and adaptable for utilization with extensive datasets.

DIAL: Interactive deep learning for remote sensing semantic segmentation [8]: DIAL concept This model is a good method for remote sensing semantic segmentation. The framework is easy to use and can be used with many deep neural networks.

Active Learning for Timely Monitoring of Ear Detection in Economic Grain Crops [9]: Advocating for practical ear detection methods, this approach offers simplicity and compatibility with various deep learning models, making it suitable for implementation on large datasets. Effective for plant detection by agricultural robots, it requires a semantic train segmentation model, offering flexibility across multiple devices. Additionally, it provides a straightforward method for land cover classification in satellite images, adaptable to various deep learning models and scalable for diverse devices. [12] The study focuses on satellite image segmentation, employing two active learning methods: Bayesian output and

The active learning technique introduced in [11] shows promise in minimizing the labeled data necessary for training semantic segmentation models in satellite image land cover classification. Its implementation is straightforward, compatible with various deep learning architectures, and capable of handling large datasets efficiently.

Active learning for object detection in high-resolution satellite images.[12] Their main objective was to apply two active learning techniques to segmentation of satellite images: Bayesian dropout and Core-Set Measurement of Uncertainty and Contribution of Process Factors, we compare them to two different baselines. The Core-Set approach seems to be the most effective method for the weak model. For the strong model, even if the Core-Set approach is also performing well the Bayesian dropout approach outperforms it by a large margin.

Dropout as Bayesian approximation: Deep learning representation of model uncertainty.[16] This influential paper demonstrates how dropout in deep neural networks can be interpreted as approximate Bayesian inference, enabling active learning methods to leverage model uncertainty estimates for effective instance selection.

Active learning based on diversity with applications to the detection of unusual classes in data streams. [23] This work investigates rare-class recognition in streaming data using diversity-based active learning techniques. In order to enhance the detection of unusual occurrences in changing data streams, it provides a novel active learning technique that maximizes the diversity of labeled cases.

A comparison study using active learning for graph classification. [20] A comparison of active learning methods for graph classification tasks is presented in this research report. It assesses and contrasts several active learning techniques on diverse graph datasets, offering insights into active learning's efficacy for graph-based tasks.

The power of ensembles in classifying images through active learning.[18] In this study, deep ensemble models and active learning are combined. It demonstrates how ensemble-based active learning techniques may be used to do image classification tasks with less annotation work and higher classification accuracy.

This research explores active learning through a novel approach utilizing deep reinforcement learning techniques. The methodology involves framing active learning as a sequential decision-making problem and employing a deep Q-network to learn the agent's policy for selecting instances to label.

The study presented in [22] investigates active learning under conditions where annotators can only provide incomplete or restricted information. It introduces novel active learning techniques designed to handle missing or ambiguous label queries, thereby enhancing the resilience and flexibility of active learning in real-world annotation scenarios.

Dataset	Accuracy in training all images	Least Confidence	Margin Sampling	Entropy Sampling
CIFAR-10	70.49%	71.5%	71.4%	70.3%
EuroSAT RGB	85.94%	83.8%	86.09%	83.88%
Fashion MNIST	86.11%	87.2%	86.9%	86.6%

Table 1.0: Accuracy for each dataset and each methods

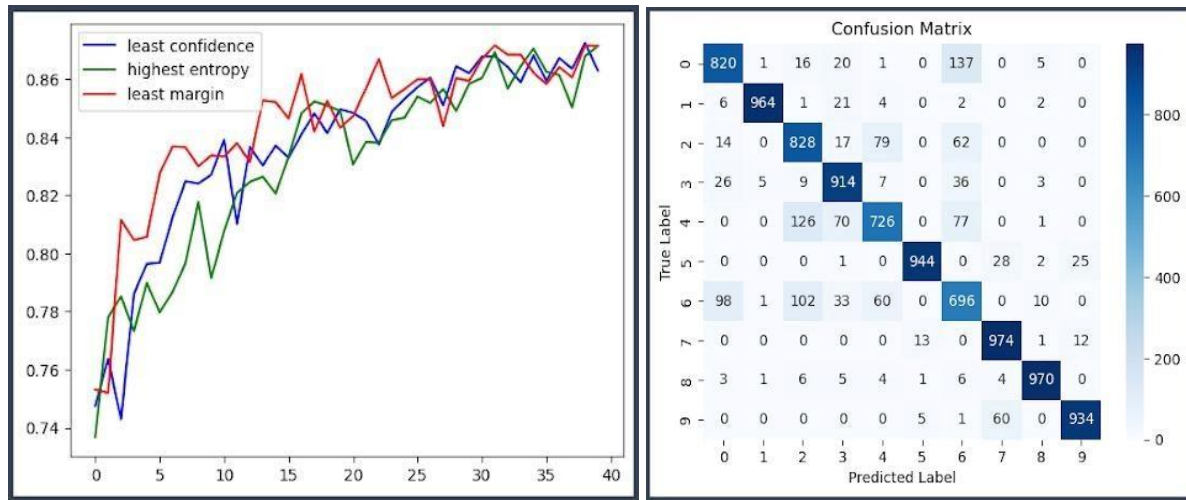


Fig.3.0:Least Confidence Method,Margin sampling Method and Entropy sampling Method Accuracy Plot and confusion matrix

In the CIFAR dataset, the accuracy achieved by the Least Confidence method is 71.5%, Margin Sampling is 71.4%, and Entropy Sampling is 70.3%. For the EuroSAT dataset, the accuracy achieved by the Least Confidence method is 83.8%, Margin Sampling is 86.09%, and Entropy Sampling is 83.88%. In the Fashion MNIST dataset, the accuracy achieved by the Least Confidence method is 87.2%, Margin Sampling is 86.9%, and Entropy Sampling is 86.6%.

V. CONCLUSION AND FUTURE WORK

Active learning is a powerful technique that can be used to decrease the amount of data labeling required. Active learning can be a useful tool for machine learning researchers looking to improve model performance while reducing the cost of data labeling. However, applying active learning to these datasets presents several challenges. One problem is that obtaining labels for data points in these datasets is expensive. Another issue arises from the lack of diversity in these datasets, posing a challenge for training models to generalize effectively to new data. However, despite these hurdles, active learning shows promise in improving the performance of machine learning models across CIFAR10, EuroSAT, and Fashion MNIST datasets. We anticipate that future research in active learning will prioritize addressing the specific challenges posed by these datasets and innovating more efficient active learning techniques.

ACKNOWLEDGMENT

We would like to express our heartfelt gratitude to our guide **Dr. Malini M Patil** Professor & Head, Department of Computer Science and Engineering, RV Institute of Technology and Management for her unwavering and consistent support at all times.

REFERENCES

[1] O. Sener and S. Savarese, "Active learning for convolutional neural networks: A core-set approach", International Conference on Learning Representations, 2018.

[2] S. Ravi and H. Larochelle, "Meta-learning for batch mode active learning", International Conference on Learning Representations workshop, 2018.

[3] L. Copa, D. Tuia, M. Volpi and M. Kanevski, "Unbiased query-by-bagging active learning for VHR image classification", Image and Signal Processing for RemoteSensing XVI, pp. 78300K, 2018.

[4] M. Li, R. Wang and K. Tang, "Combining SemiSupervised and active learning for hyperspectral image classification", Computational Intelligence and Data Mining (CIDM) 2013 IEEE Symposium on, pp. 89-94, 2013 [5] D. Wang and Y. Shang, "A new active labeling method for deep learning," 2014 International Joint Conference on Neural Networks (IJCNN), Beijing, China, 2014, pp. 112-119, doi: 10.1109/IJCNN.2014.6889457. [6] T. Kasarla, G. Nagendar, G. M. Hegde, V. Balasubramanian and C. V. Jawahar, "Region-based active learning for efficient labeling in semantic

- segmentation," 2019 IEEE Winter Conference on Applications of Computer Vision (WACV), Waikoloa, HI, USA, 2019, pp. 1109-1117, doi: 10.1109/WACV.2019.00123.
- [7] J. Z. Bengar, J. van de Weijer, B. Twardowski and B. Raducanu, "Reducing Label Effort: Self-Supervised meets Active Learning," 2021 IEEE/CVF International Conference on Computer Vision Workshops (ICCVW), Montreal, BC, Canada, 2021, pp. 1631-1639, doi: 10.1109/ICCVW54120.2021.00188.
- [8] G. Lenczner, A. Chan-Hon-Tong, B. Le Saux, N. Luminari and G. Le Besnerais, "DIAL: Deep Interactive and Active Learning for Semantic Segmentation in Remote Sensing," in IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing, vol.15, pp. 3376-3389, 2022, doi: 10.1109/JSTARS.2022.3166551.
- [9] Chandra, A.L., Desai, S.V., Balasubramanian, V.N. et al. Active learning with point supervision for cost-effective panicle detection in cereal crops. *Plant Methods* 16, 34 (2020). <https://doi.org/10.1186/s13007-020-00575-8>. [10] R. Sheikh, A. Milioto, P. Lottes, C. Stachniss, M. Bennewitz and T. Schultz, "Gradient and Log-based Active Learning for Semantic Segmentation of Crop and Weed for Agricultural Robots," 2020 IEEE International Conference on Robotics and Automation (ICRA), Paris, France, 2020, pp. 1350-1356, doi: 10.1109/ICRA40945.2020.9196722.
- [11] S. Desai and D. Ghose, "Active Learning for Improved Semi-Supervised Semantic Segmentation in Satellite Images," 2022 IEEE/CVF Winter Conference on Applications of Computer Vision (WACV), Waikoloa, HI, USA, 2022, pp. 1485-1495, doi: 10.1109/WACV51458.2022.00155.
- [12] R. Sheikh, A. Milioto, P. Lottes, C. Stachniss, M. Bennewitz and T. Schultz, "Gradient and Log-based Active Learning for Semantic Segmentation of Crop and Weed for Agricultural Robots," 2020 IEEE International Conference on Robotics and Automation (ICRA), Paris, France, 2020, pp. 1350-1356, doi: 10.1109/ICRA40945.2020.9196722.
- [13] S. Desai and D. Ghose, "Active Learning for Improved Semi-Supervised Semantic Segmentation in Satellite Images," 2022 IEEE/CVF Winter Conference on Applications of Computer Vision (WACV), Waikoloa, HI, USA, 2022, pp. 1485-1495, doi: 10.1109/WACV51458.2022.00155.
- [14] Settles, B. (2017). Active learning literature survey. Computer Sciences Technical Report 1648, University of Wisconsin-Madison.
- [15] Sener, O., & Savarese, S. (2018). Active learning for convolutional neural networks: A core-set approach. International Conference on Learning Representations (ICLR)
- [16] Gal, Y., & Ghahramani, Z. (2016). Dropout as a Bayesian approximation: Representing model uncertainty in deep learning. International Conference on Machine Learning (ICML).
- [17] Fang, Y., et al. (2017). Learning how to do active learning: A deep reinforcement learning approach. Advances in Neural Information Processing Systems (NeurIPS).
- [18] Beluch, W. H., et al. (2018). The power of ensembles for active learning in image classification. IEEE Conference on Computer Vision and Pattern Recognition (CVPR).
- [19] Wei, L., et al. (2020). Bayesian active learning for clinical decision support using electronic health records. Journal of the American Medical Informatics Association (JAMIA).
- [20] Siddiqui, S., et al. (2021). Active learning for graph classification: A comparative study. IEEE Transactions on Knowledge and Data Engineering (TKDE).
- [21] Yang, T., et al. (2022). Efficient active learning with graph neural networks. International Conference on Machine Learning (ICML).
- [22] Dasgupta, A., et al. (2020). Active learning with incomplete information. Advances in Neural Information Processing Systems (NeurIPS).
- [23] Ash, S. D., et al. (2019). Diversity-based active learning with applications to rare-class detection in data streams. Knowledge and Information Systems