模式分析与机器智能
工业和信息化部重点实验室
MIIT Key Laboratory of
Pattern Analysis & Machine Intelligence

ParNeC 模式识别与神经计算研究组
PAttern Recognition and NEural Computing

# On Learning Contrastive Representations for Learning with Noisy Labels

Li Yi[1]    Sheng Liu[2]    Qi She[3]    A. Ian McLeod[1]    Boyu Wang*[1,4]

[1]University of Western Ontario, [2]NYU Center for Data Science
[3]ByteDance Inc.    [4]Vector Institute

lyi7@uwo.ca    shengliu@nyu.edu    sheqi1991@gmail.com
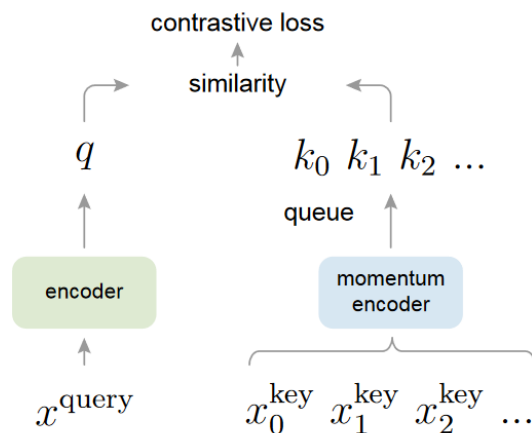aimcleod@uwo.ca    bwang@csd.uwo.ca
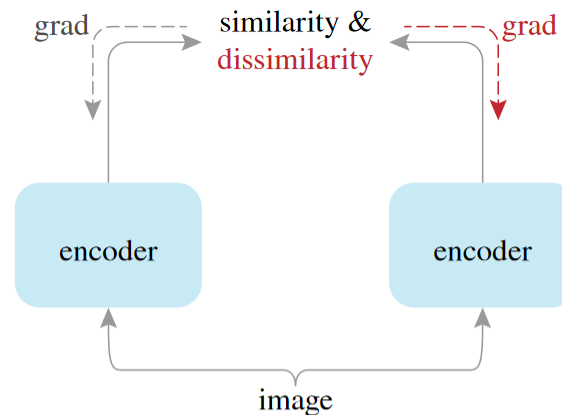
**CVPR 2022**

## Contrastive Learning Model Structure

- Momentum encoder or sharing parameters
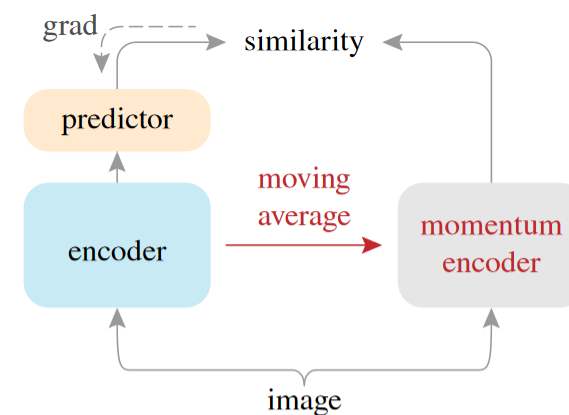
$$\theta_k \leftarrow m\theta_k + (1-m)\theta_q.$$

- Use negative samples or not
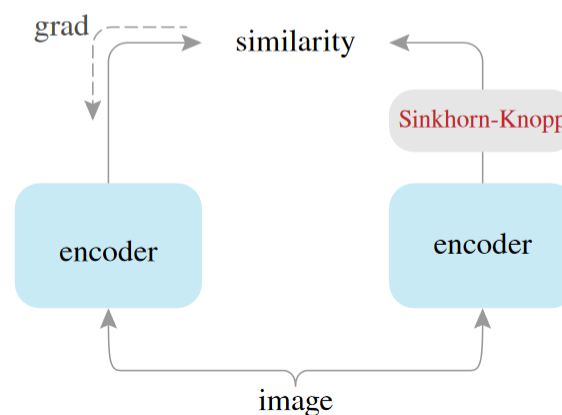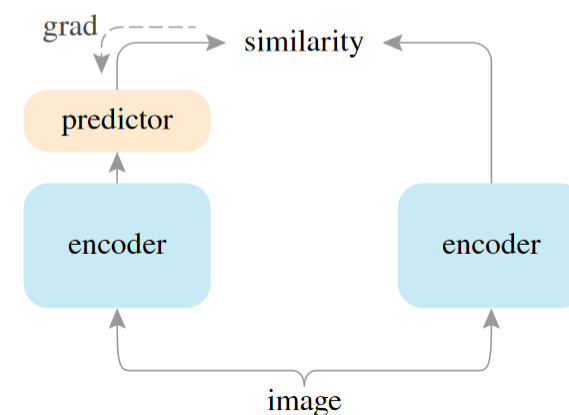- Additional predictor/projector
- Stop gradient or not



SimCLR



BYOL



MoCo



SwAV



SimSiam

## Loss design

- CE → Not robust to label noise

- Noise robust loss → suffer from the underfitting problem

- A trade-off → Explicitly or implicitly jointly used with

  the CE loss



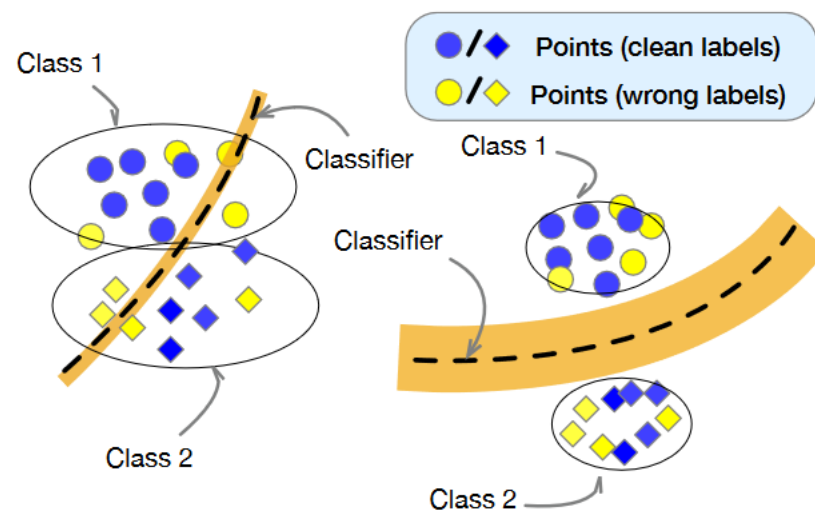Figure 1. Illustration of the proposed method with noisy labels. Black curves are the best classifiers that are learned during training. **Left:** Deep networks without contrastive regularization. **Right:** Deep networks with contrastive regularization. Two classes are better separated by deep networks that points with the same class are pulled into a tight cluster and clusters are pushed away from each other.

## Representations Induced by Contrastive Regularization

- Key component

positive contrastive pair $(x_1, x_2)$

- Unsupervised CL

Correct positive contrastive pairs are formed from <span style="color:red">two different augmentations from the same image</span>.

- Supervised CL

Correct positive contrastive pairs are formed by <span style="color:red">examples from the same class</span>.

- When encountering with noisy labels ?



Figure 2. An example of Grad-CAM [35] results of Resnet34 trained on noisy dataset with 40% symmetric label noise and clean dataset, separately. When there is label noise, information related to corrupted labels captured by the model varies from image to image (e.g. window bars in Cat 1 v.s. floor and wall in Cat 2). When there is no label noise, information related to true labels are similar for images from the same class (e.g. cat face in Cat 1 v.s. cat face in Cat 2).

Figure 5. Illustration of our framework.

$$\mathcal{L} = \mathcal{L}_{ce} + \lambda \widetilde{\mathcal{L}}_{ctr},$$

Simsiam

**Design 1**

Initial contrastive regularization function

$$\mathcal{L}_{\mathrm{ctr}}(x_i, x_j) = -\big(\langle \tilde{q}_i, \tilde{z}_j \rangle + \langle \tilde{q}_j, \tilde{z}_i \rangle\big) \mathbb{1}\{y_i = y_j\}$$

**Design 2**

Deep networks first fit examples with clean labels and the probabilistic outputs of these examples are higher than examples with corrupted labels.

$$\mathcal{L}'_{ctr}(x_i, x_j) = -\left(\langle \tilde{q}_i, \tilde{z}_j \rangle + \langle \tilde{q}_j, \tilde{z}_i \rangle\right) \mathbb{1}\{p_i^\top p_j \geq \tau\}$$

In early stage, $p_i^T p_j \approx 1$ for clean pair and 0 for noise pair

Consider two clean examples $x_i, x_j$ with clean label $y_i = y_j$
One wrongly labeled example $x_m$ with $\tilde{y}_m = y_i = y_j$

After this period?

$$\left\| \frac{\partial \mathcal{L}'_{ctr}(x_i, x_m)}{\partial q_i} \right\|_2^2 = c_i (1 - \underbrace{\tilde{q}_i^\top \tilde{q}_m}_{\approx 1})$$

$$\gg c_i (1 - \underbrace{\tilde{q}_i^\top \tilde{q}_j}_{\approx 0}) = \left\| \frac{\partial \mathcal{L}'_{ctr}(x_i, x_j)}{\partial q_i} \right\|_2^2,$$

**Design 3**

$$\widetilde{\mathcal{L}}_{\text{ctr}}(x_i, x_j) =$$

$$\left( \log \left( 1 - \langle \tilde{q}_i, \tilde{z}_j \rangle \right) + \log \left( 1 - \langle \tilde{q}_j, \tilde{z}_i \rangle \right) \right) \mathbb{1}\{ p_i^\top p_j \geq \tau \}$$

$$\left\| \frac{\partial \widetilde{\mathcal{L}}_{\text{ctr}}(x_i, x_j)}{\partial q_i} \right\|_2^2 = c_i (1 + \tilde{q}_i^\top \tilde{q}_j)$$

$$(1 + \tilde{q}_i^\top \tilde{q}_j > 1 + \tilde{q}_i^\top \tilde{q}_m \approx 1)\text{———————→} \quad \text{Does it overfit on clean data?}$$

**Table 1 — CIFAR-10**

| Method | Sym. 0% | Sym. 20% | Sym. 40% | Sym. 60% | Sym. 80% | Sym. 90% | Asym. 40% |
|---|---|---|---|---|---|---|---|
| CE | $93.97_{\pm0.22}$ | $88.51_{\pm0.17}$ | $82.73_{\pm0.16}$ | $76.26_{\pm0.29}$ | $59.25_{\pm1.01}$ | $39.43_{\pm1.17}$ | $83.23_{\pm0.59}$ |
| Forward | $93.47_{\pm0.19}$ | $88.87_{\pm0.21}$ | $83.28_{\pm0.37}$ | $75.15_{\pm0.73}$ | $58.58_{\pm1.05}$ | $38.49_{\pm1.02}$ | $82.93_{\pm0.74}$ |
| GCE | $92.38_{\pm0.32}$ | $91.22_{\pm0.25}$ | $89.26_{\pm0.34}$ | $85.76_{\pm0.58}$ | $70.57_{\pm0.83}$ | $31.25_{\pm1.04}$ | $82.23_{\pm0.61}$ |
| Co-teaching | $93.37_{\pm0.12}$ | $92.05_{\pm0.15}$ | $87.73_{\pm0.17}$ | $85.10_{\pm0.49}$ | $44.16_{\pm0.71}$ | $30.39_{\pm1.08}$ | $77.78_{\pm0.59}$ |
| LIMIT | $93.47_{\pm0.56}$ | $89.63_{\pm0.42}$ | $85.39_{\pm0.63}$ | $78.05_{\pm0.85}$ | $58.71_{\pm0.83}$ | $40.46_{\pm0.97}$ | $83.56_{\pm0.70}$ |
| SLN | $93.21_{\pm0.21}$ | $88.77_{\pm0.23}$ | $87.03_{\pm0.70}$ | $80.57_{\pm0.50}$ | $63.99_{\pm0.79}$ | $36.64_{\pm1.77}$ | $81.02_{\pm0.25}$ |
| SL | $94.21_{\pm0.13}$ | $92.45_{\pm0.08}$ | $89.22_{\pm0.08}$ | $84.63_{\pm0.21}$ | $72.59_{\pm0.23}$ | $51.13_{\pm0.27}$ | $83.58_{\pm0.60}$ |
| APL | $93.97_{\pm0.25}$ | $92.51_{\pm0.39}$ | $89.34_{\pm0.33}$ | $85.01_{\pm0.17}$ | $70.52_{\pm2.36}$ | $49.38_{\pm2.86}$ | $84.06_{\pm0.20}$ |
| CTRR | $\mathbf{94.29_{\pm0.21}}$ | $\mathbf{93.05_{\pm0.32}}$ | $\mathbf{92.16_{\pm0.31}}$ | $\mathbf{87.34_{\pm0.84}}$ | $\mathbf{83.66_{\pm0.52}}$ | $\mathbf{81.65_{\pm2.46}}$ | $\mathbf{89.00_{\pm0.56}}$ |

Table 1. Test accuracy on CIFAR-10 with different noise types and noise levels. All method use the same model PreAct ResNet18 and their best results are reported over three runs.

**Table 2 — CIFAR-100**

| Method | Sym. 0% | Sym. 20% | Sym. 40% | Sym. 60% | Sym. 80% | Asym. 40% |
|---|---|---|---|---|---|---|
| CE | $73.21_{\pm0.14}$ | $60.57_{\pm0.53}$ | $52.48_{\pm0.34}$ | $43.20_{\pm0.21}$ | $22.96_{\pm0.84}$ | $44.45_{\pm0.37}$ |
| Forward | $73.01_{\pm0.33}$ | $58.72_{\pm0.54}$ | $50.10_{\pm0.84}$ | $39.35_{\pm0.82}$ | $17.15_{\pm1.81}$ | - |
| GCE | $72.27_{\pm0.27}$ | $68.31_{\pm0.34}$ | $62.25_{\pm0.48}$ | $53.86_{\pm0.95}$ | $19.31_{\pm1.14}$ | $46.50_{\pm0.71}$ |
| Co-teaching | $73.39_{\pm0.27}$ | $65.71_{\pm0.20}$ | $57.64_{\pm0.71}$ | $31.59_{\pm0.88}$ | $15.28_{\pm1.94}$ | - |
| LIMIT | $65.53_{\pm0.91}$ | $58.02_{\pm1.93}$ | $49.71_{\pm1.81}$ | $37.05_{\pm1.39}$ | $20.01_{\pm0.11}$ | - |
| SLN | $63.13_{\pm0.21}$ | $55.35_{\pm1.26}$ | $51.39_{\pm0.48}$ | $35.53_{\pm0.58}$ | $11.96_{\pm2.03}$ | - |
| SL | $72.44_{\pm0.44}$ | $66.46_{\pm0.26}$ | $61.44_{\pm0.23}$ | $54.17_{\pm1.32}$ | $34.22_{\pm1.06}$ | $46.12_{\pm0.47}$ |
| APL | $73.88_{\pm0.99}$ | $68.09_{\pm0.15}$ | $63.46_{\pm0.17}$ | $53.63_{\pm0.45}$ | $20.00_{\pm2.02}$ | $52.80_{\pm0.52}$ |
| CTRR | $\mathbf{74.36_{\pm0.41}}$ | $\mathbf{70.09_{\pm0.45}}$ | $\mathbf{65.32_{\pm0.20}}$ | $\mathbf{54.20_{\pm0.34}}$ | $\mathbf{43.69_{\pm0.28}}$ | $\mathbf{54.47_{\pm0.37}}$ |

Table 2. Test accuracy on CIFAR-100 with different noise levels. All method use the same model PreAct ResNet18 and their best results are reported over three runs.

**Table 3**

| Method | ANIMAL-10N | Clothing1M |
|---|---|---|
| CE | $83.18_{\pm0.15}$ | $70.88_{\pm0.45}$ |
| Forward | $83.67_{\pm0.31}$ | $71.23_{\pm0.39}$ |
| GCE | $84.42_{\pm0.39}$ | $71.34_{\pm0.12}$ |
| Co-teaching | $85.73_{\pm0.27}$ | $71.68_{\pm0.21}$ |
| SLN | $83.17_{\pm0.08}$ | $71.17_{\pm0.12}$ |
| SL | $83.92_{\pm0.28}$ | $72.03_{\pm0.13}$ |
| APL | $84.25_{\pm0.11}$ | $72.18_{\pm0.21}$ |
| CTRR | $\mathbf{86.71_{\pm0.15}}$ | $\mathbf{72.71_{\pm0.19}}$ |

Table 3. Test accuracy on the real-world datasets ANIMAL-10N and Clothing1M. The results are obtained based on three different runs.

**Forward correction** corrects loss values by a estimated noise transition matrix.
**GCE** takes advantages of both MAE loss and CE and designs a robust loss function.
**Co-teaching** maintains two networks and uses small-loss examples to update.
**LIMIT** introduces noise to gradients to avoid memorization.
**SLN** adds Gaussian noise to noisy labels to combat label noise.
**SL** uses CE loss and a reverse cross entropy loss (RCE) as a robust loss function.
**APL** (NCE+RCE) combines two mutually boosted robust loss functions for training.

| Regularization Functions | CIFAR-10 | | | | | |
|:---:|:---:|:---:|:---:|:---:|:---:|:---:|
| | 0% | 20% | 40% | 60% | 80% | 90% |
| $\mathcal{L}'_{\text{ctr}}(6)$ | $93.58_{\pm 0.11}$ | $86.05_{\pm 0.33}$ | $82.34_{\pm 0.25}$ | $74.35_{\pm 0.54}$ | $54.83_{\pm 1.00}$ | $40.96_{\pm 0.99}$ |
| $\widetilde{\mathcal{L}}_{\text{ctr}}(8)$ | $\mathbf{94.29_{\pm 0.21}}$ | $\mathbf{93.05_{\pm 0.32}}$ | $\mathbf{92.16_{\pm 0.31}}$ | $\mathbf{87.34_{\pm 0.84}}$ | $\mathbf{83.66_{\pm 0.52}}$ | $\mathbf{81.65_{\pm 2.46}}$ |

Table 4. The performance of the model with respect to different regularization functions.

$$\mathcal{L}'_{\text{ctr}}(x_i, x_j) = -\big(\langle \tilde{q}_i, \tilde{z}_j \rangle + \langle \tilde{q}_j, \tilde{z}_i \rangle\big)\mathbb{1}\{p_i^\top p_j \geq \tau\}, \quad (6)$$

$$\widetilde{\mathcal{L}}_{\text{ctr}}(x_i, x_j) = \\ \left( \log\big(1 - \langle \tilde{q}_i, \tilde{z}_j \rangle\big) + \log\big(1 - \langle \tilde{q}_j, \tilde{z}_i \rangle\big) \right)\mathbb{1}\{p_i^\top p_j \geq \tau\} \\ (8)$$

| Contrastive Frameworks | CIFAR-10 | | | | |
|:---:|:---:|:---:|:---:|:---:|:---:|
| | 20% | 40% | 60% | 80% | 90% |
| CTRR (SimSiam) | $93.05_{\pm 0.32}$ | $92.16_{\pm 0.31}$ | $87.34_{\pm 0.84}$ | $83.66_{\pm 0.52}$ | $81.65_{\pm 2.46}$ |
| CTRR (SimCLR) | $92.50_{\pm 0.35}$ | $90.12_{\pm 0.43}$ | $87.41_{\pm 0.83}$ | $84.96_{\pm 0.44}$ | $79.57_{\pm 1.32}$ |
| CTRR (BYOL) | $93.31_{\pm 0.16}$ | $92.12_{\pm 0.16}$ | $88.71_{\pm 0.52}$ | $86.99_{\pm 0.59}$ | $84.31_{\pm 0.66}$ |

Table 5. Extending our method to other contrasitve learning frameworks.

Figure 4. Analysis of $\lambda$ and $\tau$ on CIFAR-10 with 60% symmetric label noise.

| Label Correction Technique | CIFAR-10 | | | |
|---|---|---|---|---|
| | 20% | 40% | 60% | 80% |
| ✗ | $93.05_{\pm 0.32}$ | $92.16_{\pm 0.31}$ | $87.34_{\pm 0.84}$ | $83.66_{\pm 0.52}$ |
| ✓ | $\mathbf{93.32_{\pm 0.11}}$ | $\mathbf{92.76_{\pm 0.67}}$ | $\mathbf{89.23_{\pm 0.18}}$ | $\mathbf{85.40_{\pm 0.93}}$ |

Table 6. ✓/✗ indicates the label correction technique is enabled/disabled.

| Method | CIFAR-10 | | | |
|---|---|---|---|---|
| | 20% | 40% | 60% | 80% |
| GCE | $91.22_{\pm 0.25}$ | $89.26_{\pm 0.34}$ | $85.76_{\pm 0.58}$ | $70.57_{\pm 0.83}$ |
| CTRR | $93.05_{\pm 0.32}$ | $92.16_{\pm 0.31}$ | $87.34_{\pm 0.84}$ | $83.66_{\pm 0.52}$ |
| CTRR+GCE | $\mathbf{93.94_{\pm 0.09}}$ | $\mathbf{93.06_{\pm 0.29}}$ | $\mathbf{92.79_{\pm 0.06}}$ | $\mathbf{90.25_{\pm 0.40}}$ |

Table 7. The performance of the model with respect to GCE, CTRR and CTRR+GCE.

# Thanks