

DIRECT: Deep Active Learning under Imbalance and Label Noise

Shyam Nuggehalli¹ Jifan Zhang¹ Lalit Jain² Robert Nowak¹

Abstract

Class imbalance is a prevalent issue in real world machine learning applications, often leading to poor performance in rare and minority classes. With an abundance of wild unlabeled data, active learning is perhaps the most effective technique in solving the problem at its root – collecting a more balanced and informative set of labeled examples during annotation. In this work, we propose a novel algorithm that first identifies the class separation threshold and then annotate the most uncertain examples from the minority classes, close to the separation threshold. Through a novel reduction to one-dimensional active learning, our algorithm DIRECT is able to leverage the classic active learning literature to address issues such as batch labeling and tolerance towards label noise. Compared to existing algorithms, our algorithm saves more than 15% of the annotation budget compared to state-of-art active learning algorithm and more than 90% of annotation budget compared to random sampling.

1. Introduction

Large-scale deep learning models is playing an increasing role across many industries. Human feedback and annotations have played a significant role in developing such systems. Increasingly, we believe the role of human would shift to annotating rare yet important cases such as safety risks. In this paper, we propose an active learning algorithm to address the class imbalance problem, where we sequentially and adaptively choose examples for annotation. Our proposed algorithm optimizes for the class-balancedness and informativeness of the annotation collected.

Inspired by previous work by Zhang et al. (2022), we propose a novel reduction of the imbalanced classification problem into an one-dimensional active learning problem by arranging unlabeled examples into a list ordered by uncer-

tainty scores. With this reduction, we apply classic active learning literature in finding the optimal separation threshold between the minority and majority classes. We then annotate examples near this threshold to gather uncertain examples.

On the other hand, existing state-of-art algorithm like GALAXY (Zhang et al., 2022) does not allow parallel annotation, and performs poorly when label noise is high. With our novel one-dimensional reduction, we inherit the noise tolerance and batch labeling advantages of classic active learning algorithms, allowing our *deep* active learning algorithm DIRECT to overcome such issues.

To summarize our main contributions:

- We propose a novel reduction that bridges the advancement in theoretical active learning literature to imbalanced active classification.
- Our novel algorithm addresses the prevalent imbalancedness issue by annotating a more class-balanced and informative set.
- Compared to state-of-art algorithm GALAXY (Zhang et al., 2022), our algorithm allows parallel annotation and is tolerant to label noise. It also saves 15% more annotation cost compared to GALAXY, and 90% of the annotation cost compared to random sampling.

2. Related Work

Deep Active Learning Under Class Balance Active Learning strategies sequentially and adaptively choose examples for annotation (Settles, 2009). Many uncertainty methods active learning methods extend the traditional active learning literature such as margin, least confidence and entropy sampling (Tong & Koller, 2001; Settles, 2009; Balcan et al., 2006; Kremer et al., 2014). These methods have also been shown to perform among top methods when fine-tuning large pretrained models and combined with semi-supervised learning algorithms (Zhang et al., 2023a). More sophisticated methods have also been proposed to optimize the chosen unlabeled examples’ uncertainty (Gal et al., 2017; Ducoffe & Precioso, 2018; Beluch et al., 2018), diversity (Sener & Savarese, 2017; Geifman & El-Yaniv, 2017; Citovsky et al., 2021), or a mix of both (Ash et al., 2019; 2021; Wang et al., 2021; Elenter et al., 2022; Mohamadi

¹University of Wisconsin-Madison ²Michael G. Foster School of Business, University of Washington. Correspondence to: Jifan Zhang <jifan@cs.wisc.edu>.

et al., 2022). However, these methods often perform poorly under prevalent and realistic scenarios such as under label noises (Khosla et al., 2022) or under class imbalance (Kothawade et al., 2021; Zhang et al., 2022; 2023a).

Deep Active Learning under Imbalance Data imbalance and rare instances are prevalent in almost all modern machine learning applications. Active learning techniques are effective in addressing the problem in its root by collecting a more class-balanced dataset (Aggarwal et al., 2020; Kothawade et al., 2021; Emam et al., 2021; Zhang et al., 2022; Coleman et al., 2022; Jin et al., 2022; Cai, 2022; Zhang et al., 2023b). To this end, Kothawade et al. (2021) propose a submodular-based method that actively annotates examples similar to known examples of rare instances. Our work is motivated by Zhang et al. (2022), where they propose GALAXY, an algorithm that construct one-dimensional linear graphs and apply graph-based active learning techniques in annotating a set of examples that are both class-balanced and uncertain. While GALAXY outperforms existing algorithms, due to a bisection procedure involved, it does not allow parallel annotation. In addition, bisection procedures are generally not robust against label noises, a prevalent challenge in real world annotation tasks. Our algorithm DIRECT mitigates all of the above shortcomings of GALAXY while outperforming it even with synchronous labeling and no label noise, beating GALAXY in its own game.

Lastly, we distinguish our work from Zhang et al. (2023b), where the paper studies the algorithm selection problem. Unlike our goal of proposing a new deep active learning algorithm, the paper proposes meta algorithms to choose the right active learning algorithm among a large number of candidate algorithms.

3. Preliminary

3.1. Notations

We study the pool-based active learning problem, where an initial unlabeled set of N examples $X = \{x_1, \dots, x_N\}$ are available for annotation. We let \mathcal{X} denote the space of possible examples. Their corresponding labels $Y = \{y_1, \dots, y_N\}$ are initially unknown. Furthermore, we study the multi-class classification problem, so the space of labels $\mathcal{Y} := [K]$, where each element denotes the class index and K is the total number of classes. Moreover, let N_1, \dots, N_K denote the number of examples in X in each class. We define the imbalance ratio as $\gamma = \frac{\min_{k \in [K]} N_k}{\max_{k \in [K]} N_k}$.

A deep active learning algorithm iteratively chooses batches of examples for annotation. During the t -th iteration, the algorithm is given labeled and unlabeled sets of examples, L_t and U_t respectively, where $L_t \cup U_t = X$ and $L_t \cap U_t = \emptyset$. The algorithm then chooses B examples from the unlabeled

set $X^{(t)} \subseteq U_t$ for their corresponding labels $Y^{(t)}$ to be annotated. The labeled and unlabeled sets are then updated, i.e., $L_{t+1} \leftarrow L_t \cup X^{(t)}$ and $U_{t+1} \leftarrow U_t \setminus X^{(t)}$. Based on new labeled set L_{t+1} and its corresponding labels, a neural network $f_t : \mathcal{X} \rightarrow \mathcal{Y}$ is trained to inform the choice for the next iteration.

The ultimate goal of deep active learning is to obtain high predictive accuracy for the trained neural network while annotating as few examples as possible.

3.2. Review of Imbalanced Active Learning: GALAXY

We first consider a binary case of imbalanced classification, where we assume $N_1 < N_2$ without loss of generality. After a total of T rounds and sampling TB examples, random sampling would annotate a subset of X with an imbalance ratio close to $\frac{N_1}{N_2}$. On the other hand, as shown previously in Zhang et al. (2022), uncertainty based methods such as confidence (Settles, 2009), margin (Tong & Koller, 2001; Balcan et al., 2006) and entropy (Kremer et al., 2014) sampling could annotate a set of examples with imbalance ratio arbitrarily low. The reasoning is simply demonstrated in Figure 1a, where choosing examples with predictive sigmoid score \hat{p} closest to .5 may end up annotating examples (all) from the majority class. It is also empirically shown that annotating examples most certain to be in the minority class is not the most effective in improving model performance.

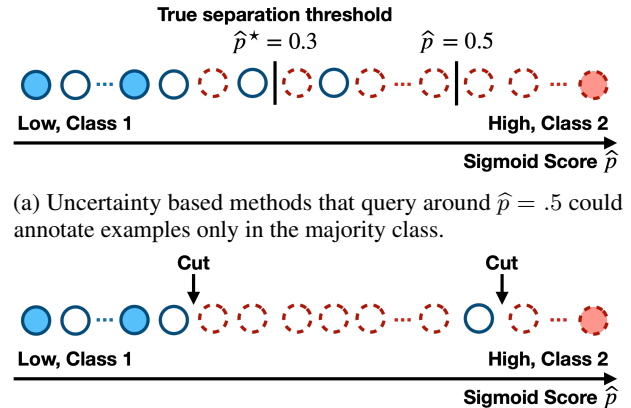


Figure 1: Figures demonstrating GALAXY. Ordered lists of examples are ranked by the predictive sigmoid score \hat{p} . The ground truth label of each example is represented by its border – solid blue lines for class 1 and dotted red lines for class 2. Annotated examples are shaded.

Therefore, the key objective in imbalanced active learning is labeling examples that are both *uncertain* and *class balanced*. The proposed algorithm GALAXY (Zhang

et al., 2022) follows the following two-phased procedure:

1. Identify the *optimal separation threshold* \hat{p}^* , as shown in Figure 1a.
2. Annotate examples in the minority class closest to \hat{p}^* .

GALAXY draws inspiration from graph-based active learning in identifying \hat{p}^* . In particular, it follows a modified bisection procedure that sequentially annotate examples to search for *all* cuts, namely thresholds with a class 1 example to the left and a class 2 example to the right (see Figure 1b). GALAXY then annotate examples close to all such cuts. However, the algorithm suffers from two weaknesses. Firstly, a bisection procedure only allows sequential labeling, which prevents multiple annotators labeling in parallel. Secondly, the problem of identifying the *optimal separation threshold* is a noisy procedure, since the predictive sigmoid scores may not perfectly separate the minority class examples from the majority. GALAXY mitigates the problem by wastefully identifying and annotating around all cuts, which could lead to annotating large number of majority class examples (shown in Figure 1b).

4. A Robust Algorithm with Batch Labeling

In this section, we formally define optimal separation threshold and pose the problem of identifying it as an 1-dimensional reduction to the agnostic active learning problem. We then propose a novel algorithm inspired by the agnostic active learning literature (Balcan et al., 2006; Dasgupta et al., 2007; Hanneke et al., 2014; Katz-Samuels et al., 2021).

4.1. An 1-D Reduction to Agnostic Active Learning

We start by considering the imbalanced binary classification setting mentioned in Section 3.2. When given a neural network model, we let $\hat{p} : \mathcal{X} \rightarrow [0, 1]$ be the predictive function mapping examples to sigmoid scores. Here, a higher sigmoid score represents a higher confidence of the example being in class 2. Let q_1, \dots, q_N denote the sigmoid scores of every example in the pool X , where $q_i = \hat{p}(x_i)$ for each $i \in [N]$. We also sort examples by the sigmoid predictive score similar to Section 3.2. Formally, we let $q_{(1)} \leq q_{(2)} \leq \dots \leq q_{(N)}$ be a sorted permutation of q_1, \dots, q_N , and let $y_{(1)}, y_{(2)}, \dots, y_{(N)}$ denote the sorted list's corresponding labels. We are now ready to define the optimal separation threshold as described in Section 3.2 below.

Definition 4.1. Given the sorted predictive sigmoid scores $0 \leq q_{(1)} \leq \dots \leq q_{(N)} \leq 1$, we define the *optimal separation threshold* as $j^* \in [N]$ such that

$$j^* = \arg \max_j \left(|\{y_{(i)} = 1 : 1 \leq i \leq j\}| - |\{y_{(i)} = 2 : 1 \leq i \leq j\}| \right). \quad (1)$$

In other words, j^* has the largest difference in the number of class 1 examples from class 2 examples to its left. Furthermore, ties are broken by choosing the largest j^* that attains the argmax if class 1 is the minority class and the lowest j^* otherwise.

1D Reduction. We now provide a reduction of finding j^* to an 1-dimensional agnostic active learning problem. In particular, let $\mathcal{H} = \{h_0, h_1, \dots, h_N\}$ be the hypothesis class. Here each hypothesis $h_j : [0, 1] \rightarrow \{1, 2\}$ maps sigmoid scores to binary target classes and

$$h_j(q) = \begin{cases} 1 & \text{if } q \leq q_{(j)} \\ 2 & \text{if } q > q_{(j)} \end{cases}.$$

The empirical zero-one loss for each hypothesis is then defined as $\mathcal{L}(h_j) = \sum_{i=1}^N \mathbf{1}\{h_j(q_{(i)}) \neq y_{(i)}\}$. In Appendix 6, we show that optimizing for the zero-one loss $\arg \min_{0 \leq j \leq N} \mathcal{L}(h_j)$ is equivalent to (1). Namely, $j^* \in \arg \min_{0 \leq j \leq N} \mathcal{L}(h_j)$, transforming the problem of finding the optimal separation threshold to an 1-D agnostic active learning problem.

Multi-Class Classification. To generalize the above problem formulation to multi-class classification, we follow a similar strategy to Zhang et al. (2022). Specifically, for each class k , one can view the problem of class- k v.s. others as a binary classification problem. In other words, class k can be seen as class 1 and classes $\{1, \dots, k-1, k+1, \dots, N\}$ can be seen as class 2. The goal therefore becomes finding all K optimal separation thresholds, which is equivalent with solving K 1-D agnostic active learning problems. Moreover, let $\hat{p} : \mathcal{X} \rightarrow \Delta^{(K-1)}$ denote the neural network prediction function mapping examples to sigmoid scores. Instead of using sigmoid scores, for class k , we use the margin scores $q_i^k := \max_{k'} [\hat{p}(x_i)]_{k'} - [\hat{p}(x_i)]_k$ to sort the examples. We also note that sorting with margin scores is equivalent to sorting with sigmoid scores in the binary classification setting, so this is a strictly more general formulation to multi-class classification.

4.2. Algorithm

We are now ready to state our algorithm as shown in Algorithm 1. The algorithm follows a two-phased, where the first phase aims to identify the optimal separation threshold for each class. The second phase then annotates examples closest to the estimated optimal separation thresholds for each class.

Algorithm 1 DIRECT: DIMension REDuction for aCTive Learning under Imbalance and Label Noise

Input: Pool X , number of rounds T , retraining batch size B_{train} , number of parallel annotations B_{parallel} .
Initialize: Uniformly sample B elements without replacement from X to form L_0 . Let $U_0 \leftarrow X \setminus L_0$.
for $t = 1, \dots, T - 1$ **do**
 Train neural network on L_{t-1} and obtain f_{t-1} .
 Find optimal separation thresholds
 Initialize labeled set $L_t \leftarrow L_{t-1}$.
 Initialize budget per class $b \leftarrow \frac{B_{\text{train}}}{2K}$.
 for k in $\text{RandPerm}(\{1, \dots, K\})$ **do**
 Compute margin scores q_1^k, \dots, q_N^k based on f_{t-1} and sort them by $q_{(1)}^k \leq \dots \leq q_{(N)}^k$.
 Find optimal separation threshold:
 $L_t \leftarrow \text{VReduce}(L_t, b, k, B_{\text{parallel}}, \{(x_{(i)}, y_{(i)})\}_{i=1}^N)$.
 end for
 Annotate rare and uncertain examples
 Compute budget per class $b \leftarrow \frac{B_{\text{train}} - |L_t|}{K}$.
 for k in $\text{RandPerm}(\{1, \dots, K\})$ **do**
 Estimate separation threshold

$$\hat{j}^k \leftarrow \arg \max_j (|\{y_{(i)} = 1 : x_{(i)} \in L_t \text{ and } i \leq j\}| - |\{y_{(i)} = 1 : x_{(i)} \in L_t \text{ and } i \leq j\}|)$$

 and break ties by choosing the index closest to $\frac{N}{2}$.
 Annotate b unlabeled examples with sorted indices closest to \hat{j}^k and insert to L_t .
 end for
end for
Return: Train final classifier f_T based on L_T .

To identify the optimal separation threshold, we loop over each class k and run a fix-budget batch active learning procedure for the one-dimensional reduction (as shown in Algorithm 2). Specifically, we keep a version space of the hypotheses with higher likelihood of being the optimal separation threshold. We annotate examples sequentially in batches and shrink the version space by a fixed rate each iteration. The shrinkage rate c is determined by the budget and batch size, so that after the final iteration, the version space has exactly one hypothesis left.

To address batch labeling, we note the number of examples one collect before retraining is usually far greater than the number of annotators annotating in parallel. In particular, we let B_{train} denote the number of examples the algorithm collects before the neural network is retrained. In practice, this number is usually determined by the practical constraints of computational training cost. On the other hand, we let B_{parallel} denote the number of examples annotated in parallel. As discussed above, $B_{\text{parallel}} \ll B_{\text{train}}$ in practice.

Lastly, as will be discussed in Section 6, a variant of our algorithm can also be proposed for asynchronous labeling.

Algorithm 2 VReduce: Version Space Reduction

Input: Labeled set L , budget b , class of interest k , parallel batch size B_{parallel} , examples and ground truth labels sorted by uncertainty $\{(x_{(i)}, y_{(i)})\}_{i=1}^N$ (ground truth labels are hidden to the learner).
Initialize: Version space: shortest segment of indices $[I, J]$, where $X_I, X_J \in L, \forall i \leq I, x_{(i)} \in L \implies y_{(i)} = k$ and $\forall j \geq J, x_{(j)} \in L \implies y_{(j)} \neq k$.
Initialize: Number of iterations $m \leftarrow \frac{b}{B_{\text{parallel}}}$. Shrinking factor $c \leftarrow \sqrt[m]{J - I}$.
for $t = 1, \dots, m$ **do**
 Sample uniformly at random B_{parallel} unlabeled examples in $x_{(I)}, \dots, x_{(J)}$ for annotation and insert to L .
 Find $[I', J'] \subset [I, J]$ such that

$$I', J' = \arg \min_{i, j: j-i = \frac{1}{c}(J-I)} \max\{\hat{\mathcal{L}}^k(i), \hat{\mathcal{L}}^k(j)\} \quad \text{where}$$

$$\hat{\mathcal{L}}^k(s) = \sum_{\substack{r \leq s: \\ x_{(r)} \in L}} \mathbf{1}\{y_{(r)} \neq k\} + \sum_{\substack{r > s: \\ x_{(r)} \in L}} \mathbf{1}\{y_{(r)} = k\}.$$

 Update version space $I, J \leftarrow I', J'$.
end for
Return: Updated labeled set L .

5. Experiments

5.1. Experiment Setup

Our experiment setup follows closely from Zhang et al. (2022). We use the CIFAR-10, CIFAR-100 (Krizhevsky et al., 2009) and SVHN (Netzer et al., 2011) image classification datasets. The original CIFAR-10 and SVHN datasets have 10 classes while the CIFAR-100 dataset has 100 classes. We construct an extremely imbalanced dataset by grouping a large number of classes into one majority class. For example, suppose we have a dataset with M classes. We generate an imbalanced dataset with K classes ($K < M$) by reusing classes $1, \dots, K - 1$ from the original dataset and combining the rest of the classes K, \dots, M into a single majority class K . Detailed imbalance ratios are shown in Table 1.

We compare our algorithm DIRECT against nine baselines: GALAXY (Zhang et al., 2022), SIMILAR (Kothawade et al., 2021), BADGE (Ash et al., 2019), BASE (Emam et al., 2021), BAIT (Ash et al., 2021), Cluster Margin (Citovsky et al., 2021), Confidence Sampling (Settles, 2009), Most Likely Positive (Jiang et al., 2018; Warmuth et al., 2001; 2003) and Random Sampling. The highlight surrounding the baselines represents the standard error plotted as the

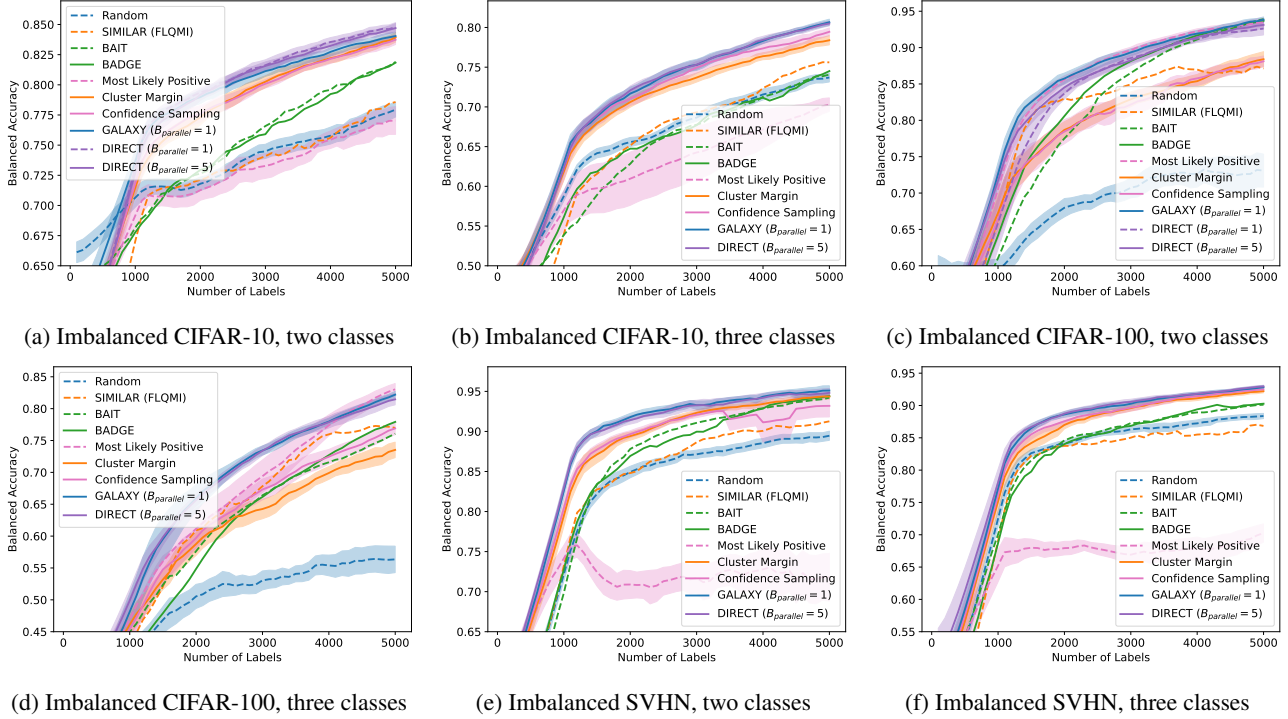


Figure 2: Performance of DIRECT vs different baselines across six different dataset settings.

NAME	K	N_K	$\sum_{k=1}^{K-1} N_k$	IMB RATIO γ
CIFAR-10	2	45000	5000	.1111
CIFAR-10	3	40000	10000	.1250
CIFAR-100	2	49500	500	.0101
CIFAR-100	3	49000	1000	.0102
SVHN	2	68309	4948	.0724
SVHN	3	54448	18809	.2546

Table 1: Dataset setting for extremely unbalanced scenario. N_K denotes the number of images in the majority class while $\sum_{k=1}^{K-1} N_k$ is the total number of images in all minority classes. ϵ is the class imbalance factor defined in Section 3.1.

confidence intervals. Due to computational constraints, we only have single runs for SIMILAR, BADGE, and BAIT. For all other baselines, we average over four runs.

For all experiments, we use PyTorch’s ResNet-18 model pretrained on ImageNet for initialization. We use the Adam optimization algorithm with 500 epochs for each L and the learning rate set to 0.01.

5.2. Results

Synchronous Labeling ($B_{\text{parallel}} = 1$). In our first set of experiments, we compare DIRECT against GALAXY in the synchronous labeling setting by setting $B_{\text{parallel}} =$

1 in our algorithm. As shown in Figure 2(a), DIRECT outperforms GALAXY by saving an additional 15% of the labeling budget. Compared to random sampling, DIRECT is able to around 90% of annotation cost.

Batch Labeling. We observe that on all six datasets shown in Figure 2, DIRECT performs comparably to GALAXY even when under parallel annotation ($B_{\text{parallel}} = 5$). However, we do observe a slight drop in label-efficiency when comparing to synchronous labeling.

6. Future Work

In this paper, we addressed the batch sampling problem of GALAXY by annotating multiple examples in parallel. However, in practice, examples are often annotated asynchronously. We believe this is a valuable future direction and a potential solution is to utilize an asynchronous variant of one-dimensional active learning algorithm.

Acknowledgements

This work has been supported in part by NSF Award 2112471.

References

- Aggarwal, U., Popescu, A., and Hudelot, C. Active learning for imbalanced datasets. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pp. 1428–1437, 2020.
- Ash, J., Goel, S., Krishnamurthy, A., and Kakade, S. Gone fishing: Neural active learning with fisher embeddings. *Advances in Neural Information Processing Systems*, 34: 8927–8939, 2021.
- Ash, J. T., Zhang, C., Krishnamurthy, A., Langford, J., and Agarwal, A. Deep batch active learning by diverse, uncertain gradient lower bounds. *arXiv preprint arXiv:1906.03671*, 2019.
- Balcan, M.-F., Beygelzimer, A., and Langford, J. Agnostic active learning. In *Proceedings of the 23rd international conference on Machine learning*, pp. 65–72, 2006.
- Beluch, W. H., Genewein, T., Nürnberger, A., and Köhler, J. M. The power of ensembles for active learning in image classification. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 9368–9377, 2018.
- Cai, X. Active learning for imbalanced data: The difficulty and proportions of class matter. *Wireless Communications and Mobile Computing*, 2022, 2022.
- Citovsky, G., DeSalvo, G., Gentile, C., Karydas, L., Rajagopalan, A., Rostamizadeh, A., and Kumar, S. Batch active learning at scale. *Advances in Neural Information Processing Systems*, 34:11933–11944, 2021.
- Coleman, C., Chou, E., Katz-Samuels, J., Culatana, S., Bailis, P., Berg, A. C., Nowak, R., Sumbaly, R., Zaharia, M., and Yalniz, I. Z. Similarity search for efficient active learning and search of rare concepts. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pp. 6402–6410, 2022.
- Dasgupta, S., Hsu, D. J., and Monteleoni, C. A general agnostic active learning algorithm. *Advances in neural information processing systems*, 20, 2007.
- Ducoffe, M. and Precioso, F. Adversarial active learning for deep networks: a margin based approach. *arXiv preprint arXiv:1802.09841*, 2018.
- Elenter, J., NaderiAlizadeh, N., and Ribeiro, A. A lagrangian duality approach to active learning. *arXiv preprint arXiv:2202.04108*, 2022.
- Emam, Z. A. S., Chu, H.-M., Chiang, P.-Y., Czaja, W., Leapman, R., Goldblum, M., and Goldstein, T. Active learning at the imagenet scale. *arXiv preprint arXiv:2111.12880*, 2021.
- Gal, Y., Islam, R., and Ghahramani, Z. Deep bayesian active learning with image data. In *International Conference on Machine Learning*, pp. 1183–1192. PMLR, 2017.
- Geifman, Y. and El-Yaniv, R. Deep active learning over the long tail. *arXiv preprint arXiv:1711.00941*, 2017.
- Hanneke, S. et al. Theory of disagreement-based active learning. *Foundations and Trends® in Machine Learning*, 7(2-3):131–309, 2014.
- Jiang, S., Malkomes, G., Abbott, M., Moseley, B., and Garnett, R. Efficient nonmyopic batch active search. In *32nd Conference on Neural Information Processing Systems (NeurIPS 2018)*, 2018.
- Jin, Q., Yuan, M., Wang, H., Wang, M., and Song, Z. Deep active learning models for imbalanced image classification. *Knowledge-Based Systems*, 257:109817, 2022.
- Katz-Samuels, J., Zhang, J., Jain, L., and Jamieson, K. Improved algorithms for agnostic pool-based active classification. In *International Conference on Machine Learning*, pp. 5334–5344. PMLR, 2021.
- Khosla, S., Whye, C. K., Ash, J. T., Zhang, C., Kawaguchi, K., and Lamb, A. Neural active learning on heteroskedastic distributions. *arXiv preprint arXiv:2211.00928*, 2022.
- Kothawade, S., Beck, N., Killamsetty, K., and Iyer, R. Similar: Submodular information measures based active learning in realistic scenarios. *Advances in Neural Information Processing Systems*, 34:18685–18697, 2021.
- Kremer, J., Steenstrup Pedersen, K., and Igel, C. Active learning with support vector machines. *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, 4(4):313–326, 2014.
- Krizhevsky, A., Hinton, G., et al. Learning multiple layers of features from tiny images. 2009.
- Mohamadi, M. A., Bae, W., and Sutherland, D. J. Making look-ahead active learning strategies feasible with neural tangent kernels. *arXiv preprint arXiv:2206.12569*, 2022.
- Netzer, Y., Wang, T., Coates, A., Bissacco, A., Wu, B., and Ng, A. Y. Reading digits in natural images with unsupervised feature learning. 2011.
- Sener, O. and Savarese, S. Active learning for convolutional neural networks: A core-set approach. *arXiv preprint arXiv:1708.00489*, 2017.
- Settles, B. Active learning literature survey. 2009.
- Tong, S. and Koller, D. Support vector machine active learning with applications to text classification. *Journal of machine learning research*, 2(Nov):45–66, 2001.

- Wang, H., Huang, W., Margenot, A., Tong, H., and He, J. Deep active learning by leveraging training dynamics. *arXiv preprint arXiv:2110.08611*, 2021.
- Warmuth, M. K., Rätsch, G., Mathieson, M., Liao, J., and Lemmen, C. Active learning in the drug discovery process. In *NIPS*, pp. 1449–1456, 2001.
- Warmuth, M. K., Liao, J., Rätsch, G., Mathieson, M., Putta, S., and Lemmen, C. Active learning with support vector machines in the drug discovery process. *Journal of chemical information and computer sciences*, 43(2):667–673, 2003.
- Zhang, J., Katz-Samuels, J., and Nowak, R. Galaxy: Graph-based active learning at the extreme. *arXiv preprint arXiv:2202.01402*, 2022.
- Zhang, J., Chen, Y., Canal, G., Mussmann, S., Das, A. M., Bhatt, G., Zhu, Y., Du, S. S., Jamieson, K., and Nowak, R. D. Labelbench: A comprehensive framework for benchmarking adaptive label-efficient learning, 2023a.
- Zhang, J., Shao, S., Verma, S., and Nowak, R. Algorithm selection for deep active learning with imbalanced datasets. *arXiv preprint arXiv:2302.07317*, 2023b.

Appendix A.

Lemma .1. *The agnostic active learning reduction is equivalently finding the optimal separation threshold. Namely,*

$$\arg \min_j \mathcal{L}(h_j) = \arg \max_j (|\{y_{(i)} = 1 : 1 \leq i \leq j\}| - |\{y_{(i)} = 2 : 1 \leq i \leq j\}|)$$

Proof. Recall the definitions: $h_j(q) = \begin{cases} 1 & \text{if } q \leq q_{(j)} \\ 2 & \text{if } q > q_{(j)} \end{cases}$ and $\mathcal{L}(h_j) = \sum_{i=1}^N \mathbf{1}\{h_j(q_{(i)}) \neq y_{(i)}\}$, we can expand the loss as follows

$$\begin{aligned} \arg \min_j \mathcal{L}(h_j) &= \arg \min_j \sum_{i=1}^N \mathbf{1}\{h_j(q_{(i)}) \neq y_{(i)}\} \\ &= \arg \min_j N - \sum_{i=1}^N \mathbf{1}\{h_j(q_{(i)}) = y_{(i)}\} \\ &= \arg \max_j \sum_{i=1}^N \mathbf{1}\{h_j(q_{(i)}) = y_{(i)}\} \\ &= \arg \max_j \left(\sum_{i=1}^j \mathbf{1}\{y_{(i)} = 1\} \right) + \left(\sum_{i=j+1}^N \mathbf{1}\{y_{(i)} = 2\} \right) \\ &= \arg \max_j \left(\sum_{i=1}^j \mathbf{1}\{y_{(i)} = 1\} \right) + \left(\sum_{i=j+1}^N \mathbf{1}\{y_{(i)} = 2\} \right) - \left(\sum_{i=1}^N \mathbf{1}\{y_{(i)} = 2\} \right) \\ &= \arg \max_j \sum_{i=1}^j (\mathbf{1}\{y_{(i)} = 1\} - \mathbf{1}\{y_{(i)} = 2\}) \end{aligned}$$

□