## Theory-Inspired Contrastive Learning with Self-Labeling Refinement

Pan Zhou, Caiming Xiong, Xiaotong Yuan, and Steven HOI

Salesforce panzhou3@gmail.com

Dec 06, 2021

## Outline

**Motivation:** why one-hot label in MoCo is not accurate?

Solution for accurate label: self-labeling refinery and momentum mixup

**Experiments:** higher classification accuracy

Conclusion

### Background: MoCo

**MoCo** pulls together crops of same image, pushs away crops of different images.



### Background: MoCo

**MoCo** pulls together crops of same image, pushs away crops of different images.



Given a minibatch samples  $\{c_i\}_{i=1}^s$ , it augments each  $c_i$  into two views  $(x_i, \tilde{x}_i)$  and optimizes:

$$\mathcal{L}_{n}(\boldsymbol{w}) = -\frac{1}{s} \sum_{i=1}^{s} \log \Bigl( \frac{\sigma(\boldsymbol{x}_{i}, \widetilde{\boldsymbol{x}}_{i})}{\sigma(\boldsymbol{x}_{i}, \widetilde{\boldsymbol{x}}_{i}) + \sum_{l=1}^{b} \sigma(\boldsymbol{x}_{i}, \boldsymbol{b}_{l})} \Bigr),$$
query positive key negative key

- $\sigma(\mathbf{x}_i, \widetilde{\mathbf{x}}_i) = \exp\left(-\frac{\langle f(\mathbf{x}_i), g(\widetilde{\mathbf{x}}_i) \rangle}{\tau \| f(\mathbf{x}_i) \|_2 \cdot \| g(\widetilde{\mathbf{x}}_i) \|_2}\right)$ : cosine similarity function
- $f_w$  : online network
- $g_{\xi}$  : target network
- $B = \{b_i\}_{i=1}^b$ : negative key buffer (previous minibatch crops)

**One-hot label assignment**: query  $m{x}_i$  has only one positive  $\widetilde{m{x}}_i$  among  $\widetilde{m{x}}_i \cup m{B}$ 

$$\mathcal{L}_{n}(\boldsymbol{w}) = -\frac{1}{s} \sum_{i=1}^{s} \log \Big( \frac{\sigma(\boldsymbol{x}_{i}, \widetilde{\boldsymbol{x}}_{i})}{\sigma(\boldsymbol{x}_{i}, \widetilde{\boldsymbol{x}}_{i}) + \sum_{l=1}^{b} \sigma(\boldsymbol{x}_{i}, \boldsymbol{b}_{l})} \Big),$$
query positive key negative key

**One-hot label assignment**: query  $m{x}_i$  has only one positive  $\widetilde{m{x}}_i$  among  $\widetilde{m{x}}_i \cup m{B}$ 



#### **Issues of hot label assignment**: imprecise & uninformative

• some negatives & query belong to same semantic class



**One-hot label assignment**: query  $m{x}_i$  has only one positive  $\widetilde{m{x}}_i$  among  $\widetilde{m{x}}_i \cup m{B}$ 



**Issues of hot label assignment**: imprecise and and uninformative

- some negatives & query belong to same semantic class
- random augmentations provides crops with different semantic information, e.g. image having several objects







. . .

**One-hot label assignment**: query  $m{x}_i$  has only one positive  $\widetilde{m{x}}_i$  among  $\widetilde{m{x}}_i \cup m{B}$ 



#### **Issues of hot label assignment**: imprecise and and uninformative

- some negatives & query belong to same semantic class
- random augmentations provides crops with different semantic information, e.g. image having several objects
- different negatives have different similarity to query



**One-hot label assignment**: query  $m{x}_i$  has only one positive  $\widetilde{m{x}}_i$  among  $\widetilde{m{x}}_i \cup m{B}$ 

$$\mathcal{L}_{n}(\boldsymbol{w}) = -\frac{1}{s} \sum_{i=1}^{s} \log \Bigl( \frac{\sigma(\boldsymbol{x}_{i}, \widetilde{\boldsymbol{x}}_{i})}{\sigma(\boldsymbol{x}_{i}, \widetilde{\boldsymbol{x}}_{i}) + \sum_{l=1}^{b} \sigma(\boldsymbol{x}_{i}, \boldsymbol{b}_{l})} \Bigr),$$
query positive key negative key

#### **Issues of hot label assignment**: imprecise and and uninformative

- some negatives & query belong to same semantic class
- random augmentations provides crops with different semantic information, e.g. image having several objects
- different negatives have different similarity to query

#### **Result:** one-hot label cannot guarantee semantically similar samples to close

Support that the pair  $(x_i, \tilde{x}_i)$  in the training dataset  $\mathcal{D} = \{(x_i, \tilde{x}_i)\}_{i=1}^n$  sampled from an unknown distribution  $\mathcal{S}$  denotes the positive pair in MoCo.

Assume the query  $x_i$  has ground truth soft label  $y_i^* \in b+1$  over the key set  $B_i = \{ \tilde{x}_i \cup B \}$ where  $y_{it}^*$  measures the semantic similarity between  $x_i$  and the t-th key  $b'_t$  in buffer  $B_i$ 

Theorem 1 (upper bound of generalization error, informal).

Under proper assumptions, for MoCo, with probability  $1 - \nu$ , the generalization error on instance discrimination task can be upper bounded as :

$$\underline{\text{generalization error}} \leq \mathcal{O}\big(\mathbb{E}_{\boldsymbol{\mathcal{D}}\sim\boldsymbol{\mathcal{S}}}\left[\|\boldsymbol{y}-\boldsymbol{y}^*\|_2\right]\big) + \mathcal{O}\Big(\sqrt{\frac{V_{\boldsymbol{\mathcal{D}}}\ln(|\mathcal{F}|/\nu)}{n}} + \frac{\ln(|\mathcal{F}|/\nu)}{n}\Big),$$

training & test error gap

where  $V_{\mathcal{D}}$  is the variance of  $f_w$  on data  $\mathcal{D}$ ,  $\mathcal{F}$  is the covering number of encoder  $f_w$ 

Support that the pair  $(x_i, \tilde{x}_i)$  in the training dataset  $\mathcal{D} = \{(x_i, \tilde{x}_i)\}_{i=1}^n$  sampled from an unknown distribution  $\mathcal{S}$  denotes the positive pair in MoCo.

Assume the query  $x_i$  has ground truth soft label  $y_i^* \in b+1$  over the key set  $B_i = \{ \tilde{x}_i \cup B \}$ where  $y_{it}^*$  measures the semantic similarity between  $x_i$  and the t-th key  $b'_t$  in buffer  $B_i$ 

Theorem 1 (upper bound of generalization error, informal).

Under proper assumptions, for MoCo, with probability  $1 - \nu$ , the generalization error on instance discrimination task can be upper bounded as :

$$\begin{array}{c|c} \text{one-hot label} & \text{true soft label} \\ \hline \text{generalization error} \leq \mathcal{O}\left(\mathbb{E}_{\mathcal{D}\sim\mathcal{S}}\left[\|\boldsymbol{y}-\boldsymbol{y}^{*}\|_{2}\right]\right) + \mathcal{O}\left(\sqrt{\frac{V_{\mathcal{D}}\ln(|\mathcal{F}|/\nu)}{n}} + \frac{\ln(|\mathcal{F}|/\nu)}{n}\right), \end{array}$$

training & test error gap

where  $V_{\mathcal{D}}$  is the variance of  $f_w$  on data  $\mathcal{D}$ ,  $\mathcal{F}$  is the covering number of encoder  $f_w$ 

Support that the pair  $(x_i, \tilde{x}_i)$  in the training dataset  $\mathcal{D} = \{(x_i, \tilde{x}_i)\}_{i=1}^n$  sampled from an unknown distribution  $\mathcal{S}$  denotes the positive pair in MoCo.

Assume the query  $x_i$  has ground truth soft label  $y_i^* \in b+1$  over the key set  $B_i = \{ \tilde{x}_i \cup B \}$ where  $y_{it}^*$  measures the semantic similarity between  $x_i$  and the t-th key  $b'_t$  in buffer  $B_i$ 

Theorem 1 (upper bound of generalization error, informal).

Under proper assumptions, for MoCo, with probability  $1 - \nu$ , the generalization error on instance discrimination task can be upper bounded as :

where  $V_{\mathcal{D}}$  is the variance of  $f_w$  on data  $\mathcal{D}$ ,  $\mathcal{F}$  is the covering number of encoder  $f_w$ 

Support that the pair  $(x_i, \tilde{x}_i)$  in the training dataset  $\mathcal{D} = \{(x_i, \tilde{x}_i)\}_{i=1}^n$  sampled from an unknown distribution  $\mathcal{S}$  denotes the positive pair in MoCo.

Assume the query  $x_i$  has ground truth soft label  $y_i^* \in b+1$  over the key set  $B_i = \{\tilde{x}_i \cup B\}$ where  $y_{it}^*$  measures the semantic similarity between  $x_i$  and the t-th key  $b'_t$  in buffer  $B_i$ 

Theorem 2 (lower bound of generalization error, informal).

Under proper assumptions, for MoCo, there exists a contrastive learning problem such that the generalization error on instance discrimination task is lower bounded as :

one-hot label true soft label  
generalization error 
$$\geq \mathcal{O}(\mathbb{E}_{\mathcal{D}\sim \mathcal{S}}[\|\boldsymbol{y} - \boldsymbol{y}^*\|_2)).$$

Lower and upper bounds show that generalization error  $\sim \mathbb{E}_{D \sim S} \left[ \| y - y^* \|_2 \right]$ .

#### the more accurate of the label y, the better the generalization

## Outline

Motivation: why one-hot label in MoCo is not accurate?

#### Solution for accurate label: self-labeling refinery and momentum mixup

**Experiments:** higher classification accuracy

Conclusion

#### **Self-Labeling Refinement** :

- self-labeling refinery: soft label replaces one-hot label to directly improve label accuracy
- **momentum mixup:** increase similarity of positive pair to indirectly improve label accuracy

#### **Self-Labeling Refinement** :

- self-labeling refinery: soft label replaces one-hot label to directly improve label accuracy
- **momentum mixup:** increase similarity of positive pair to indirectly improve label accuracy

**Reformulation of MoCo:** for query  $x_i$  in minibatch  $\{(x_i, \tilde{x}_i)\}_{i=1}^s$ , we maximize its similarity to its positive  $\tilde{x}_i$  in the key set  $\bar{B} = \{\tilde{x}_i\}_{i=1}^s \cup \{b_i\}_{i=1}^b$  and push it away from samples in  $\bar{B}$ :

$$\mathcal{L}_{\mathrm{c}}ig(oldsymbol{w},\{(oldsymbol{x}_i,oldsymbol{y}_i)\}ig) = -rac{1}{s} \sum_{i=1}^{s} \sum_{k=1}^{s+b} oldsymbol{y}_{ik} \logigg(rac{\sigma(oldsymbol{x}_i,ar{oldsymbol{b}}_k)}{\sum_{l=1}^{s+b} \sigma(oldsymbol{x}_i,ar{oldsymbol{b}}_l)}igg),$$

where  $ar{m{b}}_k$  is the k-th sample in  $ar{m{B}}$ ,  $m{y}_i$  is the one-hot label of query  $m{x}_i$  whose i-th entry  $m{y}_{ii}$  is 1.

#### **Self-Labeling Refinement** :

- self-labeling refinery: soft label replaces one-hot label to directly improve label accuracy
- **momentum mixup:** increase similarity of positive pair to indirectly improve label accuracy

**Reformulation of MoCo:** for query  $x_i$  in minibatch  $\{(x_i, \tilde{x}_i)\}_{i=1}^s$ , we maximize its similarity to its positive  $\tilde{x}_i$  in the key set  $\bar{B} = \{\tilde{x}_i\}_{i=1}^s \cup \{b_i\}_{i=1}^b$  and push it away from samples in  $\bar{B}$ :

$$\mathcal{L}_{c}(\boldsymbol{w},\{(\boldsymbol{x}_{i},\boldsymbol{y}_{i})\}) = -rac{1}{s}\sum_{i=1}^{s}\sum_{k=1}^{s+b} \boldsymbol{y}_{ik}\logigg(rac{\sigma(\boldsymbol{x}_{i},ar{m{b}}_{k})}{\sum_{l=1}^{s+b}\sigma(\boldsymbol{x}_{i},ar{m{b}}_{l})}igg),$$

where  $ar{m{b}}_k$  is the k-th sample in  $ar{m{B}}$ ,  $m{y}_i$  is the one-hot label of query  $m{x}_i$  whose i-th entry  $m{y}_{ii}$  is 1.

Benefit of Reformulation: labels of different samples are defined on a shared dictionary, and thus can be linearly combined.

**Self-Labeling Refinement** : (1) self-labeling refinery; (2) momentum mixup

Self-labeling refinery iteratively employs network and data to improve labels during training.

• Step1. for query  $x_i$ , we use its positive  $\tilde{x}_i$  to estimate semantic similarity between  $x_i$  and instances in  $\bar{B} = {\{\tilde{x}_i\}_{i=1}^s \cup {\{b_i\}_{i=1}^b}}$ , since  $x_i$  and  $\tilde{x}_i$  come from the same image:

$$\boldsymbol{p}_{ik}^{t} = \sigma^{1/\tau'}(\boldsymbol{\widetilde{x}}_{i}, \boldsymbol{\overline{b}}_{k}) / \sum_{l=1}^{s+b} \sigma^{1/\tau'}(\boldsymbol{\widetilde{x}}_{i}, \boldsymbol{\overline{b}}_{l}),$$

**Self-Labeling Refinement** : (1) self-labeling refinery; (2) momentum mixup

Self-labeling refinery iteratively employs network and data to improve labels during training.

• Step1. for query  $x_i$ , we use its positive  $\tilde{x}_i$  to estimate semantic similarity between  $x_i$  and instances in  $\bar{B} = {\{\tilde{x}_i\}_{i=1}^s \cup {\{b_i\}_{i=1}^b}}$ , since  $x_i$  and  $\tilde{x}_i$  come from the same image:

$$\boldsymbol{p}_{ik}^{t} = \sigma^{1/\tau'}(\boldsymbol{\widetilde{x}}_{i}, \boldsymbol{\overline{b}}_{k}) / \sum_{l=1}^{s+b} \sigma^{1/\tau'}(\boldsymbol{\widetilde{x}}_{i}, \boldsymbol{\overline{b}}_{l}),$$

• Step2. as  $\tilde{x}_i$  is highly similar to itself in  $\bar{B}$ ,  $p_{ii}^t$  will be much larger than others and conceals the similarity of other semantically similar instances in $\bar{B}$ . So we remove  $\tilde{x}_i$  from  $\bar{B}$ 

$$\boldsymbol{q}_{ik}^{t} = \sigma^{1/\tau'}(\widetilde{\boldsymbol{x}}_{i}, \overline{\boldsymbol{b}}_{k}) / \sum_{l=1, l \neq i}^{s+b} \sigma^{1/\tau'}(\widetilde{\boldsymbol{x}}_{i}, \overline{\boldsymbol{b}}_{l}), \quad \boldsymbol{q}_{ii}^{t} = 0.$$

**Self-Labeling Refinement** : (1) self-labeling refinery; (2) momentum mixup

Self-labeling refinery iteratively employs network and data to improve labels during training.

• Step1. for query  $x_i$ , we use its positive  $\tilde{x}_i$  to estimate semantic similarity between  $x_i$  and instances in  $\overline{B} = {\{\tilde{x}_i\}_{i=1}^s \cup {\{b_i\}_{i=1}^b}}$ , since  $x_i$  and  $\tilde{x}_i$  come from the same image:

$$\boldsymbol{p}_{ik}^{t} = \sigma^{1/\tau'}(\widetilde{\boldsymbol{x}}_{i}, \overline{\boldsymbol{b}}_{k}) / \sum_{l=1}^{s+b} \sigma^{1/\tau'}(\widetilde{\boldsymbol{x}}_{i}, \overline{\boldsymbol{b}}_{l}),$$

• Step2. as  $\tilde{x}_i$  is highly similar to itself in  $\bar{B}$ ,  $p_{ii}^t$  will be much larger than others and conceals the similarity of other semantically similar instances in  $\bar{B}$ . So we remove  $\tilde{x}_i$  from  $\bar{B}$ 

$$\boldsymbol{q}_{ik}^{t} = \sigma^{1/\tau'}(\widetilde{\boldsymbol{x}}_{i}, \overline{\boldsymbol{b}}_{k}) / \sum_{l=1, l \neq i}^{s+b} \sigma^{1/\tau'}(\widetilde{\boldsymbol{x}}_{i}, \overline{\boldsymbol{b}}_{l}), \quad \boldsymbol{q}_{ii}^{t} = 0.$$

Linear combination:

$$\bar{\boldsymbol{y}}_i^t = (1 - \alpha_t - \beta_t) \boldsymbol{y}_i + \alpha_t \boldsymbol{p}_i^t + \beta_t \boldsymbol{q}_i^t,$$

#### Label-corrupted dataset

Let  $\{(\boldsymbol{x}_i, \boldsymbol{y}_i^*)\}_{i=1}^n$  denote the pairs of crops and ground-truth semantic label

- crop  $\boldsymbol{x}_i$  generated from vanilla sample  $\boldsymbol{c}_t$  obeys  $\|\boldsymbol{x}_i \boldsymbol{c}_t\|_2 \leq \varepsilon$
- ground-truth semantic label  $y_i^* \in \{\gamma_t\}_{t=1}^K$  of  $x_i$  is decided by its corresponding
- the classes are separated:  $|\gamma_i \gamma_k| \ge \delta$ ,  $\|c_i c_k\|_2 \ge 2\varepsilon$ ,  $(\forall i \ne k)$ ,
- for each sample  $c_i$ , at most  $\rho n_i$  augmentations are assigned to wrong labels, where  $n_i$  denotes the crop sample number of  $c_i$

#### Label-corrupted dataset

Let  $\{(\boldsymbol{x}_i, \boldsymbol{y}_i^*)\}_{i=1}^n$  denote the pairs of crops and ground-truth semantic label

- crop  $x_i$  generated from vanilla sample  $c_t$  obeys  $\|x_i c_t\|_2 \leq \varepsilon$
- ground-truth semantic label  $y_i^* \in \{\gamma_t\}_{t=1}^K$  of  $x_i$  is decided by its corresponding
- the classes are separated:  $|\gamma_i \gamma_k| \ge \delta$ ,  $\|c_i c_k\|_2 \ge 2\varepsilon$ ,  $(\forall i \neq k)$ ,
- for each sample  $c_i$ , at most  $\rho n_i$  augmentations are assigned to wrong labels, where  $n_i$ denotes the crop sample number of  $c_i$

#### Theorem 3 (exact label recovery guarantee, informal).

Under proper assumptions, the discrepancy between the label  $ar{m{y}}^t$  estimated by our refinery and the true label  $y^*$  of data  $\{x_i\}_{i=1}^n$  is bounded: estimated label  $\frac{1}{\sqrt{n}} \| \mathbf{\bar{y}}^t - \mathbf{y}^* \|_2 \le \frac{1 - \alpha_t}{\sqrt{n}} \| \mathbf{y} - \mathbf{y}^* \|_2 + \alpha_t (6\rho + \zeta),$ 

true soft label

where  $\zeta = c_6 \varepsilon K^2 \Gamma^5 \xi_3^3 \sqrt{\log K} / \lambda^2(C), \ \boldsymbol{y}^* = [\boldsymbol{y}_1^*, \cdots, \boldsymbol{y}_n^*]$ 

22

Theorem 3 (exact label recovery guarantee, informal).

Under proper assumptions, the discrepancy between the label  $\bar{y}^t$  estimated by our refinery and the true label  $y^*$  of data  $\{x_i\}_{i=1}^n$  is bounded: estimated label

$$\frac{1}{\sqrt{n}} \| \mathbf{\bar{y}}^t - \mathbf{y}^* \|_2 \le \frac{1 - \alpha_t}{\sqrt{n}} \| \mathbf{y} - \mathbf{y}^* \|_2 + \alpha_t (6\rho + \zeta),$$
  
true soft label

where  $\zeta = c_6 \varepsilon K^2 \Gamma^5 \xi_3^3 \sqrt{\log K} / \lambda^2(\boldsymbol{C}), \ \boldsymbol{y}^* = [\boldsymbol{y}_1^*, \cdots, \boldsymbol{y}_n^*]$ 

If  $\rho \leq \frac{\delta}{24}$ ,  $(1 - \alpha_0)|y_i - y_i^*| + \frac{1}{3}\alpha_0 \delta < \frac{1}{2}\delta$ , the estimated label  $\bar{y}_i^t$  predicts true label of  $y_i^*$  any crop  $x_i$ 

**Exact label recovery:**  $\gamma_{k^*} = \boldsymbol{y}_i^*$  with  $k^* = \operatorname{argmin}_{1 \le k \le \bar{K}} |\bar{\boldsymbol{y}}_i^t - \gamma_k|$ .

Theorem 3 (exact prediction of network when using label refinery, informal) Under proper assumptions, by using the refined label  $\bar{y}^t$  to train network, the error of network prediction on  $\{x_i\}_{i=1}^n$  is upper bounded predicted label

predicted label  $\frac{1}{\sqrt{n}} \|f(\boldsymbol{W}_t, \boldsymbol{X}) - \boldsymbol{y}^*\|_2 \le 6\rho + \frac{\zeta\lambda(\boldsymbol{C})}{K\Gamma^2\xi_3^2},$ true label

where  $\zeta = c_6 \varepsilon K^2 \Gamma^5 \xi_3^3 \sqrt{\log K} / \lambda^2(C), \ \boldsymbol{y}^* = [\boldsymbol{y}_1^*, \cdots, \boldsymbol{y}_n^*]$ 

If  $\rho \leq \frac{\delta}{24}$ ,  $(1 - \alpha_0)|y_i - y_i^*| + \frac{1}{3}\alpha_0 \delta < \frac{1}{2}\delta$ , for any vanilla sample  $c_k$ , network  $f(W_t, \cdot)$  predicts the true semantic label  $y_i^*$  of any augmentation x that obeys  $||x - c_k||_2 \leq \varepsilon$ :

**Exact label prediction:**  $\gamma_{k^*} = \gamma_k$  with  $k^* = \operatorname{argmin}_{1 \le i \le \overline{K}} |f(\boldsymbol{W}_t, \boldsymbol{x}) - \gamma_i|.$ 

**Self-Labeling Refinement** : (1) self-label refinery; (2) momentum mixup

**Momentum mixup** constructs virtual instance as follows:

 $\boldsymbol{x}_{i}^{\prime} = \theta \boldsymbol{x}_{i} + (1-\theta) \widetilde{\boldsymbol{x}}_{k}, \quad \boldsymbol{y}_{i}^{\prime} = \theta \bar{\boldsymbol{y}}_{i} + (1-\theta) \bar{\boldsymbol{y}}_{k}, \quad (1)$ 

where  $\tilde{x}_k$  is randomly sampled from the key set  $\{\tilde{x}_i\}_{i=1}^s$ ,  $\bar{y}_i$  denotes the refined label by selflabeling refinery,  $\theta \in [0, 1]$  obeys the beta distribution

**Benefits**: the component  $\tilde{x}_k$  in  $x'_i = \theta x_i + (1 - \theta) \tilde{x}_k$  directly increases the similarity between the query  $x'_i$  and its positive key  $\tilde{x}_k$  in  $\bar{B}$ 

momentum mixup can improve the accuracy of the label

## Outline

**Background:** what is contrastive learning and why label is not accurate?

**Solution:** self-labeling refinery and momentum mixup

#### **Experiments:** higher classification accuracy

Conclusion

## Outline

Motivation: why one-hot label in MoCo is not accurate?

Solution for accurate label: self-labeling refinery and momentum mixup

**Experiments:** higher classification accuracy

Conclusion

### **Experimental Results of Proposed NAS Method**

#### CIFAR10 and ImageNet: much lower Top-1 error





#### **Downstream tasks:** higher performance on VOC classification and detection



## Outline

Motivation: why one-hot label in MoCo is not accurate?

Solution for accurate label: self-labeling refinery and momentum mixup

**Experiments:** higher classification accuracy

#### Conclusion

### Conclusion

#### • Problems:

(1) what relationship between one-hot label and the generalization performance?more precise of labels in contrastive learning, the better the generalization

(2) how to estimate more precise labels?

we propose self-labeling refinery and momentum mixup

# Thanks!