

# Scalable Batch-Mode Deep Bayesian Active Learning via Equivalence Class Annealing

Renyu Zhang, Aly A. Khan, Robert L. Grossman, Yuxin Chen



## Summary

- ▶ **BALanCe**: an efficient deep Bayesian active learning framework motivated by a decision-theoretic selection criterion
- ▶ **Batch-BALanCe**: batch-mode version

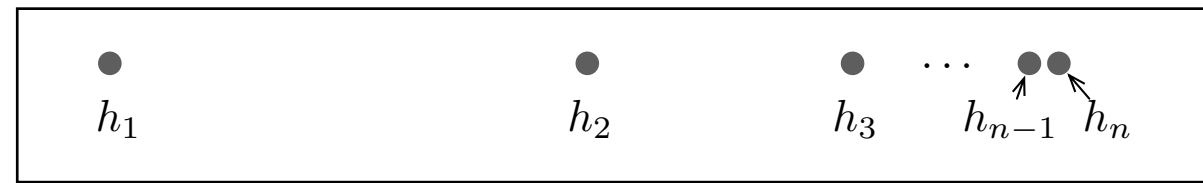
## Motivation

### Numerical example

The Most Information Selection criterion uses mutual information between the predicted label and model parameters as the acquisition function:

$$\Delta_{\text{BALD}}(x | \mathcal{D}_{\text{train}}) \triangleq \mathbb{I}(y; \omega | x, \mathcal{D}_{\text{train}})$$

MIS/BALD can be ineffective



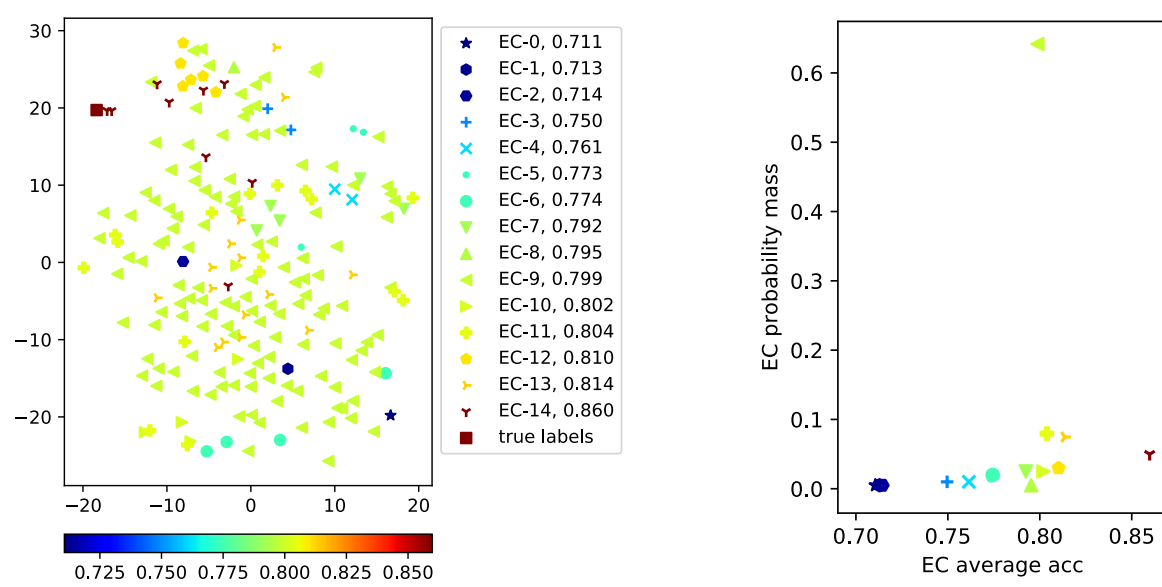
The hypothesis class  $\mathcal{H} = \{h_1, \dots, h_n\}$  is structured such that

$$d_{\text{H}}(h_i, h_j) = \begin{cases} 2^{1-i} - 2^{1-j} & \text{if } i < j, \\ 2^{1-j} - 2^{1-i} & \text{o.w.} \end{cases}$$

- ▶ MIS/BALD on average require  $O(\log n)$ .

- ▶ A “smarter” policy could query examples to sequentially check the consistency of  $h_1, h_2, \dots, h_n$  until all remaining hypotheses are within distance  $\sigma$ . It requires  $\log(1/\sigma)$ .

### Empirical validation



Samples from posterior BNN via MC dropout; embedding is generated by applying t-SNE on the disagreement between hypotheses; Colorbar indicates the (approximate) test accuracy of the sampled neural networks on the MNIST dataset.

## Problem statement

### Assume

- ▶ labeling each query  $x$  incurs a unit cost.
- ▶ Labeled dataset  $\mathcal{D}_{\text{train}}$
- ▶ Unlabelled dataset  $\mathcal{D}_{\text{pool}}$  drawn i.i.d. from some underlying data distribution
- ▶ A set of hypotheses  $\mathcal{H} = \{h_1, \dots, h_n\}$ ;
- ▶ In this work, we consider BNN hypothesis class with parameters  $\omega \sim p(\omega | \mathcal{D}_{\text{train}})$ .

### Want

an active learning algorithm that allows us to find a hypotheses with a target error rate  $\sigma \in [0, 1]$  while minimizing the total query cost.

### Acquisition function

$$\Delta_{\text{BALanCe}}(x_{1:b} | \mathcal{D}_{\text{train}}) \approx \sum_{\hat{y}_{1:b}} \left( \sum_{k=1}^K \frac{p(\hat{y}_{1:b} | \hat{\omega}_k) + p(\hat{y}_{1:b} | \hat{\omega}'_k)}{2K} \right) \cdot \left[ \sum_{k=1}^K \frac{\mathbb{1}_{d_{\text{H}}(\hat{\omega}_k, \hat{\omega}'_k) > \tau} \cdot (1 - \lambda_{\hat{\omega}_k, \hat{y}_{1:b}} \lambda_{\hat{\omega}'_k, \hat{y}_{1:b}})}{K} \right]$$

Threshold adaptively annealed based on validation prediction

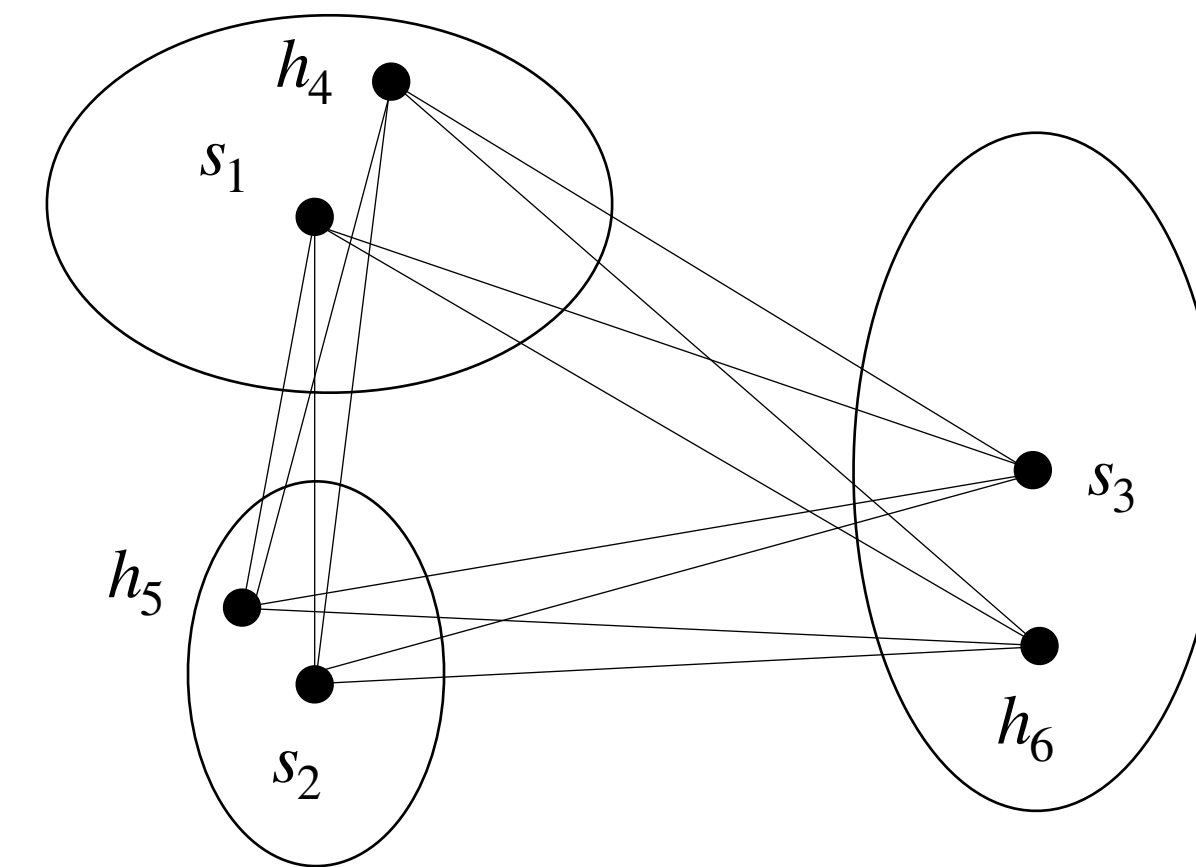
### New metric for clustering

$$\mathbb{I}_{\Delta_{\text{BALanCe}}}(x_1, x_2) = \Delta_{\text{BALanCe}}(x_1) + \Delta_{\text{BALanCe}}(x_2) - \Delta_{\text{BALanCe}}(\{x_1, x_2\})$$

## Equivalence class

### Definition 1.

Let  $(\mathcal{H}, d)$  be a metric space where  $\mathcal{H}$  is a hypothesis class and  $d$  is a metric. For a given set  $\mathcal{V} \subseteq \mathcal{H}$  and centers  $\mathcal{S} = \{s_1, \dots, s_k\} \subseteq \mathcal{V}$  of size  $k$ , let  $r^{\mathcal{S}}: \mathcal{V} \rightarrow [k]$  be a partition function over  $\mathcal{V}$ , and  $\mathcal{D}_i := \{h \in \mathcal{V} | r^{\mathcal{S}}(h) = i\}$ , such that  $\forall i, j \in [k], r^{\mathcal{S}}(s_j) = i$  and  $\forall h \in \mathcal{D}_i, d(h, s_i) \leq d(h, s_j)$ . Each  $\mathcal{D}_i \subseteq \mathcal{V}$  is called an equivalence class induced by  $s_i \in \mathcal{S}$ .



## Two selection strategies

### Algorithm I: Greedy selection

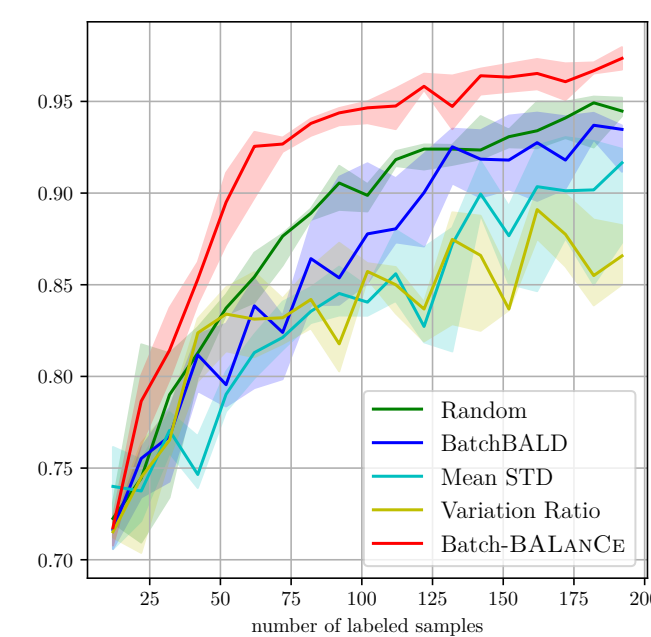
- 1 **input:**  $\mathcal{D}_{\text{pool}}, \tilde{\mathcal{D}}_{\text{pool}}$ : a trained BNN, and threshold  $\tau, \{\hat{\omega}_k, \hat{\omega}'_k\}_{k=1}^K, B$
- 2  $\mathcal{A}_0 = \emptyset$
- 3 **for**  $b \in [B]$  **do**
- 4 **for**  $x \in \mathcal{D}_{\text{pool}} \setminus \mathcal{A}_{b-1}$  **do**
- 5  $s_x \leftarrow \Delta_{\text{BALanCe}}(\mathcal{A}_{b-1} \cup \{x\})$
- 6  $x_b \leftarrow \text{argmax}_{x \in \mathcal{D}_{\text{pool}} \setminus \mathcal{A}_{b-1}} s_x$
- 7  $\mathcal{A}_b \leftarrow \mathcal{A}_{b-1} \cup \{x_b\}$
- 8 **output:** batch  $\mathcal{A}_B$

For batch size  $< 40$

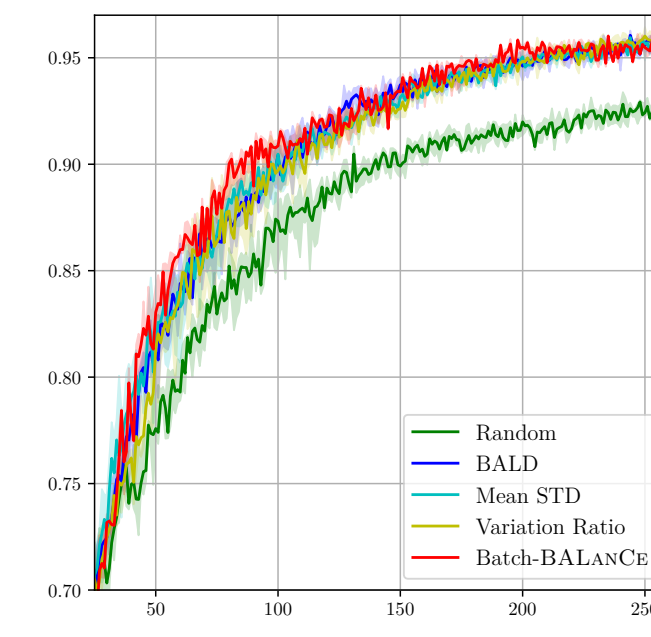
### Algorithm II: BALanCe-Clustering

- 1 **input:** subset  $\mathcal{C} \subseteq \mathcal{D}_{\text{pool}}, \tilde{\mathcal{D}}_{\text{pool}}$ ; threshold  $\tau, \{\hat{\omega}_k, \hat{\omega}'_k\}_{k=1}^K$ , coldness parameter  $\beta$ , and  $B$
- 2 sample initial centroids  $\mathcal{O} = \{\mu_j\}_{j=1}^B \subseteq \mathcal{C}$  with  $p(x) \sim \Delta_{\text{BALanCe}}(x)^\beta$
- 3 **while**  $\mathcal{O}$  not converged **do**
- 4 **for**  $x \in \mathcal{C}$  **do**
- 5  $a_x \leftarrow \text{argmax}_j \mathbb{I}_{\Delta_{\text{BALanCe}}}(x, \mu_j)$
- 6  $\mathcal{S}_j \leftarrow \{x \in \mathcal{C} : a_x = j\}$
- 7 **for**  $j \in [B]$  **do**
- 8  $\mu_j \leftarrow \text{argmax}_{x_2 \in \mathcal{S}_j} \sum_{x_1 \in \mathcal{S}_j} \mathbb{I}_{\Delta_{\text{BALanCe}}}(x_1, x_2)$
- 9 **output:**  $\mathcal{S}_{1:B}, \mu_{1:B}$

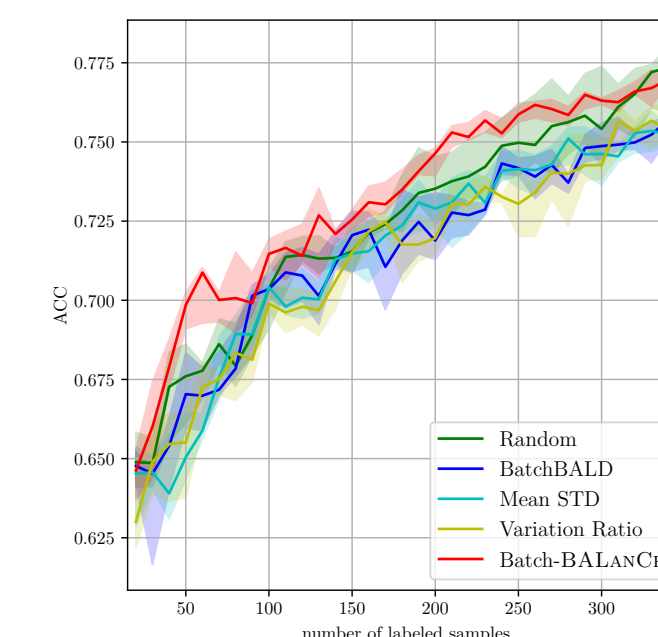
For batch size  $\geq 40$



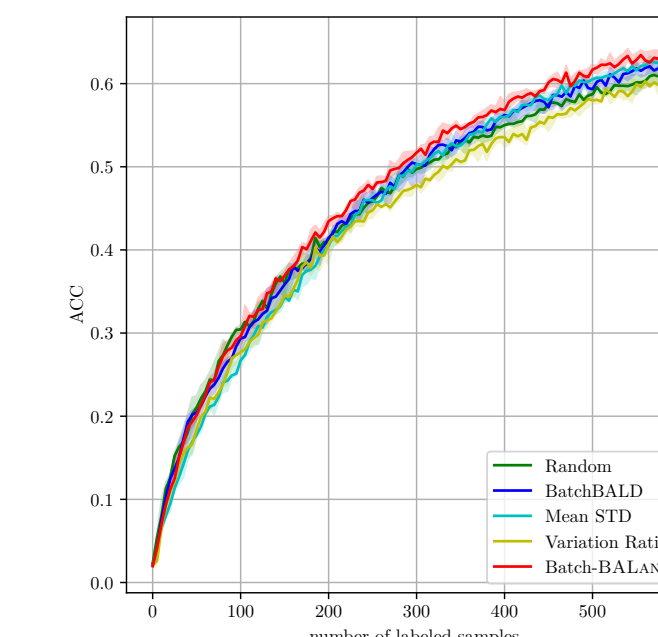
DRIFT; MC-dropout



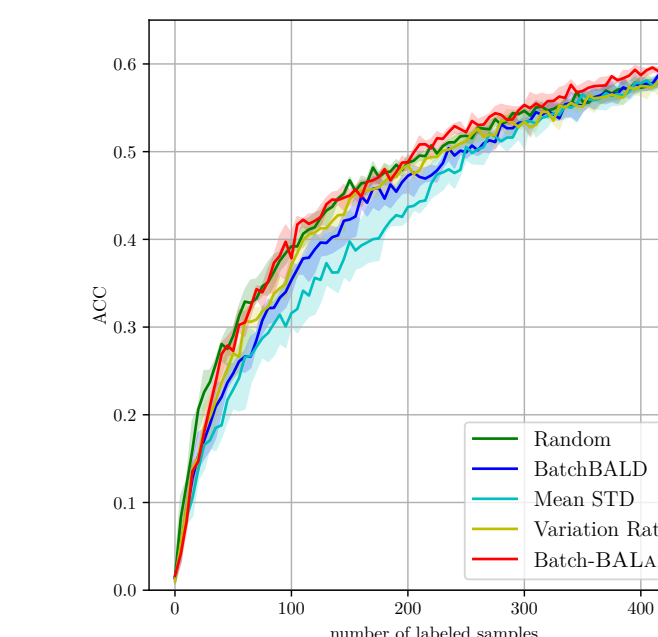
MNIST; MC-dropout



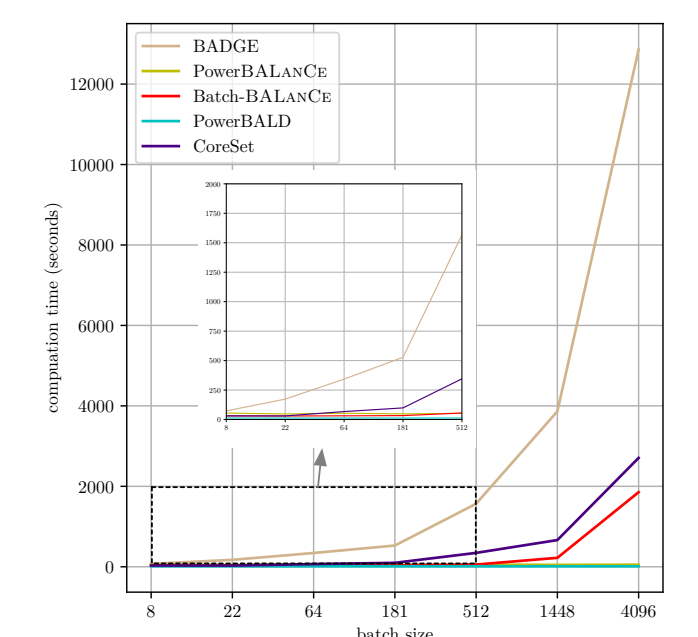
Fashion-MNIST; MC-dropout



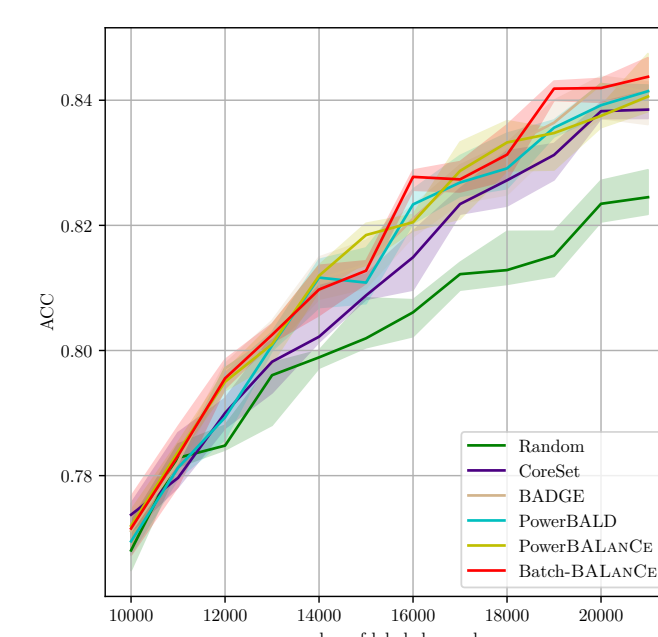
EMNIST-Balanced; MC-dropout



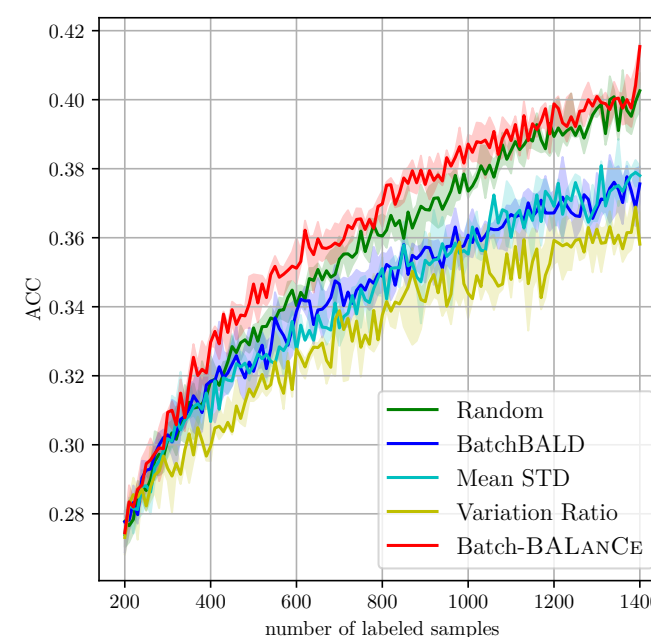
EMNIST-ByClass; MC-dropout



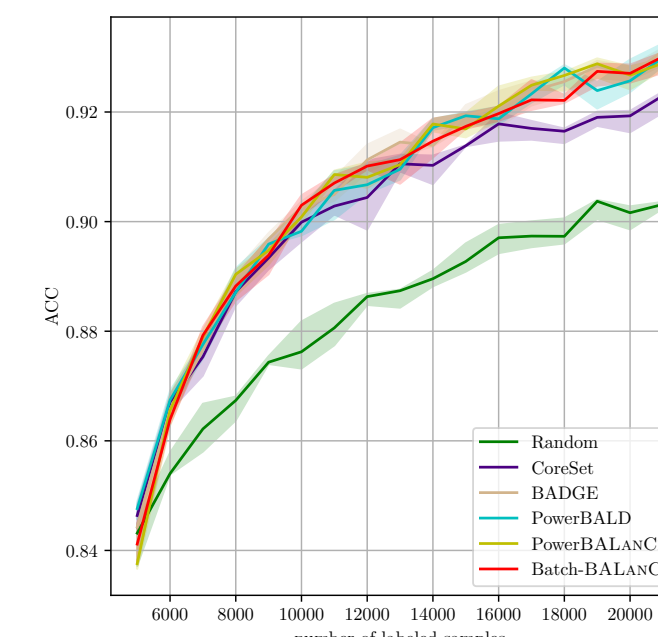
Computation time vs. batch size



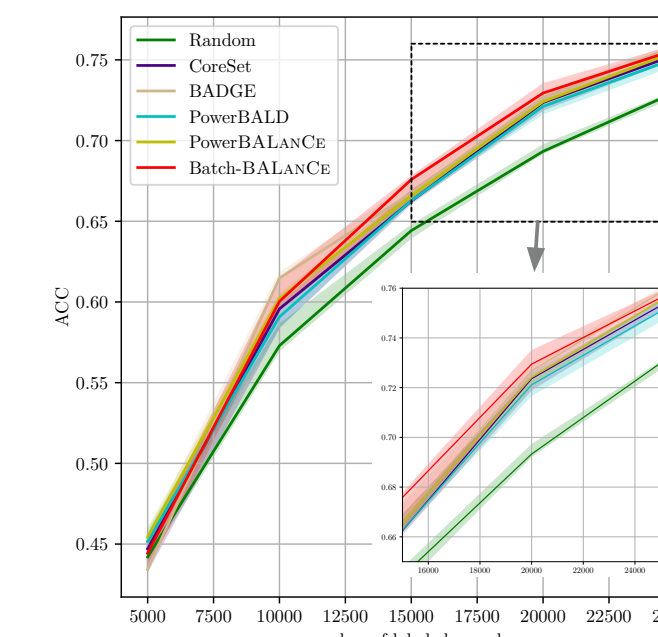
CIFAR-10; MC-dropout



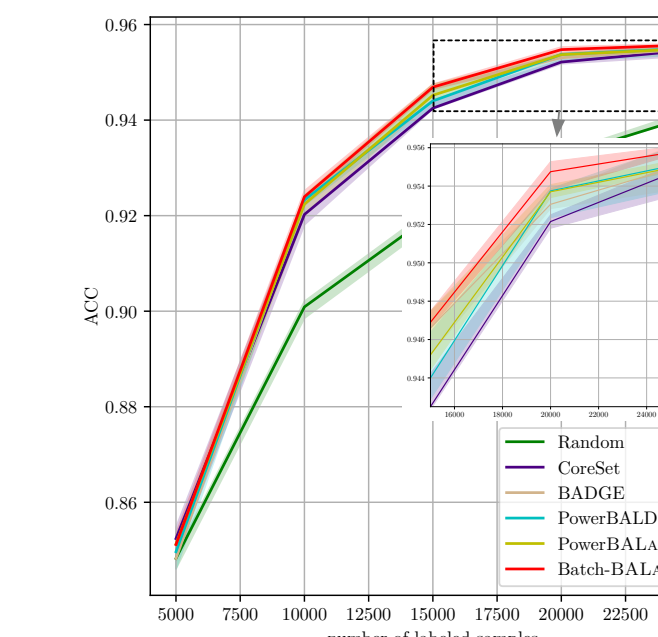
CINIC; MC-dropout



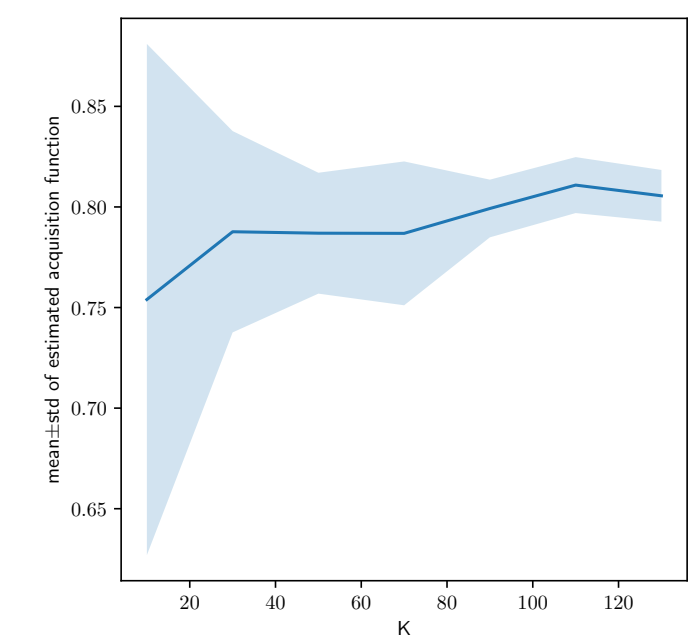
SVHN; MC-dropout



CIFAR-100; cSG-MCMC



CIFAR-10; cSG-MCMC



Acquisition function vs #MC samples