

# CLINICAL: Targeted Active Learning for Imbalanced Medical Image Classification

Suraj Kothawade<sup>1</sup>, Atharv Savarkar<sup>2</sup>, Venkat Iyer<sup>2</sup>, Lakshman Tamil<sup>1</sup>, Ganesh Ramakrishnan<sup>2</sup>, and Rishabh Iyer<sup>1</sup>

<sup>1</sup> University of Texas at Dallas, USA

<sup>2</sup> Indian Institute of Technology, Bombay, India

suraj.kothawade@utdallas.edu

**Abstract.** Training deep learning models on medical datasets that perform well for all classes is a challenging task. It is often the case that a suboptimal performance is obtained on some classes due to the natural class imbalance issue that comes with medical data. An effective way to tackle this problem is by using *targeted active learning*, where we iteratively add data points that belong to the rare classes, to the training data. However, existing active learning methods are ineffective in targeting rare classes in medical datasets. In this work, we propose CLINICAL (targeted aCtive Learning for ImbalANced medICAL imAge cLassification) a framework that uses submodular mutual information functions as acquisition functions to mine critical data points from rare classes. We apply our framework to a wide-array of medical imaging datasets on a variety of real-world class imbalance scenarios - namely, *binary* imbalance and *long-tail* imbalance. We show that CLINICAL outperforms the state-of-the-art active learning methods by acquiring a diverse set of data points that belong to the rare classes.

## 1 Introduction

Owing to the advancement of deep learning, medical image classification has made tremendous advances in the past decade. However, medical datasets are naturally imbalanced at the class level, *i.e.*, some classes are comparatively rarer than the others. For instance, cancerous classes are naturally rarer than non-cancerous ones. In such scenarios, the over-represented classes *overpower* the training process and the model ends up learning a biased representation. Deploying such biased models results in incorrect predictions, which can be catastrophic and even lead to loss of life. Active learning (AL) is a promising solution to mitigate this imbalance in the training dataset. The goal of AL is to select data points from an unlabeled set for addition to the training dataset at an additional labeling cost. The model is then retrained with the new training set and the process is repeated. Reducing the labeling cost using the AL paradigm is crucial in domains like medical imaging, where labeling data requires expert supervision (*e.g.*, doctors), which makes the process extremely expensive. However, current AL methods are inefficient in selecting data points from the rare classes in medical image datasets. Broadly,

they use acquisition functions that are either: i) based on the uncertainty scores of the model, which are used to select the top uncertain data points [26], or ii) based on diversity scores, where data points having diverse gradients are selected [3, 25]. They mainly focus on improving the overall performance of the model, and thereby fail to target these rare yet critical classes. Unfortunately, this leads to a wastage of expensive labeling resources when the goal is to improve performance on these rare classes.

In this work, we consider two types of class imbalance that recur in a wide array of medical imaging datasets. The first scenario is *binary* imbalance, where a subset of classes is rare/infrequent and the remaining subset is relatively frequent. The second scenario is that of *long-tail* imbalance, where the frequency of data points from each class keeps *steeply* reducing as we go from the most frequent class to the rarest class (see Fig. 1). Such class imbalance scenarios are particularly challenging in the medical imaging domain since there exist subtle differences which are barely visually evident (see Fig. 1). In Sec. 3, we discuss CLINICAL, a targeted active learning algorithm that acquires a subset by maximizing the submodular mutual information with a set of *misclassified* data points from the rare classes. This enables us to focus on data points from the unlabeled set that are critical and belong to the rare classes.

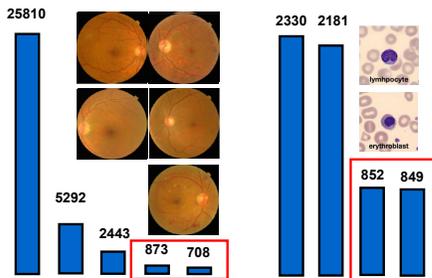


Fig. 1: Motivating examples of two main class imbalance scenarios occurring in medical imaging. **Left:** Long-tail imbalance (Diabetic retinopathy grading from retinal images in APTOS-2019 [10]). **Right:** Binary imbalance (Microscopic peripheral blood cell image classification in Blood-MNIST [1]). Red boxes in both scenario denote targeted rare classes.

### 1.1 Related work

**Uncertainty based Active Learning.** Uncertainty based methods aim to select the most uncertain data points according to a model for labeling. The most common techniques are - 1) ENTROPY [26] selects data points with maximum entropy, 2) LEAST CONFIDENCE [29] selects data points with the lowest confidence, and 3) MARGIN [24] selects data points such that the difference between the top two predictions is minimum.

**Diversity based Active Learning.** The main drawback of uncertainty based methods is that they lack diversity within the acquired subset. To mitigate this, a number of approaches have proposed to incorporate diversity. The CORESET method [25] minimizes a coresets loss to form coresets that represent the geometric structure of the original dataset. They do so using a greedy  $k$ -center clustering. A

recent approach called BADGE [3] uses the last linear layer gradients to represent data points and runs K-MEANS++ [2] to obtain centers that have a high gradient magnitude. The centers being representative and having high gradient magnitude ensures uncertainty and diversity at the same time. However, for batch AL, BADGE models diversity and uncertainty only within the batch and *not* across all batches. Another method, BATCHBALD [15] requires a large number of Monte Carlo dropout samples to obtain significant mutual information which limits its application to medical domains where data is scarce.

**Class Imbalanced and Personalized Active Learning.** Closely related to our method CLINICAL, are methods which optimize an objective that involves a held-out set. GRADMATCH [13] uses an orthogonal matching pursuit algorithm to select a subset whose gradient closely matches the gradient of a validation set. Another method, GLISTER-ACTIVE [14] formulates an acquisition function that maximizes the log-likelihood on a held-out validation set. We adopt GRADMATCH and GLISTER-ACTIVE as baselines that *targets* rare classes in our class imbalance setting and refer to it T-GRADMATCH and T-GLISTER in Sec. 4. Recently, [16] proposed the use of submodular information measures for active learning in realistic scenarios, while [17] used them to find rare objects in an autonomous driving object detection dataset. Finally, [19] use the submodular mutual information functions (used here) for personalized speech recognition. Our proposed method uses the submodular mutual information to target selecting data points from the rare classes via using a small set of *misclassified* data points as exemplars, which makes our method applicable to binary as well as long-tail imbalance scenarios.

## 1.2 Our contributions

We summarize our contributions as follows: **1)** We emphasize on the issue of binary and long-tail class imbalance in medical datasets that leads to poor performance on rare yet critical classes. **2)** Given the limitations of current AL methods on medical datasets, we propose CLINICAL, a novel AL framework that can be applied to any class imbalance scenario. **3)** We demonstrate the effectiveness of our framework for a diverse set of image classification tasks and modalities on Pneumonia-MNIST [12], Path-MNIST [11], Blood-MNIST [1], APTOS-2019 [10], and ISIC-2018 [4] datasets. Furthermore, we show that CLINICAL outperforms the state-of-the-art AL methods by up to  $\approx 6\% - 10\%$  on an average in terms of the average rare classes accuracy for binary imbalance scenarios and long-tail imbalance scenarios. **4)** We provide valuable insights about the *choice* of submodular functions to be used for subset selection based on the *modality* of medical data.

## 2 Preliminaries

**Submodular Functions:** We let  $\mathcal{V}$  denote the *ground-set* of  $n$  data points  $\mathcal{V} = \{1, 2, 3, \dots, n\}$  and a set function  $f : 2^{\mathcal{V}} \rightarrow \mathbb{R}$ . The function  $f$  is submodular [5] if it satisfies diminishing returns, namely  $f(j|\mathcal{X}) \geq f(j|\mathcal{Y})$  for all  $\mathcal{X} \subseteq \mathcal{Y} \subseteq \mathcal{V}, j \notin \mathcal{Y}$ . Facility location, graph cut, log determinants, *etc.* are some examples [9].

**Submodular Mutual Information (SMI):** Given a set of items  $\mathcal{A}, \mathcal{Q} \subseteq \mathcal{V}$ , the submodular mutual information (MI) [6, 8] is defined as  $I_f(\mathcal{A}; \mathcal{Q}) = f(\mathcal{A}) + f(\mathcal{Q}) - f(\mathcal{A} \cup \mathcal{Q})$ . Intuitively, this measures the similarity between  $\mathcal{Q}$  and  $\mathcal{A}$  and we refer to  $\mathcal{Q}$  as the query set. [18] extend SMI to handle the case when the *target* can come from a different set  $\mathcal{V}'$  apart from the ground set  $\mathcal{V}$ . In the context of imbalanced medical image classification,  $\mathcal{V}$  is the source set of images and the query set  $\mathcal{Q}$  is the target set containing the rare class images. To find an optimal subset given a query set  $\mathcal{Q} \subseteq \mathcal{V}'$ , we can define  $g_{\mathcal{Q}}(\mathcal{A}) = I_f(\mathcal{A}; \mathcal{Q})$ ,  $\mathcal{A} \subseteq \mathcal{V}$  and maximize the same.

## 2.1 Examples of SMI functions

For targeted active learning, we use the recently introduced SMI functions in [8, 6] and their extensions introduced in [18] as acquisition functions. For any two data points  $i \in \mathcal{V}$  and  $j \in \mathcal{Q}$ , let  $s_{ij}$  denote the similarity between them.

**Graph Cut MI (GCMI):** The SMI instantiation of graph-cut (GCMI) is defined as:  $I_{GC}(\mathcal{A}; \mathcal{Q}) = 2 \sum_{i \in \mathcal{A}} \sum_{j \in \mathcal{Q}} s_{ij}$ . Since maximizing GCMI maximizes the joint pairwise sum with the query set, it will lead to a summary similar to the query set  $\mathcal{Q}$ . In fact, specific instantiations of GCMI have been intuitively used for query-focused summarization for videos [28] and documents [21, 20].

**Facility Location MI (FLMI):** We consider two variants of FLMI. The first variant is defined over  $\mathcal{V}$  (FLVMI), the SMI instantiation can be defined as:  $I_{FLV}(\mathcal{A}; \mathcal{Q}) = \sum_{i \in \mathcal{V}} \min(\max_{j \in \mathcal{A}} s_{ij}, \max_{j \in \mathcal{Q}} s_{ij})$ . The first term in the  $\min(\cdot)$  of FLVMI models diversity, and the second term models query relevance.

For the second variant, which is defined over  $\mathcal{Q}$  (FLQMI), the SMI instantiation can be defined as:  $I_{FLQ}(\mathcal{A}; \mathcal{Q}) = \sum_{i \in \mathcal{Q}} \max_{j \in \mathcal{A}} s_{ij} + \sum_{i \in \mathcal{A}} \max_{j \in \mathcal{Q}} s_{ij}$ . FLQMI is very intuitive for query relevance as well. It measures the representation of data points that are the most relevant to the query set and vice versa.

**Log Determinant MI (LOGDETMI):** The SMI instantiation of LOGDETMI can be defined as:  $I_{LogDet}(\mathcal{A}; \mathcal{Q}) = \log \det(S_{\mathcal{A}}) - \log \det(S_{\mathcal{A}} - S_{\mathcal{A}, \mathcal{Q}} S_{\mathcal{Q}}^{-1} S_{\mathcal{A}, \mathcal{Q}}^T)$ .  $S_{\mathcal{A}, \mathcal{Q}}$  denotes the cross-similarity matrix between the items in sets  $\mathcal{A}$  and  $\mathcal{Q}$ .

## 3 CLINICAL: Our Targeted Active Learning framework for Binary and Long-tail Imbalance

In this section, we propose our targeted active learning framework, CLINICAL (see Fig. 2), and show how it can be applied to datasets with class imbalance. Concretely, we apply the SMI functions as acquisition functions for improving a model’s accuracy on rare classes at a given additional labeling cost ( $B$  instances) without compromising on the overall accuracy. The main idea in CLINICAL, is to use *only the misclassified* data points from a held-out target set  $\mathcal{T}$  containing data points from the rare classes. Let  $\hat{\mathcal{T}} \subseteq \mathcal{T}$  be the subset of misclassified data points. Then, we optimize the SMI function  $I_f(\mathcal{A}; \hat{\mathcal{T}})$  using a greedy strategy [23].

Note that since  $\hat{\mathcal{T}}$  contains only the misclassified data points, it would contain more data points from classes that are comparatively *rarer* or the worst perform-

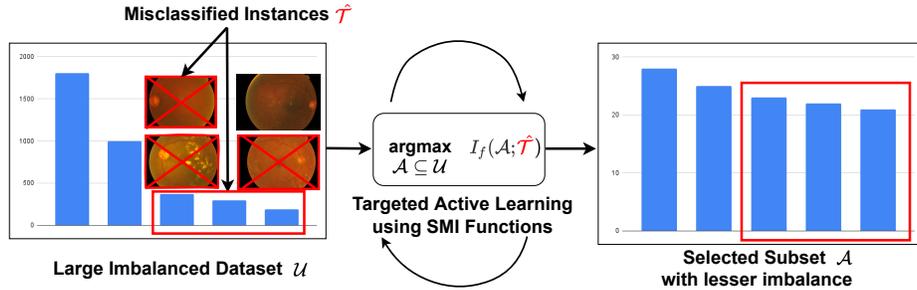


Fig. 2: The CLINICAL framework. We use a set of misclassified instances  $\hat{\mathcal{T}}$  as the query set  $\mathcal{Q}$  in the SMI function. We then maximize  $I_f(\mathcal{A}; \hat{\mathcal{T}})$  in an AL loop to target the imbalance and gradually mine data points from the rare classes.

ing. Moreover,  $\hat{\mathcal{T}}$  is updated in every AL round, this mechanism helps the SMI functions to focus on classes that require the most attention. For instance, in the long-tail imbalance scenario (see Fig. 1), CLINICAL would focus more on the tail classes in the initial rounds of AL. Next, we discuss the CLINICAL algorithm in detail:

**Algorithm:** Let  $\mathcal{L}$  be an initial training set of labeled instances and  $\mathcal{T}$  be the target set containing examples from the rare classes. Let  $\mathcal{U}$  be a large unlabeled dataset and  $\mathcal{M}$  be the trained model using  $\mathcal{L}$ . Next, we compute  $\hat{\mathcal{T}}$  as the subset of data points from  $\mathcal{T}$  that were misclassified by  $\mathcal{M}$ . Using last layer gradients as a representation for each data point which are extracted from  $\mathcal{M}$ , we compute similarity kernels of elements within  $\mathcal{U}$ , within  $\hat{\mathcal{T}}$  and between  $\mathcal{U}$  and  $\hat{\mathcal{T}}$  to instantiate an SMI function  $I_f(\mathcal{A}; \hat{\mathcal{T}})$  and maximize it to compute an optimal subset  $\mathcal{A} \subseteq \mathcal{U}$  of size  $B$  given  $\hat{\mathcal{T}}$  as target (query) set. We then augment  $\mathcal{L}$  with labeled  $\mathcal{A}$  (i.e.  $L(\mathcal{A})$ ) and re-train the model to improve the model on the rare classes. We summarize CLINICAL in Algorithm 1 and discuss its scalability aspects in Appendix. C.

---

**Algorithm 1** CLINICAL: Targeted AL for binary and long-tail imbalance

---

**Require:** Initial Labeled set of data points:  $\mathcal{L}$ , unlabeled dataset:  $\mathcal{U}$ , target set:  $\mathcal{T}$ , Loss function  $\mathcal{H}$  for learning model  $\mathcal{M}$ , batch size:  $B$ , number of selection rounds:  $N$

- 1: **for** selection round  $i = 1 : N$  **do**
  - 2:   Train  $\mathcal{M}$  with loss  $\mathcal{H}$  on the current labeled set  $\mathcal{L}$  and obtain parameters  $\theta_i$
  - 3:   Compute  $\hat{\mathcal{T}} \subseteq \mathcal{T}$  that were misclassified by the trained model  $\mathcal{M}$
  - 4:   Use  $\mathcal{M}_{\theta_i}$  to compute gradients using hypothesized labels  $\{\nabla_{\theta} \mathcal{H}(x_j, \hat{y}_j, \theta), \forall j \in \mathcal{U}\}$  and obtain a pairwise similarity matrix  $X$ . **{where  $X_{ij} = \langle \nabla_{\theta} \mathcal{H}_i(\theta), \nabla_{\theta} \mathcal{H}_j(\theta) \rangle$ }**
  - 5:   Instantiate a submodular function  $f$  based on  $X$ .
  - 6:    $\mathcal{A}_i \leftarrow \operatorname{argmax}_{\mathcal{A} \subseteq \mathcal{U}, |\mathcal{A}| \leq B} I_f(\mathcal{A}; \hat{\mathcal{T}})$
  - 7:   Get labels  $L(\mathcal{A}_i)$  for batch  $\mathcal{A}_i$ , and  $\mathcal{L} \leftarrow \mathcal{L} \cup L(\mathcal{A}_i)$ ,  $\mathcal{U} \leftarrow \mathcal{U} - \mathcal{A}_i$
  - 8:    $\mathcal{T} \leftarrow \mathcal{T} \cup \mathcal{A}_i^T$ , augment  $\mathcal{T}$  with new data points that belong to target classes.
  - 9: **end for**
  - 10: Return trained model  $\mathcal{M}$  and parameters  $\theta_N$ .
-

## 4 Experiments

In this section, we evaluate the effectiveness of CLINICAL on binary imbalance (Sec. 4.1) and long-tail imbalance (Sec. 4.2) scenarios. We do so by comparing the accuracy and class selections of various SMI functions with the existing state-of-the-art AL approaches. In our experiments, we observe that different SMI functions outperform existing approaches depending on the modality of the medical data. We show that the choice of the SMI based acquisition function is imperative and varies based on the imbalance scenario and the modality of medical data. Furthermore, in Appendix. D.2, we provide penalty matrices which show that CLINICAL statistically significantly outperforms the existing methods in all scenarios for multiple modalities.

**Baselines in all scenarios.** We compare the performance on CLINICAL against a variety of state-of-the-art uncertainty, diversity and targeted selection methods. The uncertainty based methods include ENTROPY, LEAST CONFIDENCE (LEAST-CONF), and MARGIN. The diversity based methods include CORESET and BADGE. The targeted selection methods include T-GLISTER and T-GRADMATCH. We discuss the details of all baselines in Sec. 1.1. For a fair comparison with CLINICAL, we use the same target set of misclassified data points  $\mathcal{T}$  as the held out validation set used in T-GLISTER and T-GRADMATCH. Lastly, we compare with random sampling (RANDOM).

**Experimental setup:** We use the same training procedure and hyperparameters for all AL methods to ensure a fair comparison. For all experiments, we train a ResNet-18 [7] model using an SGD optimizer with an initial learning rate of 0.001, the momentum of 0.9, and a weight decay of 5e-4. For each AL round, the weights are reinitialized using Xavier initialization and the model is trained till 99% training accuracy. The learning rate is decayed using cosine annealing [22] in every epoch. We run each experiment 5 $\times$  on a V100 GPU and provide the error bars (std deviation). We discuss dataset splits for each our experiments below and provide more details in Appendix. B.

### 4.1 Binary Imbalance

**Datasets:** We apply our framework to **1**)Pneumonia-MNIST (pediatric chest X-ray) [30, 12], **2**)Path-MNIST (colorectal cancer histology) [30, 11], and **3**)Blood-MNIST (blood cell microscope) [30, 1] medical image classification datasets. To create a more realistic medical scenario, we create a custom dataset that simulates binary class imbalance for each of these datasets for our experiments. Let  $\mathcal{C}$  be the set of data points from the rare classes and  $\mathcal{D}$  be the set of data points from the over-represented classes. We create the initial labeled set  $\mathcal{L}$  (seed set) in AL,  $|\mathcal{D}_{\mathcal{L}}| = \rho|\mathcal{C}_{\mathcal{L}}|$  and an unlabeled set  $\mathcal{U}$  such that  $|\mathcal{D}_{\mathcal{U}}| = \rho|\mathcal{C}_{\mathcal{U}}|$ , where  $\rho$  is the imbalance factor. We use a small held out target set  $\mathcal{T}$  which contains data points from the rare classes. For Path-MNIST and PneumoniaMNIST, we use  $\rho = 20$ , and for Blood-MNIST, we use  $\rho = 7$  due to the small size of the dataset. For Pneumonia-MNIST,  $|\mathcal{C}_{\mathcal{L}}| + |\mathcal{D}_{\mathcal{L}}| = 105$ ,  $|\mathcal{C}_{\mathcal{U}}| + |\mathcal{D}_{\mathcal{U}}| = 1100$ ,  $B = 10$  (AL batch size) and,  $|\mathcal{T}| = 5$ . Following the natural class imbalance, we use the ‘pneumonia’

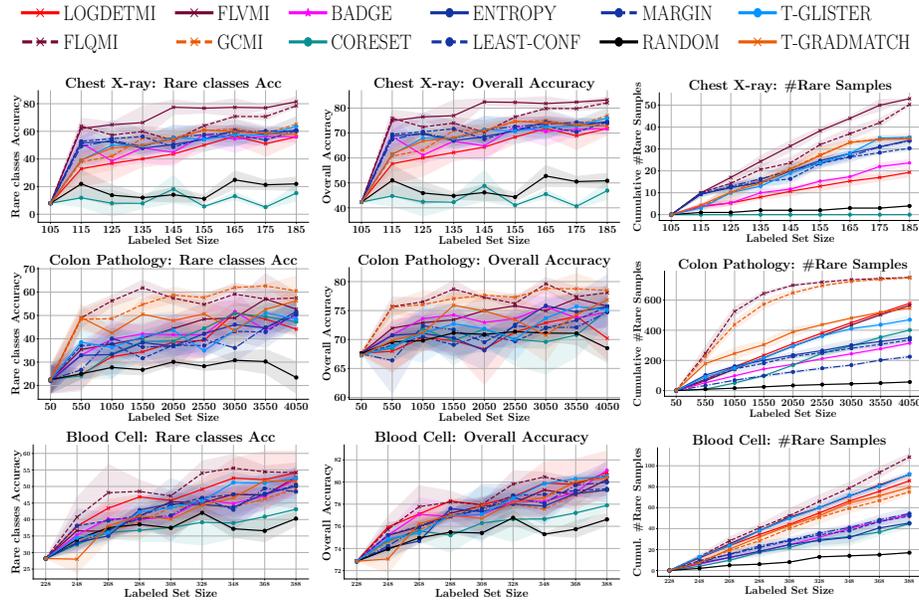


Fig. 3: AL for binary imbalanced medical image classification on Pneumonia-MNIST [12] (**first** row), Path-MNIST [11] (**second** row), and Blood-MNIST [1] (**third** row). CLINICAL outperforms the existing AL methods by  $\approx 2\% - 12\%$  on the rare classes acc. (**left** col.) and  $\approx 2\% - 6\%$  on overall acc. (**center** col.). SMI functions select the most number of rare class samples (**right** col.)

class as the rare class. For Path-MNIST,  $|\mathcal{C}_L| + |\mathcal{D}_L| = 3550$ ,  $|\mathcal{C}_U| + |\mathcal{D}_U| = 56.8K$ ,  $B = 500$  and,  $|\mathcal{T}| = 20$ . Following the natural class imbalance, we use two classes from the dataset (‘mucus’, ‘normal colon mucosa’) as rare classes. For Blood-MNIST,  $|\mathcal{C}_L| + |\mathcal{D}_L| = 228$ ,  $|\mathcal{C}_U| + |\mathcal{D}_U| = 1824$ ,  $B = 20$  and,  $|\mathcal{T}| = 20$ . Following the natural class imbalance, we use four classes from the dataset (‘basophil’, ‘eosinophil’, ‘lymphocyte’, ‘neutrophil’) as rare classes.

**Results:** The results for the binary imbalance scenario are shown in Fig. 3. We observe that the CLINICAL consistently outperform other methods by  $\approx 2\% - 12\%$  on the rare classes accuracy (Fig. 3(left column)) and  $\approx 2\% - 6\%$  on overall accuracy (Fig. 3(center column)). This is due to the fact that the SMI functions are able to select significantly more data points that belong to the rare classes (Fig. 3(right column)). Particularly, we observe that when the data modality is *X-ray* (Pneumonia-MNIST), the facility location based SMI variants, FLVMI and FLQMI perform significantly better than other acquisition functions due to their ability to model *representation*. For the *colon pathology* modality (Path-MNIST), GCM and FLQMI functions that model *query-relevance* significantly outperform other methods. Lastly, for the blood cell microscope modality (Blood-MNIST), we observe some improvement using FLQMI, although it selects many points from the rare classes.

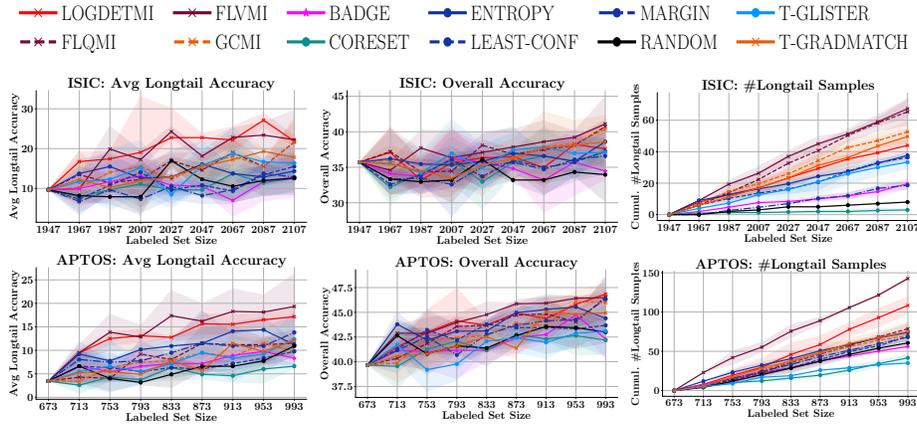


Fig. 4: Active learning for long-tail imbalanced medical image classification on ISIC-2018 [4] (**first row**) and APTOS-2019 [10] (**second row**). CLINICAL outperforms the state-of-the-art AL methods by  $\approx 10\% - 12\%$  on the average long-tail accuracy (**left col.**) and  $\approx 2\% - 5\%$  on overall accuracy (**center col.**). SMI functions select the most number of long-tail class samples (**right col.**)

## 4.2 Long-tail Imbalance

**Datasets:** We apply CLINICAL to two datasets that naturally show a long-tail distribution: 1) The ISIC-2018 skin lesion diagnosis dataset [4] and 2) APTOS-2019 [10] for diabetic retinopathy (DR) grading from retinal fundus images. We evaluate all AL methods on a balanced test set to obtain a fair estimate of accuracy across all classes. We split the remaining data randomly with 20% into the initial labeled set  $\mathcal{L}$  and 80% into the unlabeled set  $\mathcal{U}$ . We use a small held-out target set  $\mathcal{T}$  with data points from the classes at the tail of the distribution (long-tail classes, see Fig. 1). For ISIC-2018, we use the bottom three infrequent skin lesions from the tail of the distribution as long-tail classes (‘bowen’s disease’, ‘vascular lesions’, and ‘dermatofibroma’). We set  $B = 40$  and  $|\mathcal{T}| = 15$ . For APTOS-2019 we use the bottom two infrequent DR gradations as long-tail classes (‘severe DR’ and ‘proliferative DR’) (see Fig. 1). We set  $B = 20$  and  $|\mathcal{T}| = 10$ .

**Results:** We present the results for the long-tail imbalance scenario in Fig. 4. We observe that CLINICAL consistently outperforms other methods by  $\approx 10\% - 12\%$  on the average long-tail classes accuracy (Fig. 4(left column)) and  $\approx 2\% - 5\%$  on the overall accuracy (Fig. 4(center column)). This is because the SMI functions select significantly more data points from the long-tail classes (Fig. 4(right column)). On both datasets, we observe that the functions modeling query-relevance *and* diversity (FLVMI and LOGDETMI) outperform the functions modeling *only* query-relevance (FLQMI and GCMi). In addition to the ISIC-2018 dataset, we conduct additional class imbalance experiments on the Derma-MNIST [27] dataset (see

Appendix. D.1), and observe that FLVMI and LOGDETM perform significantly better than other functions on the dermatoscopy modality.

## 5 Conclusion

We demonstrate the effectiveness of CLINICAL for a wide range of medical data modalities for binary and long-tail imbalance. We empirically observe that the current methods in active learning cannot be directly applied to medical datasets with rare classes, and show that a targeting mechanism like SMI can greatly improve the performance on rare classes.

## Bibliography

- [1] Acevedo, A., Merino, A., Alférez, S., Molina, Á., Boldú, L., Rodellar, J.: A dataset of microscopic peripheral blood cell images for development of automatic recognition systems. *Data in Brief*, ISSN: 23523409, Vol. 30,(2020) (2020)
- [2] Arthur, D., Vassilvitskii, S.: k-means++: the advantages of careful seeding. In: *SODA '07: Proceedings of the eighteenth annual ACM-SIAM symposium on Discrete algorithms*. pp. 1027–1035. Society for Industrial and Applied Mathematics, Philadelphia, PA, USA (2007)
- [3] Ash, J.T., Zhang, C., Krishnamurthy, A., Langford, J., Agarwal, A.: Deep batch active learning by diverse, uncertain gradient lower bounds. In: *ICLR (2020)* (2020)
- [4] Codella, N., Rotemberg, V., Tschandl, P., Celebi, M.E., Dusza, S., Gutman, D., Helba, B., Kalloo, A., Liopyris, K., Marchetti, M., et al.: Skin lesion analysis toward melanoma detection 2018: A challenge hosted by the international skin imaging collaboration (isic). *arXiv preprint arXiv:1902.03368* (2019)
- [5] Fujishige, S.: *Submodular functions and optimization*. Elsevier (2005)
- [6] Gupta, A., Levin, R.: The online submodular cover problem. In: *ACM-SIAM Symposium on Discrete Algorithms (2020)* (2020)
- [7] He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. pp. 770–778 (2016)
- [8] Iyer, R., Khargoankar, N., Bilmes, J., Asnani, H.: Submodular combinatorial information measures with applications in machine learning. *arXiv preprint arXiv:2006.15412* (2020)
- [9] Iyer, R.K.: *Submodular optimization and machine learning: Theoretical results, unifying and scalable algorithms, and applications*. Ph.D. thesis (2015)
- [10] Kaggle: Aptos 2019 blindness detection (2019), <https://www.kaggle.com/c/aptos2019-blindness-detection/data>
- [11] Kather, J.N., Krisam, J., Charoentong, P., Luedde, T., Herpel, E., Weis, C.A., Gaiser, T., Marx, A., Valous, N.A., Ferber, D., et al.: Predicting survival from colorectal cancer histology slides using deep learning: A retrospective multicenter study. *PLoS medicine* **16**(1), e1002730 (2019)
- [12] Kermany, D.S., Goldbaum, M., Cai, W., Valentim, C.C., Liang, H., Baxter, S.L., McKeown, A., Yang, G., Wu, X., Yan, F., et al.: Identifying medical diagnoses and treatable diseases by image-based deep learning. *Cell* **172**(5), 1122–1131 (2018)
- [13] Killamsetty, K., Durga, S., Ramakrishnan, G., De, A., Iyer, R.: Grad-match: Gradient matching based data subset selection for efficient deep model training. In: *International Conference on Machine Learning*. pp. 5464–5474. PMLR (2021)

- [14] Killamsetty, K., Sivasubramanian, D., Ramakrishnan, G., Iyer, R.: Glistler: Generalization based data subset selection for efficient and robust learning. In *AAAI* (2021)
- [15] Kirsch, A., Van Amersfoort, J., Gal, Y.: Batchbald: Efficient and diverse batch acquisition for deep bayesian active learning. *arXiv preprint arXiv:1906.08158* (2019)
- [16] Kothawade, S., Beck, N., Killamsetty, K., Iyer, R.: Similar: Submodular information measures based active learning in realistic scenarios. *Advances in Neural Information Processing Systems* **34** (2021)
- [17] Kothawade, S., Ghosh, S., Shekhar, S., Xiang, Y., Iyer, R.: Talisman: Targeted active learning for object detection with rare classes and slices using submodular mutual information. *arXiv preprint arXiv:2112.00166* (2021)
- [18] Kothawade, S., Kaushal, V., Ramakrishnan, G., Bilmes, J., Iyer, R.: Prism: A rich class of parameterized submodular information measures for guided subset selection. *arXiv preprint arXiv:2103.00128* (2021)
- [19] Kothiyari, M., Mekala, A.R., Iyer, R., Ramakrishnan, G., Jyothi, P.: Personalizing asr with limited data using targeted subset selection. *arXiv preprint arXiv:2110.04908* (2021)
- [20] Li, J., Li, L., Li, T.: Multi-document summarization via submodularity. *Applied Intelligence* **37**(3), 420–430 (2012)
- [21] Lin, H.: Submodularity in natural language processing: algorithms and applications. Ph.D. thesis (2012)
- [22] Loshchilov, I., Hutter, F.: Sgdr: Stochastic gradient descent with warm restarts. *arXiv preprint arXiv:1608.03983* (2016)
- [23] Mirzasoleiman, B., Badanidiyuru, A., Karbasi, A., Vondrák, J., Krause, A.: Lazier than lazy greedy. In: *Proceedings of the AAAI Conference on Artificial Intelligence*. vol. 29 (2015)
- [24] Roth, D., Small, K.: Margin-based active learning for structured output spaces. In: *European Conference on Machine Learning*. pp. 413–424. Springer (2006)
- [25] Sener, O., Savarese, S.: Active learning for convolutional neural networks: A core-set approach. In: *International Conference on Learning Representations* (2018)
- [26] Settles, B.: Active learning literature survey. Tech. rep., University of Wisconsin-Madison Department of Computer Sciences (2009)
- [27] Tschandl, P., Rosendahl, C., Kittler, H.: The ham10000 dataset, a large collection of multi-source dermatoscopic images of common pigmented skin lesions. *Scientific data* **5**(1), 1–9 (2018)
- [28] Vasudevan, A.B., Gygli, M., Volokitin, A., Van Gool, L.: Query-adaptive video summarization via quality-aware relevance estimation. In: *Proceedings of the 25th ACM international conference on Multimedia*. pp. 582–590 (2017)
- [29] Wang, D., Shang, Y.: A new active labeling method for deep learning. In: *2014 International joint conference on neural networks (IJCNN)*. pp. 112–119. IEEE (2014)
- [30] Yang, J., Shi, R., Wei, D., Liu, Z., Zhao, L., Ke, B., Pfister, H., Ni, B.: Medmnist v2: A large-scale lightweight benchmark for 2d and 3d biomedical image classification. *arXiv preprint arXiv:2008* (2021)

## Supplementary Material for Targeted Active Learning for Imbalanced Medical Image Classification

### A Summary of Notations

Topic	Notation	Explanation
CLINICAL (Sec. 3)	$\mathcal{U}$	Unlabeled set of $ \mathcal{U} $ instances
	$\mathcal{A}$	A subset of $\mathcal{U}$
	$s_{ij}$	Similarity between any two data points $i$ and $j$
	$f$	A submodular function
	$\mathcal{L}$	Labeled set of data points
	$\mathcal{Q}$	Query set
	$\mathcal{M}$	Deep model
	$B$	Active learning selection budget
	$N$	Number of selection rounds in active learning
	$\mathcal{T}$	Held-out target set containing data points from the rare classes
	$\hat{\mathcal{T}}$	Subset of $\mathcal{T}$ containing only the misclassified data points
	$\mathcal{H}$	Loss function used to train model $\mathcal{M}$
	$X$	Pairwise similarity matrix computed using gradients
Experiments (Sec. 4)	$\mathcal{C}_{\mathcal{L}}$	Rare classes data points in the labeled set $\mathcal{L}$
	$\mathcal{C}_{\mathcal{U}}$	Rare classes data points in the unlabeled set $\mathcal{U}$
	$\mathcal{D}_{\mathcal{L}}$	Non-rare (frequent) classes data points in the labeled set $\mathcal{L}$
	$\mathcal{D}_{\mathcal{U}}$	Non-rare (frequent) classes data points in the unlabeled set $\mathcal{U}$

Table 1: Summary of notations used throughout this paper

### B Details of Datasets used

#### B.1 ISIC [4]

- ISIC dataset is a representative collection of all important diagnostic categories in the realm of pigmented lesions
- It contains, 10015 28x28 color images in 7 different classes

- Classes represent various types of skin cancer diseases like Actinic keratoses and intraepithelial carcinoma / Bowen’s disease (akiec), basal cell carcinoma (bcc), benign keratosis-like lesions (solar lentiginos / seborrheic keratoses and lichen-planus like keratoses, bkl), dermatofibroma (df), melanoma (mel), melanocytic nevi (nv) and vascular lesions (angiomas, angiokeratomas, pyogenic granulomas and hemorrhage, vasc)

## B.2 APTOS [10]

- APTOS dataset contains retina images taken using fundus photography under a variety of imaging conditions
- It contains 3662 28x28 color images in 5 different classes
- Classes represent severity of diabetic retinopathy on a scale of 0 to 4

## B.3 PathMNIST [11]

- A dataset based on a prior study for predicting survival from colorectal cancer histology slides, which provides a dataset NCT-CRC-HE-100K of 100,000 non-overlapping image patches from hematoxylin and eosin stained histological images, and a test dataset CRC-VAL-HE-7K of 7,180 image patches from a different clinical center.
- There are 9 types of tissues are involved, resulting in a multi-class classification task.
- The images are resized from 3 x 224 x 224 into 3 x 28 x 28 as in [30].
- Classes represent 9 types of tissues: ‘adipose’, ‘background’, ‘debris’, ‘lymphocytes’, ‘mucus’, ‘smooth muscle’, ‘normal colon mucosa’, ‘cancer-associated stroma’, ‘colorectal adenocarcinoma epithelium’.

## B.4 PneumoniaMNIST [12]

- A dataset based on a prior dataset of 5,856 pediatric chest X-ray images.
- The task is binary-class classification of pneumonia and normal.
- We split the source training set with a ratio of 9:1 into training and validation set, and use its source validation set as the test set.
- The source images are single-channel, and their sizes range from (384-2,916) x (127-2,713). We use data from [30] where they center-crop the images and resize them into 1 x 28 x 28.

## B.5 BloodMNIST [1]

- BloodMNIST dataset is based on a dataset of individual normal blood cells, captured from individuals without any kind of infection or disease
- It contains 17092 32x32 color images in 8 different classes
- Classes represent various types of blood cells (without any infection) like basophil, eosinophil, erythroblast, ig, lymphocyte, monocyte, neutrophil, platelet

### B.6 DermaMNIST [4]

- DermaMNIST dataset is a large collection of multi-source dermatoscopic images of common pigmented skin lesions
- It contains, 10015 32x32 color images in 7 different classes.
- Classes represent various types of skin lesions like actinic keratoses, basal cell carcinoma, benign keratosis-like lesions, dermatofibroma, melanoma, melanocytic nevi, vascular lesions

## C Scalability of CLINICAL

Below, we provide a detailed analysis of the complexity of creating and optimizing the different SMI functions. Denote  $|\mathcal{X}|$  as the size of set  $\mathcal{X}$ . Also, let  $|\mathcal{U}| = n$  (the ground set size, which is the size of the unlabeled set in this case).

- **Facility Location:** We start with FLVMI. The complexity of creating the kernel matrix is  $O(n^2)$ . The complexity of optimizing it is  $\tilde{O}(n^2)$  (using memoization)<sup>3</sup> if we use the stochastic greedy algorithm [23] and  $O(n^2k)$  with the naive greedy algorithm. The overall complexity is  $\tilde{O}(n^2)$ . For FLQMI, the cost of creating the kernel matrix is  $O(n|\mathcal{Q}|)$ , and the cost of optimization is also  $\tilde{O}(n|\mathcal{Q}|)$  (with naive greedy, it is  $O(nB|\mathcal{Q}|)$ ).
- **Log-Determinant:** For LogDetMI, the complexity of the kernel matrix computation (and storage) is  $O(n^2)$ . The complexity of optimizing the LogDet function using the stochastic greedy algorithm is  $\tilde{O}(B^2n)$ , so the overall complexity is  $\tilde{O}(n^2 + B^2n)$ .
- **Graph-Cut:** For GCMI, we require a  $O(n|\mathcal{Q}|)$  kernel matrix, and the complexity of the stochastic greedy algorithm is also  $\tilde{O}(n|\mathcal{Q}|)$ .

We end with a few comments. First, most of the complexity analysis above is with the stochastic greedy algorithm [23]. If we use the naive or lazy greedy algorithm, the worst-case complexity is a factor  $B$  larger. Secondly, we ignore log-factors in the complexity of stochastic greedy since the complexity is actually  $O(n \log 1/\epsilon)$ , which achieves an  $1 - 1/e - \epsilon$  approximation.

## D Additional Results

### D.1 Dermatoscope Binary Imbalance Results

We present the results for binary imbalance on the DermaMNIST [4] dataset in Fig. 5.

---

<sup>3</sup>  $\tilde{O}$ : Ignoring log-factors

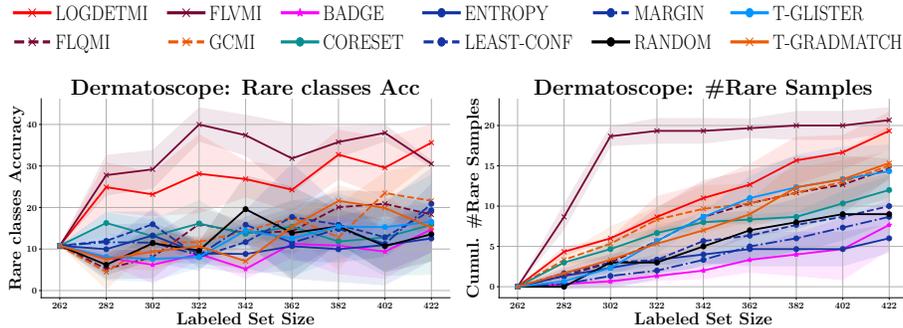


Fig. 5: Dermatoscope Binary Imbalance Results. We observe that FLVMI and LOGDETM significantly outperform other baselines. Particularly, FLVMI selects significantly more rare class samples than all methods.

## D.2 Statistical Significance Penalty Matrices

The penalty matrices computed in this paper follow the strategy used in [3]. In their strategy, a penalty matrix is constructed for each dataset-model pair. Each cell  $(i, j)$  of the matrix reflects the fraction of training rounds that AL with selection algorithm  $i$  has higher test accuracy than AL with selection algorithm,  $j$  with statistical significance. As such, the average difference between the test accuracies of  $i$  and  $j$  and the standard error of that difference are computed for each training round. A two-tailed  $t$ -test is then performed for each training round: If  $t > t_\alpha$ , then  $\frac{1}{N_{train}}$  is added to cell  $(i, j)$ . If  $t < -t_\alpha$ , then  $\frac{1}{N_{train}}$  is added to cell  $(j, i)$ . Hence, the full penalty matrix gives a holistic understanding of how each selection algorithm compares against the others: A row with mostly high values signals that the associated selection algorithm performs better than the others; however, a column with mostly high values signals that the associated selection algorithm performs worse than the others. As a final note, [3] takes an additional step where they consolidate the matrices for each dataset-model pair into one matrix by taking the sum across these matrices, giving a summary of the AL performance for their entire paper that is fairly weighted to each experiment. Below, we present the penalty matrices for each of the settings.

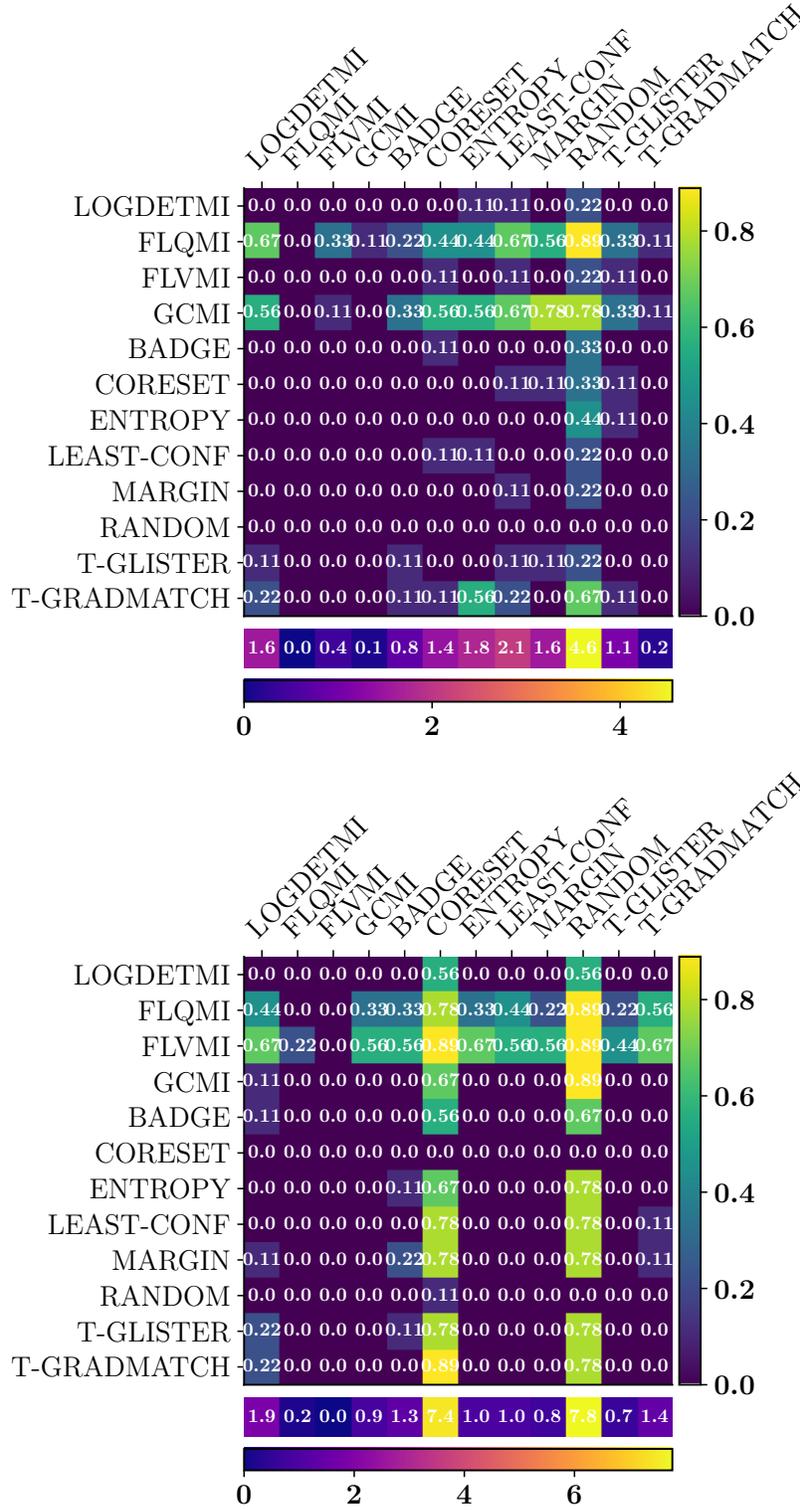


Fig. 6: Penalty Matrix comparing the average accuracy of binary imbalance for Path-MNIST (**top**) datasets and Pneumonia-MNIST (**bottom**) using targeted active learning across multiple runs. We observe that the SMI functions have a much lower column sum compared to other approaches.

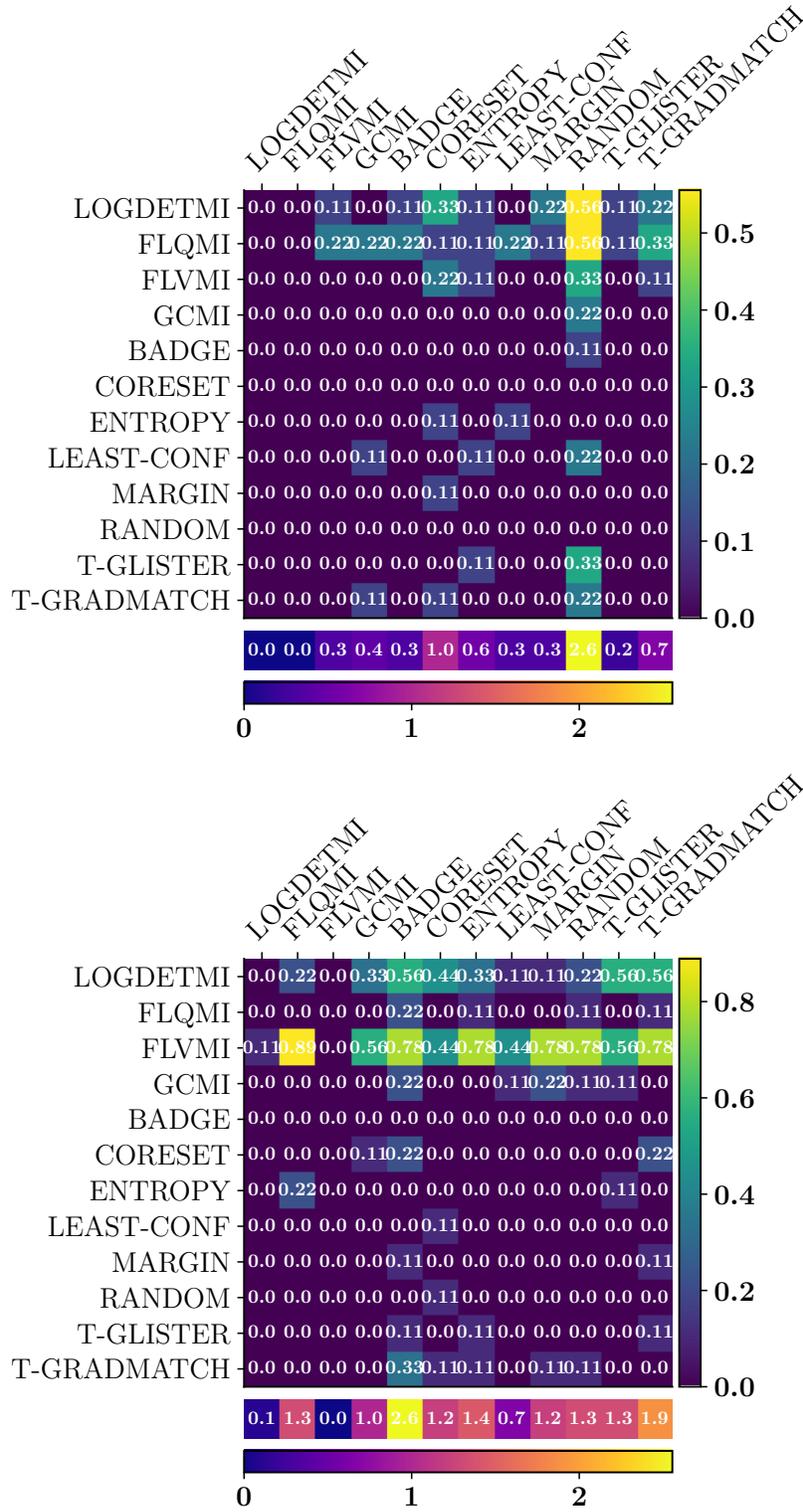


Fig. 7: Penalty Matrix comparing the average accuracy for rare classes in the binary imbalance for Blood-MNIST (top) and Derma-MNIST (bottom) datasets using targeted active learning across multiple runs. We observe that the SMI functions have a much lower column sum compared to other approaches.

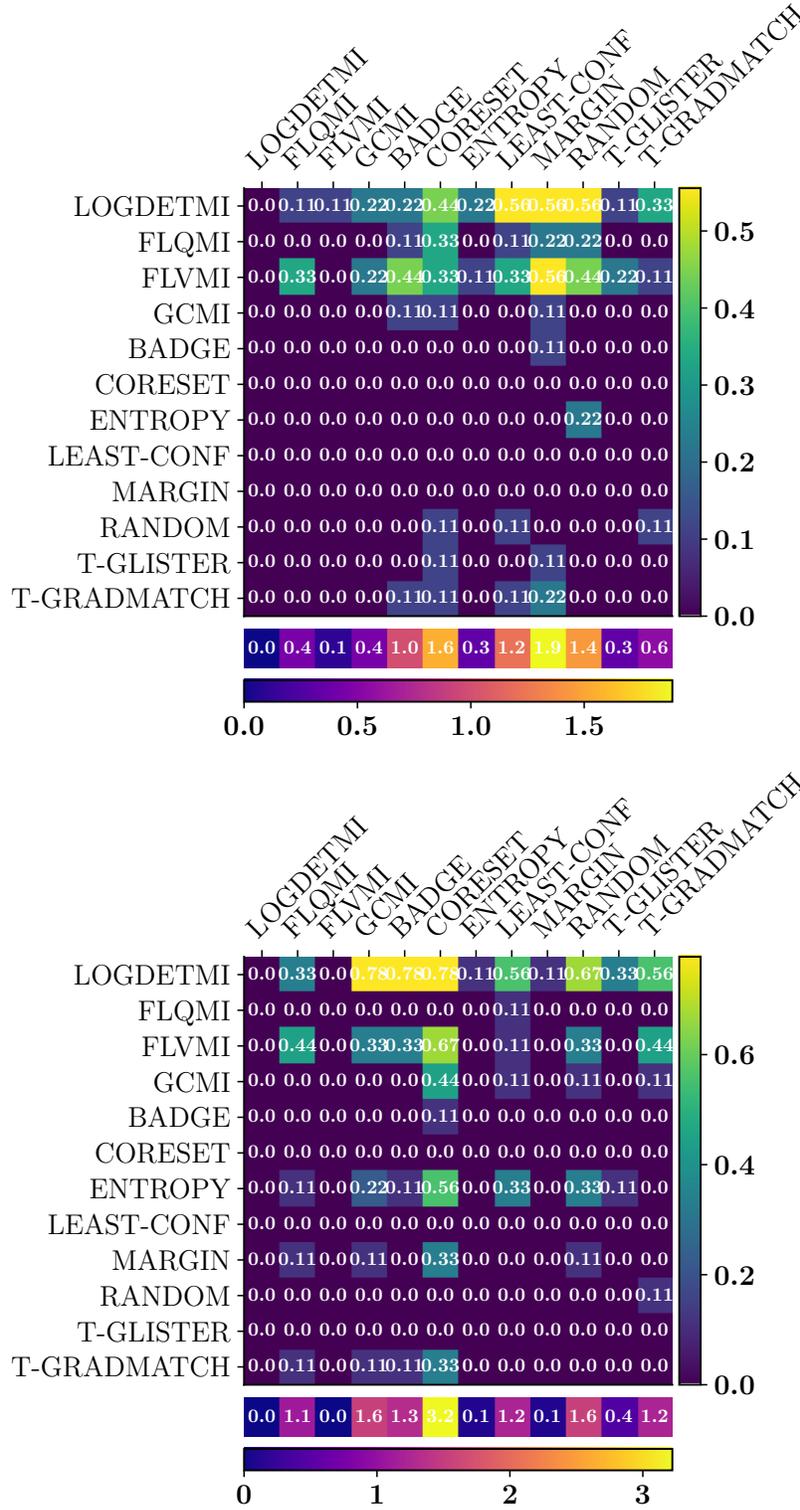


Fig. 8: Penalty Matrix comparing the average accuracy of rare classes in the long-tail imbalance for ISIC-2018 (top) and APTOS (bottom) datasets using targeted active learning across multiple runs. We observe that the SMI functions have a much lower column sum compared to other approaches.