

Dataset Distillation

Related works

- Knowledge Distillation[1]
- Dataset Distillation[2]
- Coreset Construction[3, 4]
- Dataset Condensation

[1] Distilling the Knowledge in a Neural Network. (NIPS'15)

[2] Dataset distillation. (Arxiv'18)

[3] Coresets for Data-efficient Training of Machine Learning Models. (ICML'20)

[4] Coresets for Robust Training of Neural Networks against Noisy Labels. (NIPS'21)

Background

$$\mathcal{L}(f_\theta, \mathcal{D}) = \sum_{(x, y) \sim \mathcal{D}} \ell(f_\theta(x), y)$$

□ Dataset Distillation

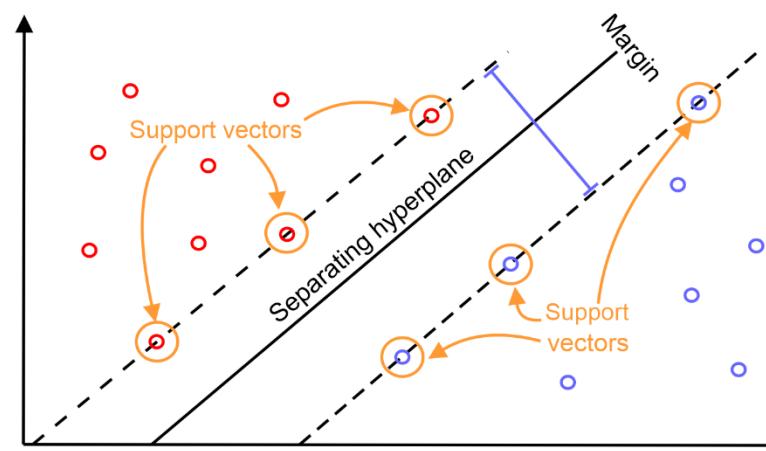
$$\begin{aligned} & \min_{\mathcal{D} \in \mathcal{S}} \quad \mathcal{L}(f_{\theta^*(\mathcal{D})}, \mathcal{D}_{test}) \\ \text{s.t.} \quad & \theta^*(\mathcal{D}) \in \operatorname{argmin}_{\theta \in \Theta} \mathcal{L}(f_\theta, \mathcal{D}) \end{aligned}$$

$$\mathcal{S} = \{\mathcal{D} : |\mathcal{D}| \leq m\}$$

□ Coreset

$$\begin{aligned} & \min_{\mathcal{D} \in \mathcal{S}} \quad \mathcal{L}(f_{\theta^*(\mathcal{D})}, \mathcal{D}_{test}) \\ \text{s.t.} \quad & \theta^*(\mathcal{D}) \in \operatorname{argmin}_{\theta \in \Theta} \mathcal{L}(f_\theta, \mathcal{D}) \end{aligned}$$

$$\mathcal{S} = \{\mathcal{D} : |\mathcal{D}| \leq m, \mathcal{D} \subseteq \mathcal{D}_{train}\}$$



DATASET CONDENSATION WITH GRADIENT MATCHING

Bo Zhao, Konda Reddy Mopuri, Hakan Bilen

School of Informatics, The University of Edinburgh
`{bo.zhao, kmopuri, hbilen}@ed.ac.uk`

Dataset Condensation



Figure 2: Visualization of condensed 1 image/class with ConvNet for MNIST, Fashion-MNIST, SVHN and CIFAR10.

Dataset Condensation

	Img/Clss	Ratio %	Random	Coreset Selection			Ours	Whole Dataset
				Herding	K-Center	Forgetting		
MNIST	1	0.017	64.9±3.5	89.2±1.6	89.3±1.5	35.5±5.6	91.7±0.5	
	10	0.17	95.1±0.9	93.7±0.3	84.4±1.7	68.1±3.3	97.4±0.2	99.6±0.0
	50	0.83	97.9±0.2	94.9±0.2	97.4±0.3	88.2±1.2	98.8±0.2	
FashionMNIST	1	0.017	51.4±3.8	67.0±1.9	66.9±1.8	42.0±5.5	70.5±0.6	
	10	0.17	73.8±0.7	71.1±0.7	54.7±1.5	53.9±2.0	82.3±0.4	93.5±0.1
	50	0.83	82.5±0.7	71.9±0.8	68.3±0.8	55.0±1.1	83.6±0.4	
SVHN	1	0.014	14.6±1.6	20.9±1.3	21.0±1.5	12.1±1.7	31.2±1.4	
	10	0.14	35.1±4.1	50.5±3.3	14.0±1.3	16.8±1.2	76.1±0.6	95.4±0.1
	50	0.7	70.9±0.9	72.6±0.8	20.1±1.4	27.2±1.5	82.3±0.3	
CIFAR10	1	0.02	14.4±2.0	21.5±1.2	21.5±1.3	13.5±1.2	28.3±0.5	
	10	0.2	26.0±1.2	31.6±0.7	14.7±0.9	23.3±1.0	44.9±0.5	84.8±0.1
	50	1	43.4±1.0	40.4±0.6	27.0±1.4	23.3±1.1	53.9±0.5	

Table 1: The performance comparison to coresnet methods. This table shows the testing accuracies (%) of different methods on four datasets. ConvNet is used for training and testing. Img/Clss: image(s) per class, Ratio (%): the ratio of condensed images to whole training set.

Dataset Condensation

Dataset	Img/Cls	DD	Ours	Whole Dataset
MNIST	1 10	- 79.5 ± 8.1	85.0 ± 1.6 93.9 ± 0.6	99.5 ± 0.0
CIFAR10	1 10	- 36.8 ± 1.2	24.2 ± 0.9 39.1 ± 1.2	83.1 ± 0.2

Table 3: Comparison to DD (Wang et al., 2018) in terms of testing accuracy (%).

Methods

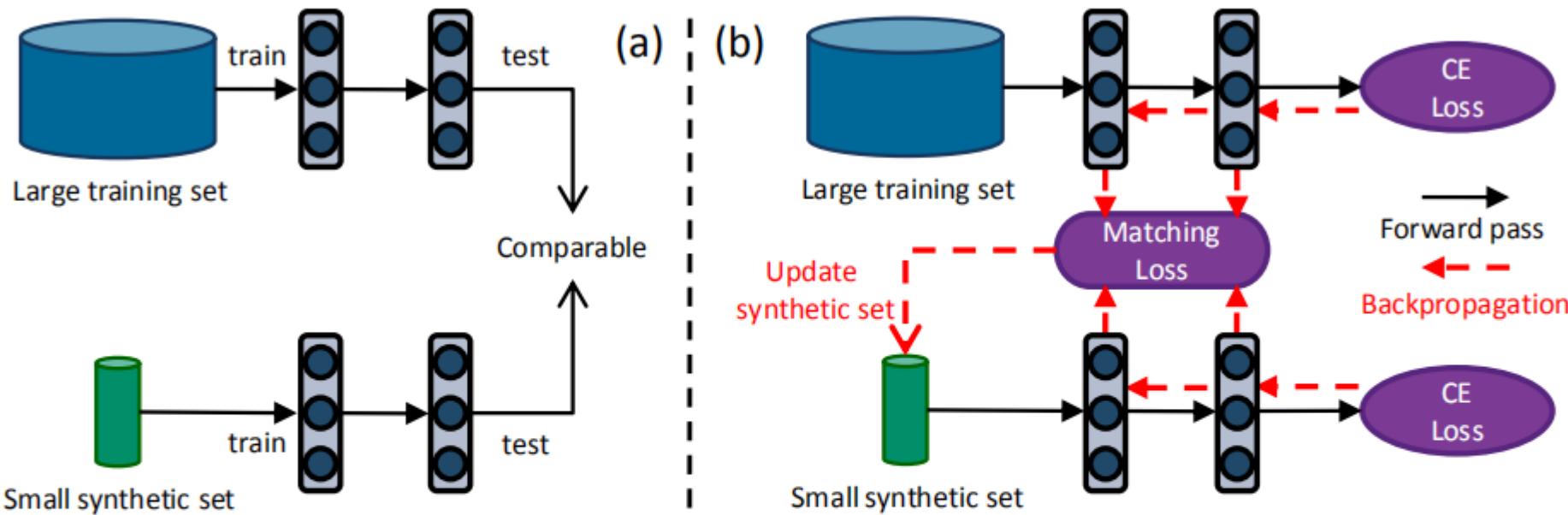


Figure 1: Dataset Condensation (left) aims to generate a small set of synthetic images that can match the performance of a network trained on a large image dataset. Our method (right) realizes this goal by learning a synthetic set such that a deep network trained on it and the large set produces similar gradients w.r.t. its weights. The synthetic data can later be used to train a network from scratch in a small fraction of the original computational load. CE denotes Cross-Entropy.

Methods

$$\mathcal{L}^{\mathcal{T}}(\boldsymbol{\theta}) = \frac{1}{|\mathcal{T}|} \sum_{(\mathbf{x}, y) \in \mathcal{T}} \ell(\phi_{\boldsymbol{\theta}}(\mathbf{x}), y)$$

$$\mathcal{T} = \{(\mathbf{x}_i, y_i)\}_{i=1}^{|\mathcal{T}|} \quad \boldsymbol{\theta}^{\mathcal{T}} = \arg \min_{\boldsymbol{\theta}} \mathcal{L}^{\mathcal{T}}(\boldsymbol{\theta})$$

$$\mathcal{S} = \{(\mathbf{s}_i, y_i)\}_{i=1}^{|\mathcal{S}|} \quad \boldsymbol{\theta}^{\mathcal{S}} = \arg \min_{\boldsymbol{\theta}} \mathcal{L}^{\mathcal{S}}(\boldsymbol{\theta})$$

$$\mathbb{E}_{\mathbf{x} \sim P_{\mathcal{D}}} [\ell(\phi_{\boldsymbol{\theta}^{\mathcal{T}}}(\mathbf{x}), y)] \simeq \mathbb{E}_{\mathbf{x} \sim P_{\mathcal{D}}} [\ell(\phi_{\boldsymbol{\theta}^{\mathcal{S}}}(\mathbf{x}), y)]$$



$$\mathcal{S}^* = \arg \min_{\mathcal{S}} \mathcal{L}^{\mathcal{T}}(\boldsymbol{\theta}^{\mathcal{S}}(\mathcal{S})) \quad \text{subject to} \quad \boldsymbol{\theta}^{\mathcal{S}}(\mathcal{S}) = \arg \min_{\boldsymbol{\theta}} \mathcal{L}^{\mathcal{S}}(\boldsymbol{\theta}).$$

Methods

$$\min_{\mathcal{S}} D(\boldsymbol{\theta}^{\mathcal{S}}, \boldsymbol{\theta}^{\mathcal{T}}) \quad \text{subject to} \quad \boldsymbol{\theta}^{\mathcal{S}}(\mathcal{S}) = \arg \min_{\boldsymbol{\theta}} \mathcal{L}^{\mathcal{S}}(\boldsymbol{\theta})$$



$$\min_{\mathcal{S}} \mathbb{E}_{\boldsymbol{\theta}_0 \sim P_{\boldsymbol{\theta}_0}} [D(\boldsymbol{\theta}^{\mathcal{S}}(\boldsymbol{\theta}_0), \boldsymbol{\theta}^{\mathcal{T}}(\boldsymbol{\theta}_0))] \quad \text{subject to} \quad \boldsymbol{\theta}^{\mathcal{S}}(\mathcal{S}) = \arg \min_{\boldsymbol{\theta}} \mathcal{L}^{\mathcal{S}}(\boldsymbol{\theta}(\boldsymbol{\theta}_0))$$



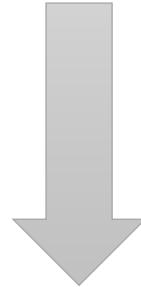
$$\min_{\mathcal{S}} \mathbb{E}_{\boldsymbol{\theta}_0 \sim P_{\boldsymbol{\theta}_0}} \left[\sum_{t=0}^{T-1} D(\boldsymbol{\theta}_t^{\mathcal{S}}, \boldsymbol{\theta}_t^{\mathcal{T}}) \right] \quad \text{subject to}$$

$$\boldsymbol{\theta}_{t+1}^{\mathcal{S}}(\mathcal{S}) = \text{opt-alg}_{\boldsymbol{\theta}}(\mathcal{L}^{\mathcal{S}}(\boldsymbol{\theta}_t^{\mathcal{S}}), \varsigma^{\mathcal{S}}) \quad \text{and} \quad \boldsymbol{\theta}_{t+1}^{\mathcal{T}} = \text{opt-alg}_{\boldsymbol{\theta}}(\mathcal{L}^{\mathcal{T}}(\boldsymbol{\theta}_t^{\mathcal{T}}), \varsigma^{\mathcal{T}})$$

Methods

$$\min_S \mathbb{E}_{\theta_0 \sim P_{\theta_0}} \left[\sum_{t=0}^{T-1} D(\theta_t^S, \theta_t^T) \right] \quad \text{subject to}$$

$$\theta_{t+1}^S(S) = \text{opt-alg}_{\theta}(\mathcal{L}^S(\theta_t^S), \varsigma^S) \quad \text{and} \quad \theta_{t+1}^T = \text{opt-alg}_{\theta}(\mathcal{L}^T(\theta_t^T), \varsigma^T)$$



$$\theta_{t+1}^S \leftarrow \theta_t^S - \eta_{\theta} \nabla_{\theta} \mathcal{L}^S(\theta_t^S) \quad \text{and} \quad \theta_{t+1}^T \leftarrow \theta_t^T - \eta_{\theta} \nabla_{\theta} \mathcal{L}^T(\theta_t^T),$$

$$D(\theta_t^S, \theta_t^T) \approx 0$$

$$\min_S \mathbb{E}_{\theta_0 \sim P_{\theta_0}} \left[\sum_{t=0}^{T-1} D(\nabla_{\theta} \mathcal{L}^S(\theta_t), \nabla_{\theta} \mathcal{L}^T(\theta_t)) \right].$$

The key idea is that we wish θ^S to be close to not only the final θ^T but also to follow a similar path to θ^T throughout the optimization

Methods

Algorithm 1: Dataset condensation with gradient matching

Input: Training set \mathcal{T}

- 1 **Required:** Randomly initialized set of synthetic samples \mathcal{S} for C classes, probability distribution over randomly initialized weights P_{θ_0} , deep neural network ϕ_θ , number of outer-loop steps K , number of inner-loop steps T , number of steps for updating weights ς_θ and synthetic samples ς_S in each inner-loop step respectively, learning rates for updating weights η_θ and synthetic samples η_S .

2 **for** $k = 0, \dots, K - 1$ **do**

3 Initialize $\theta_0 \sim P_{\theta_0}$

4 **for** $t = 0, \dots, T - 1$ **do**

5 **for** $c = 0, \dots, C - 1$ **do**

6 Sample a minibatch pair $B_c^T \sim \mathcal{T}$ and $B_c^S \sim \mathcal{S}$ $\triangleright B_c^T$ and B_c^S are of the same class c .

7 Compute $\mathcal{L}_c^T = \frac{1}{|B_c^T|} \sum_{(\mathbf{x}, y) \in B_c^T} \ell(\phi_{\theta_t}(\mathbf{x}), y)$ and $\mathcal{L}_c^S = \frac{1}{|B_c^S|} \sum_{(\mathbf{s}, y) \in B_c^S} \ell(\phi_{\theta_t}(\mathbf{s}), y)$

8 Update $\mathcal{S}_c \leftarrow \text{opt-}\text{alg}_{\mathcal{S}}(D(\nabla_{\theta} \mathcal{L}_c^S(\theta_t), \nabla_{\theta} \mathcal{L}_c^T(\theta_t)), \varsigma_S, \eta_S)$

9 Update $\theta_{t+1} \leftarrow \text{opt-}\text{alg}_{\theta}(\mathcal{L}^S(\theta_t), \varsigma_\theta, \eta_\theta)$ \triangleright Use the whole \mathcal{S}

Output: S

$$D(\nabla_{\theta}\mathcal{L}^{\mathcal{S}}, \nabla_{\theta}\mathcal{L}^{\mathcal{T}}) = \sum_{l=1}^L d(\nabla_{\theta^{(l)}}\mathcal{L}^{\mathcal{S}}, \nabla_{\theta^{(l)}}\mathcal{L}^{\mathcal{T}})$$

$$d(\mathbf{A}, \mathbf{B}) = \sum_{i=1}^{\text{out}} \left(1 - \frac{\mathbf{A}_{i\cdot} \cdot \mathbf{B}_{i\cdot}}{\|\mathbf{A}_{i\cdot}\| \|\mathbf{B}_{i\cdot}\|} \right)$$

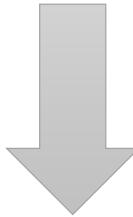
Dataset Condensation with Differentiable Siamese Augmentation

Bo Zhao¹ Hakan Bilen¹

Arxiv 2021

Methods

$$\min_{\mathcal{S}} D(\nabla_{\boldsymbol{\theta}} \mathcal{L}(\mathcal{S}, \boldsymbol{\theta}_t), \nabla_{\boldsymbol{\theta}} \mathcal{L}(\mathcal{T}, \boldsymbol{\theta}_t)),$$



$$\min_{\mathcal{S}} D(\nabla_{\boldsymbol{\theta}} \mathcal{L}(\mathcal{A}(\mathcal{S}, \omega^{\mathcal{S}}), \boldsymbol{\theta}_t), \nabla_{\boldsymbol{\theta}} \mathcal{L}(\mathcal{A}(\mathcal{T}, \omega^{\mathcal{T}}), \boldsymbol{\theta}_t)),$$

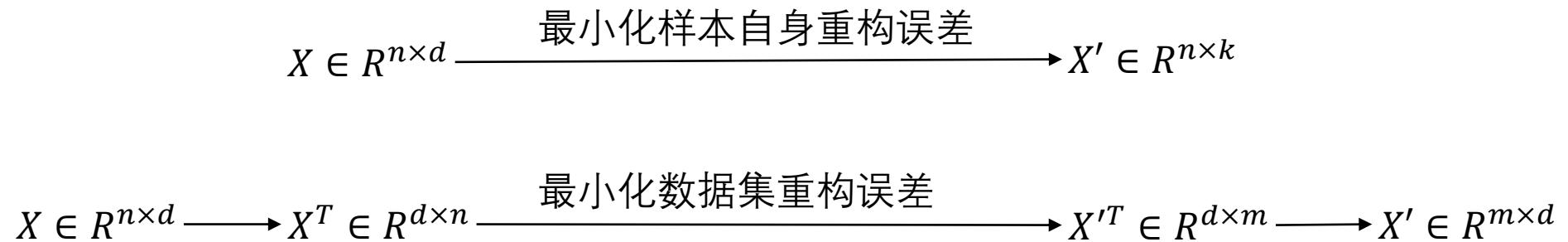
Experiments

	Img/Clss	Ratio %	Coreset Selection			DD [†]	Training Set Synthesis			Whole Dataset
			Random	Herding	Forgetting		LD [†]	DC	DSA	
MNIST	1	0.017	64.9±3.5	89.2±1.6	35.5±5.6	79.5±8.1	60.9±3.2	91.7±0.5	88.7±0.6	99.6±0.0
	10	0.17	95.1±0.9	93.7±0.3	68.1±3.3		87.3±0.7	97.4±0.2	97.8±0.1	
	50	0.83	97.9±0.2	94.8±0.2	88.2±1.2		-	93.3±0.3	98.8±0.1	99.2±0.1
FashionMNIST	1	0.017	51.4±3.8	67.0±1.9	42.0±5.5	-	-	70.5±0.6	70.6±0.6	93.5±0.1
	10	0.17	73.8±0.7	71.1±0.7	53.9±2.0	-	-	82.3±0.4	84.6±0.3	
	50	0.83	82.5±0.7	71.9±0.8	55.0±1.1	-	-	83.6±0.4	88.7±0.2	
SVHN	1	0.014	14.6±1.6	20.9±1.3	12.1±1.7	-	-	31.2±1.4	27.5±1.4	95.4±0.1
	10	0.14	35.1±4.1	50.5±3.3	16.8±1.2	-	-	76.1±0.6	79.2±0.5	
	50	0.7	70.9±0.9	72.6±0.8	27.2±1.5	-	-	82.3±0.3	84.4±0.4	
CIFAR10	1	0.02	14.4±2.0	21.5±1.2	13.5±1.2	-	25.7±0.7	28.3±0.5	28.8±0.7	84.8±0.1
	10	0.2	26.0±1.2	31.6±0.7	23.3±1.0	36.8±1.2	38.3±0.4	44.9±0.5	52.1±0.5	
	50	1	43.4±1.0	40.4±0.6	23.3±1.1	-	42.5±0.4	53.9±0.5	60.6±0.5	

Table 1. The performance comparison to coresnet selection and training set synthesis methods. This table shows the testing accuracies (%) of models trained from scratch on the small coresnet or synthetic set. Img/Clss: image(s) per class, Ratio (%): the ratio of condensed images to whole training set. DD[†] and LD[†] use LeNet for MNIST and AlexNet for CIFAR10, while the rest use ConvNet for training and testing.

Discussion

□ PCA



□ Dataset Distillation + Reconstruction Error

LDA+PCA

Idea

□ Label loss buffer & Predicted loss buffer

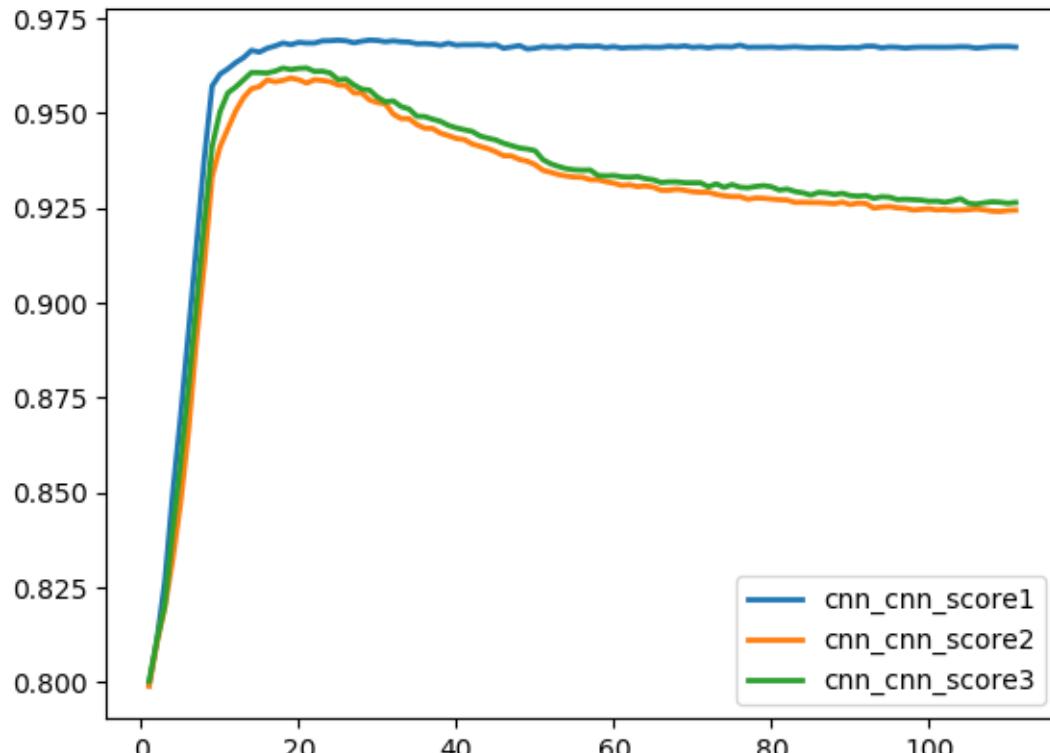
class	Label Loss					class	Predicted Loss				
1	0.13	0.27	0.36	...	0.77	1	0.13	0.17	0.16	...	0.37
2	0.68	0.70	0.82	...	0.99	2	0.30	0.21	0.11	...	0.09
3	0.11	0.13	0.34	...	0.56	3	0.11	0.13	0.24	...	0.36
...
100	0.01	0.02	0.03	...	0.07	100	0.01	0.02	0.03	...	0.07

$$\text{score}(x_i) = \frac{\ell(f(x_i), y_i)}{\frac{1}{K_{y_i}} \sum_k 1_{y_k=y_i} \cdot \ell(f(x_k), y_k)}$$

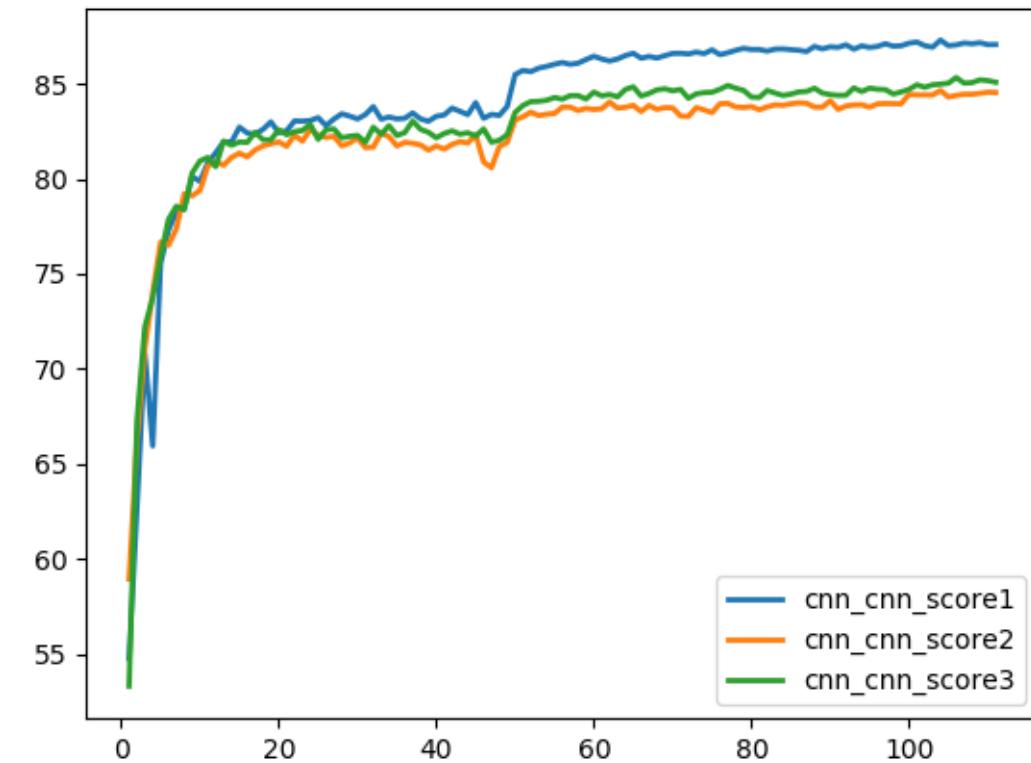
$$\text{score2}(x_i) = \frac{\ell_{label}(f(x_i), y_i) - \ell_{pred}(f(x_i), y_i^{pred})}{\frac{1}{K_{y_i}} \sum_k 1_{y_k=y_i} \cdot \ell_{label}(f(x_k), y_k)}$$

$$\text{score3}(x_i) = \frac{\ell_{label}(f(x_i), y_i) - \ell_{pred}(f(x_i), y_i^{pred})}{\frac{1}{K_{y_i}} \sum_k 1_{y_k=y_i} \cdot (\ell_{label}(f(x_k), y_k) - \ell_{pred}(f(x_i), y_k^{pred}))}$$

Idea



Pure Rate



Test Acc

Idea

□ Label loss buffer & Predicted loss buffer & Predicted class-loss buffer

class	Label Loss				
1	0.13	0.27	0.36	...	0.77
2	0.68	0.70	0.82	...	0.99
3	0.11	0.13	0.34	...	0.56
...
100	0.01	0.02	0.03	...	0.07

class	Predicted Loss				
1	0.13	0.17	0.16	...	0.37
2	0.30	0.21	0.11	...	0.09
3	0.11	0.13	0.24	...	0.36
...
100	0.01	0.02	0.03	...	0.07

Predicted class	Predicted Loss				
1	0.13	0.17	0.16	...	0.37
2	0.30	0.21	0.11	...	0.09
3	0.11	0.13	0.24	...	0.36
...
100	0.01	0.02	0.03	...	0.07

Thanks
