

# Deep Bayesian Active Learning via Equivalence Class Annealing

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



## 1 Summary

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

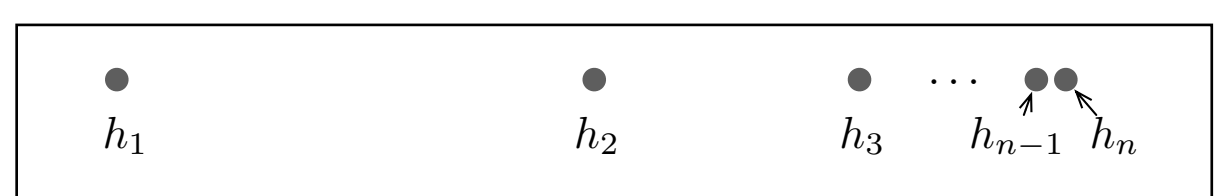
## 2 Motivation

### 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



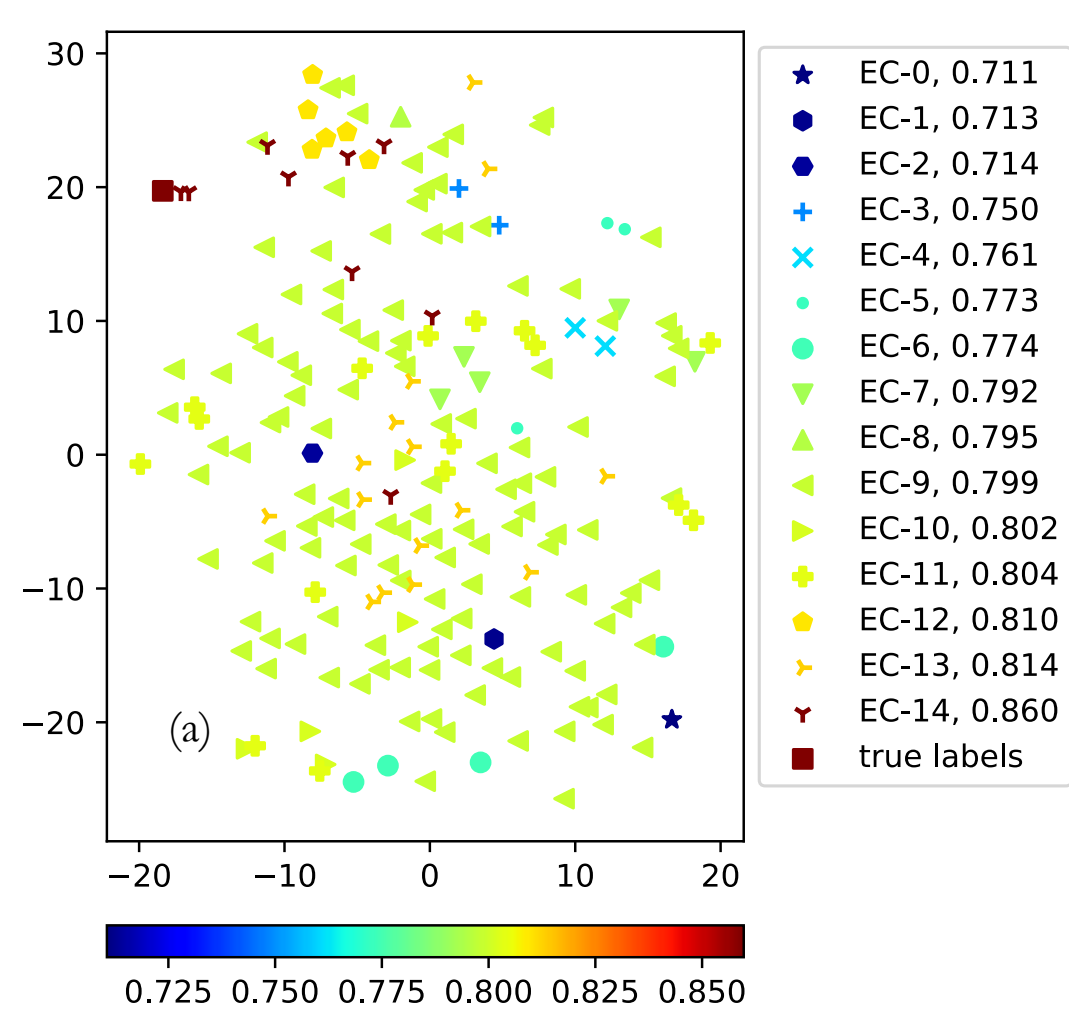
The hypothesis class  $\mathcal{H} = \{h_1, \dots, h_n\}$  is structured such that

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

- ▶ 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.

## 3 Notations

- ▶ Labeled dataset  $\mathcal{D}_{\text{train}}$
- ▶ Unlabeled 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  $\omega \sim p(\omega | \mathcal{D}_{\text{train}})$ .

## 4 Problem statement

### Assume

labeling each query  $x$  incurs a unit cost.

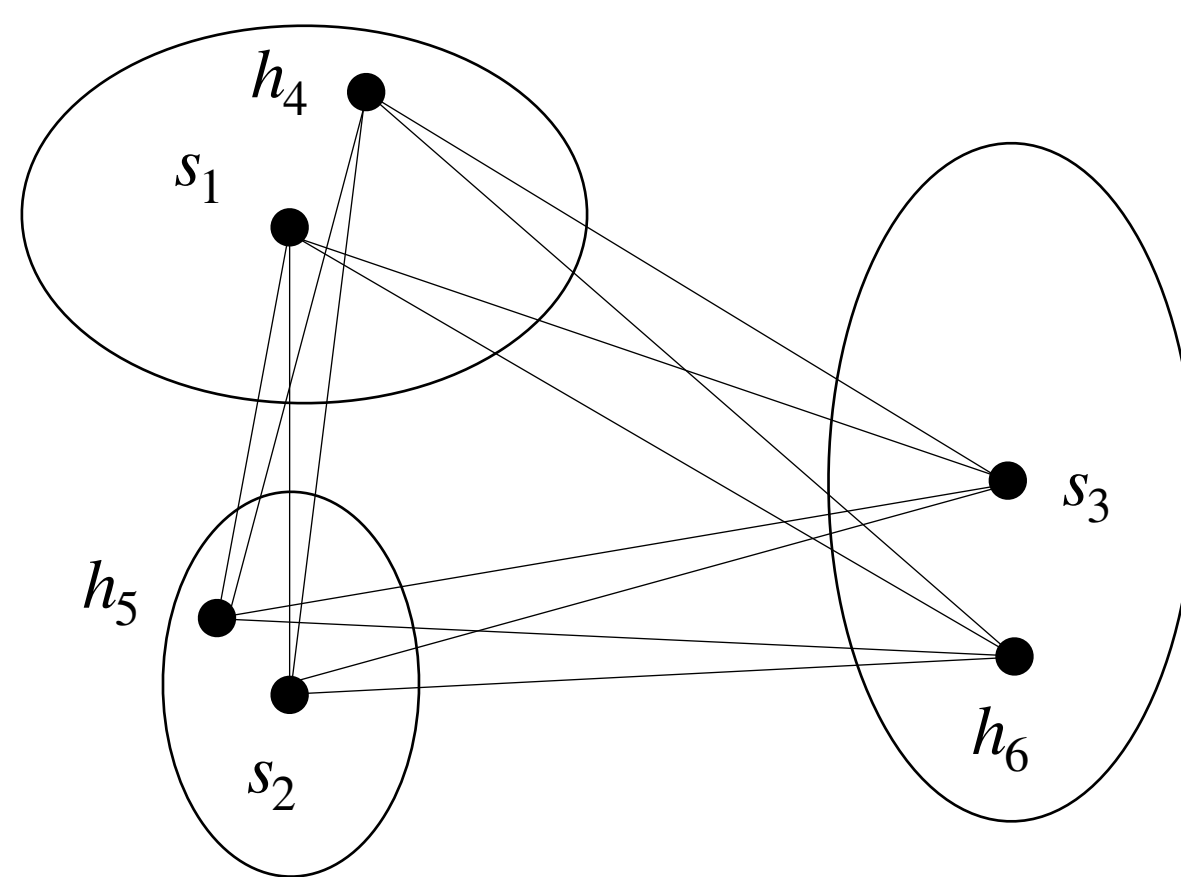
### Want

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.

## 5 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  $\forall i, j \in [k], r^{\mathcal{S}}(s_i) = i$  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  $s_i \in \mathcal{S}$ .



## 6 Our approach

$$\Delta_{\text{BALanCe}}(x | \mathcal{D}_{\text{train}}) \triangleq \mathbb{E}_{y \sim p(\omega, \omega' | \mathcal{D}_{\text{train}})} \left[ \mathbb{1}_{d_{\mathcal{H}}(\omega, \omega') > \tau} \cdot \text{WEIGHT-DISCOUNTED}(y | x, \mathcal{D}_{\text{train}}) \right]$$

Threshold adaptively annealed based on validation prediction

### Fully sequential setting:

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

### Batch mode setting:

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

## The BALanCe Algorithm

### Algorithm 1: BALanCe

- 1 **input:**  $\mathcal{D}_{\text{pool}}, \mathcal{D}_{\text{pool}}$ , a trained BNN and threshold  $\tau$
- 2 draw  $K$  pairs of MC dropout samples  $\{\hat{\omega}_k, \hat{\omega}'_k\}_{k=1}^K$  from the trained BNN
- 3 **for**  $k \in [K]$  **do**
- 4 Calculate indicator function  $\mathbb{1}_{d_{\mathcal{H}}(\hat{\omega}_k, \hat{\omega}'_k) > \tau}$  with prediction on  $\mathcal{D}_{\text{pool}}$
- 5 **for** each  $x \in \mathcal{D}_{\text{pool}}$  **do**
- 6  $\Delta(x) = 0$
- 7 **for** each  $\hat{y}$  **do**
- 8  $p(\hat{y}) = \frac{1}{2K} \sum_{k=1}^K (p(\hat{y} | \hat{\omega}_k) + p(\hat{y} | \hat{\omega}'_k))$
- 9  $\delta(x, \hat{y}) = 0$
- 10 **for**  $k \in [K]$  **do** *Only consider those pairs that have large distances*
- 11 **if**  $\mathbb{1}_{d_{\mathcal{H}}(\hat{\omega}_k, \hat{\omega}'_k) > \tau} = 1$  **then**
- 12  $\delta(x, \hat{y}) \leftarrow \delta(x, \hat{y}) + \frac{1}{K} (1 - \lambda_{\hat{\omega}_k, \hat{y}} \lambda_{\hat{\omega}'_k, \hat{y}})$
- 13  $\Delta(x) \leftarrow \Delta(x) + p(\hat{y}) \cdot \delta(x, \hat{y})$
- 14  $x \leftarrow \text{argmax}_{x \in \mathcal{D}_{\text{pool}}} \Delta(x)$
- 15 **output:**  $\{x\}$

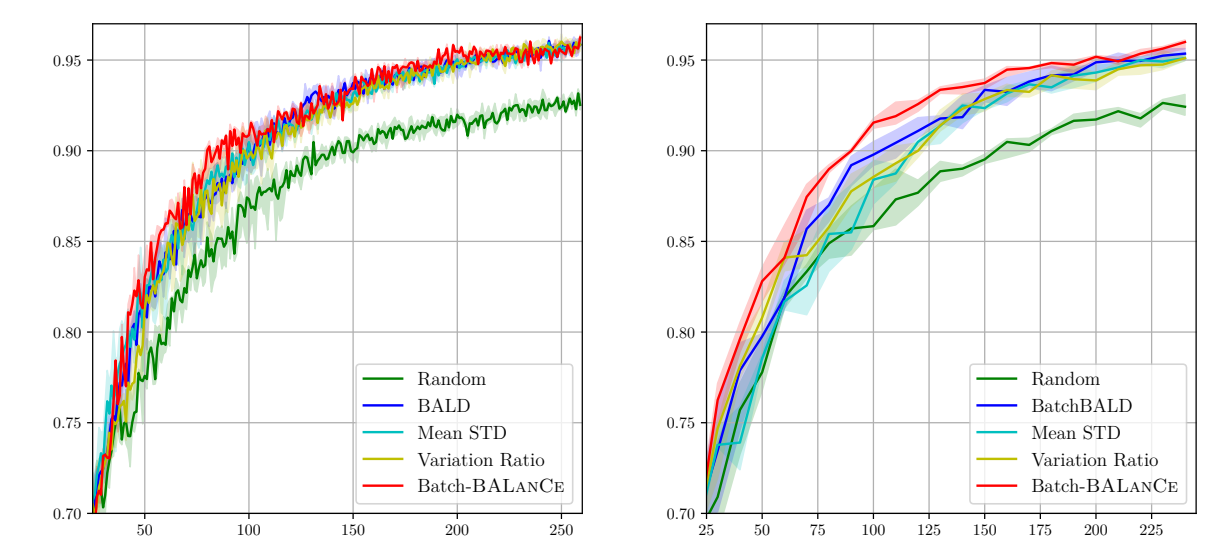
Anneal  $\tau$  at each step

## 7 Experiments

### Dataset details

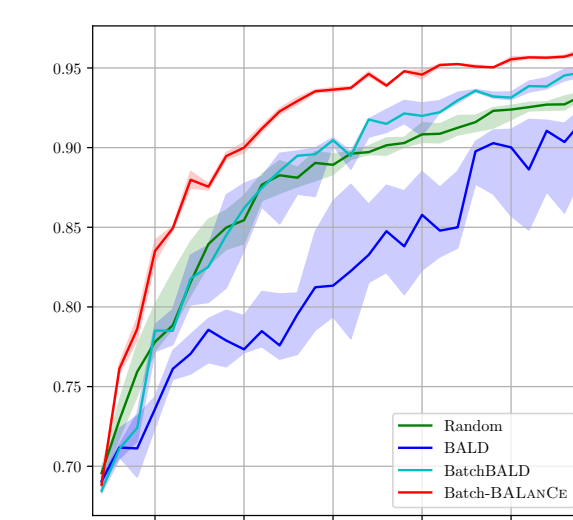
	CINIC-10	EMNIST-Balanced	EMNIST-ByMerge	EMNIST-ByClass	Fashion-MNIST	MNIST
# classes	10	47	47	62	10	10
# samples	270,000	131,600	814,255	814,255	70,000	70,000
Budget	1,400	600	600	450	350	260
$\tau$	$\frac{\text{val err}}{4}$	$\frac{\text{val err}}{4}$	$\frac{\text{val err}}{4}$	$\frac{\text{val err}}{4}$	$\frac{\text{val err}}{4}$	$\frac{\text{val err}}{4}$
$ \mathcal{D}_{\text{pool}} $	40,000	18,800	18,800	18,800	39,980	10,000

### Balanced datasets

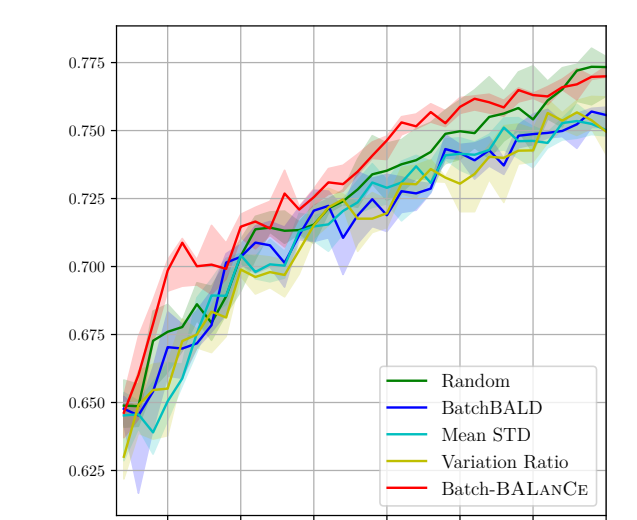


(a) ACC on MNIST B=1, K=100

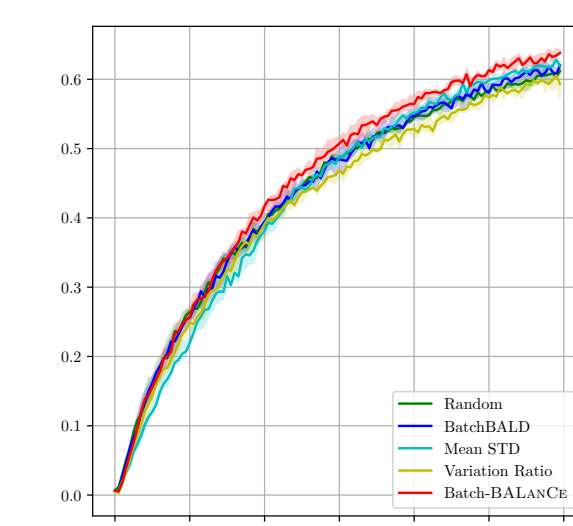
(b) ACC on MNIST B=10, K=100



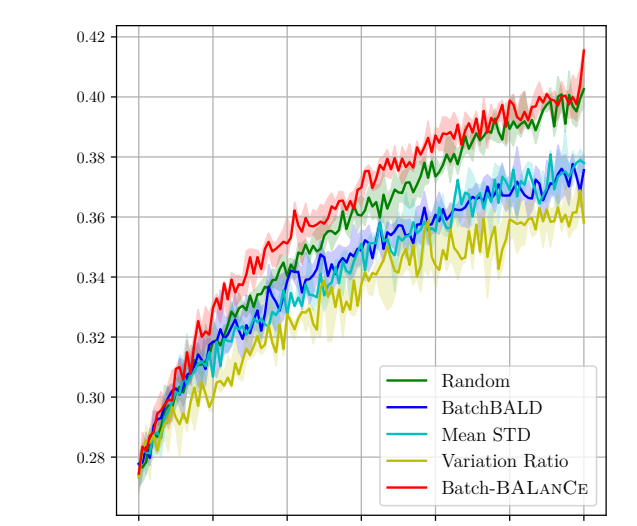
(c) ACC on Repeated-MNIST B=10, K=100



(d) ACC on Fashion-MNIST B=10, K=100

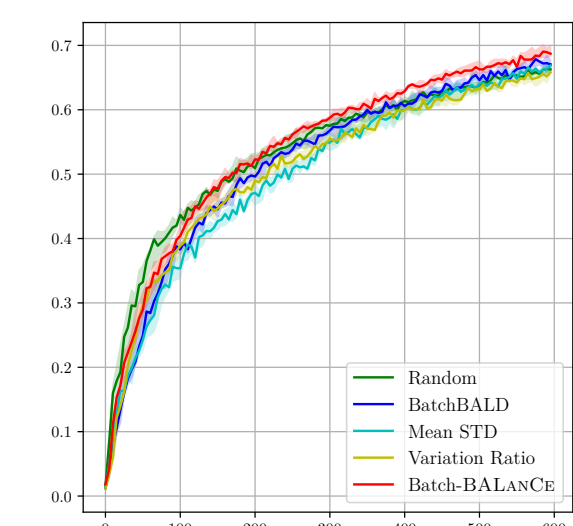


(e) F1 on EMNIST-balanced B=5, K=10

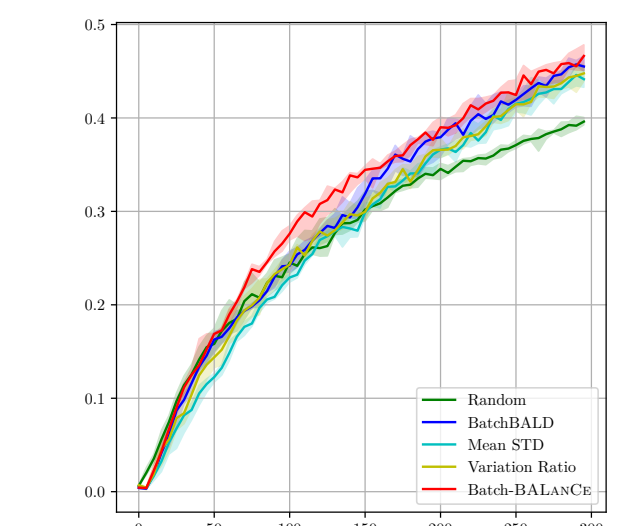


(f) ACC on CINIC-10 B=10, K=50

### Imbalanced datasets



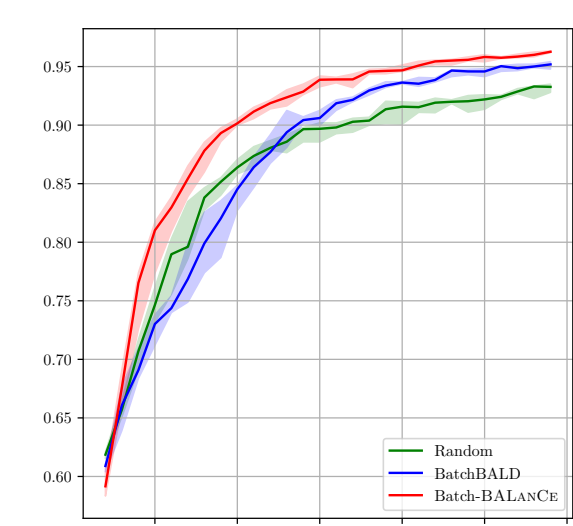
(g) ACC on EMNIST-Bymerge B=5, K=10



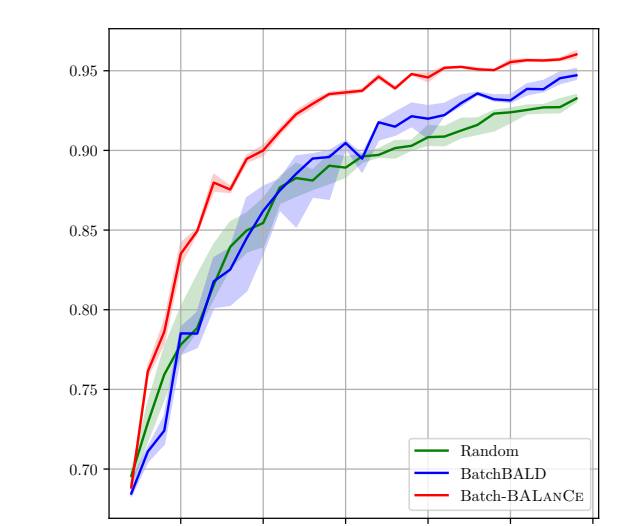
(h) F1 on EMNIST-Bymerge B=5, K=10

### Increase repeat number

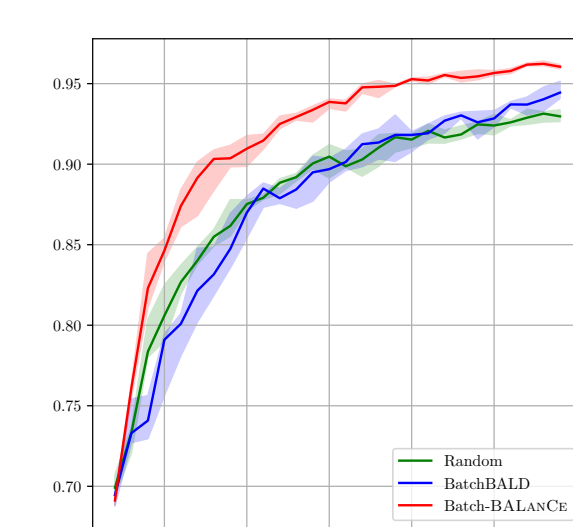
Performance on Repeated-MNIST with different repeat number



(a) Repeat 1 time



(b) Repeat 2 times



(c) Repeat 3 times



QR code of our paper