Mutual Conditional Probability for Self-Supervised Learning







Contributions

- We propose a novel self-supervised learning (SSL) method to exploit multiple view samples.
- Mutual conditional probability embeds compact and discriminative feature representation into SSL.
- Our SSL produces favorable performance even an anly 100 anach Image Mat training.



Feature space

Multiple view samples are probabilistically modeled by means of kernel density estimation (**KDE**).

-								
Generative (GenPro)	orob	$\mathbf{p}(\boldsymbol{z} \mathcal{I})$	$= \frac{1}{MC_{\sigma}} \sum_{i=1}^{M} \exp$	$o\left(-\frac{1}{2\sigma^2}\right)$				
	loss	$\ell_{gen}(\mathbf{z}) = -$	$-\sigma^2 \log p(z \mathcal{I}) =$	$\Rightarrow \nabla_{\mathbf{z}} \ell_{gen} =$				
DiscPro	ve p	rob. p(2	$\mathcal{I} \boldsymbol{z}) = rac{\mathrm{p}(\boldsymbol{z} \mathcal{I})}{\mathrm{p}(\boldsymbol{z})}$	$\propto \frac{\sum_{i y_i=\mathcal{I}} e_i}{\sum_{i=1}^N e_i}$				
(DISCEIO)	loss	$\ell_{disc}(\mathbf{z}) =$	$-\sigma^2\log p(\mathcal{I} z)$	Z)				
Mutual cond	litior	nal prob). $p(\boldsymbol{z} \mathcal{I})p(\mathcal{I})$	z)				
(MuConPro)	loss	$\ell(z) = -c$	$\sigma^2 \log[p(z \mathcal{I})p]$	$o(\mathcal{I} \boldsymbol{z})] = \ell$				
MuConPro effectively describes								
compact & d	iscrin	ninative ideal	distributic less discriminative)n. Ie				
	Ĭ	I	ĬI	Ĭ				
GenPro p $(\boldsymbol{z} \mathcal{I})$)							
DISCPRO p $(\mathcal{I} \boldsymbol{z})$)							

1

MUCONPRO $p(\boldsymbol{z}|\mathcal{I})p(\mathcal{I}|\boldsymbol{z})$

Takumi Kobayashi

Proposed SSL

Our SSL trains a model based on MuConPro loss.

Feature representation is updated by maximizing MuConPro toward compact and discriminative ones, while disregarding irrelevant view samples via KDE.



$$\mathcal{L}_{SSL}(\boldsymbol{z}_{bi}) = -2T\log\sum_{j=1|j\neq i}^{M}\exp\left(\boldsymbol{z}_{bi}^{\top}\boldsymbol{\bar{x}}_{bj}/T\right) +$$

Compactness within image

Multiple view samples are effectively aggregated in a MuConPro manner which naturally induces compact and discriminative feature representation.

Temperature T controls the bandwidth in KDE.

Analysis by loss gradients.

• $T \rightarrow \infty$ results in BYOL [Grill+20].

$$\frac{\partial \ell_{gen}(\boldsymbol{z}_{bi})}{\partial \boldsymbol{z}_{bi}} = -\frac{\sum_{j \neq i} \exp(\boldsymbol{z}_{bi}^{\top} \bar{\boldsymbol{x}}_{bj}/T) \bar{\boldsymbol{x}}_{bj}}{\sum_{j \neq i} \exp(\boldsymbol{z}_{bi}^{\top} \bar{\boldsymbol{x}}_{bj}/T)} \xrightarrow{T \to \infty} -\frac{1}{M-1} \sum_{j \neq i} \bar{\boldsymbol{x}}_{bj} = \frac{\partial}{\partial \boldsymbol{z}} \mathop{\mathbb{E}}_{j \neq i} [-\boldsymbol{z}_{bi}^{\top} \bar{\boldsymbol{x}}_{bj}]$$

• $T \rightarrow 0$ results in too sparse updating.

 $\sum_{j\neq i} \exp(\mathbf{z}_{bi}^{\dagger} \bar{\mathbf{x}}_{bj} / T) \bar{\mathbf{x}}_{bj}$ $\frac{\partial \ell_{gen}(\mathbf{z}_{bi})}{dt} =$ $\sum_{j\neq i} \exp(\mathbf{z}_{bi}^{\top} \bar{\mathbf{x}}_{bj}/T)$ ∂z_{bi}

We simply apply T = 1.



$$\mathbf{z} - \mathbf{x}_i \|_2^2 \Big)$$

MeanShift-vec.

 $\exp(-\frac{1}{2\sigma^2}\|\boldsymbol{z}-\boldsymbol{x}_i\|_2^2)$ $\exp\left(-\frac{1}{2\sigma^2}\|\boldsymbol{z}-\boldsymbol{x}_i\|_2^2\right)$



ess compact



National Institute of AIST, Japan University of Tsukuba

B,M $\exp\left(\boldsymbol{z}_{bi}^{\top}\boldsymbol{\bar{x}}_{cj}/T\right)$ $T\log$ $c=1, j=1|c\neq b \lor j\neq i$

Discrimination across images

$$\stackrel{bj}{\longrightarrow} - rg \max_{\{\bar{\boldsymbol{x}}_{bj}\}_{j\neq i}} \boldsymbol{z}_{bi}^{\top} \bar{\boldsymbol{x}}_{bj}$$

Experimental Results





Probabilistic models are compared with various

	(a) Linear evaluation		(b) k-1	(b) k-NN evaluation		
Model	T = 0	1	∞	T = 0	1	∞
GENPRO p $(\boldsymbol{z} \mathcal{I})$	79.24	81.38	80.90	71.50	73.38	72.98
DISCPRO $p(\mathcal{I} \boldsymbol{z})$	17.86	81.56	81.08	16.58	73.46	73.38
MUCONPRO $p(\boldsymbol{z} \mathcal{I})p(\mathcal{I} \boldsymbol{z})$	78.38	81.64	81.50	70.16	74.06	74.00

Performance comparison with various SSL on 100-training epochs.

MuConPro with $4 \cdot (160 \times 160)$ views and T = 1.

Model	Acc. (%)	Model	Acc. (%)
SimCLR [5]	66.5	Supervised [16]	76.2
BYOL [14]	66.5	SwAV [3]	72.1
MoCo-v2 [7]	67.4	MS-SSL [22]	72.4
MoCo-v3 [8]	68.9	GenPro p $(\boldsymbol{z} \mathcal{I})$	70.9
SimSiam [6]	68.1	DISCPRO p $(\mathcal{I} z)$	72.1
W-MSE [12]	69.4	MUCONPRO p $(z \mathcal{I})$ p $(\mathcal{I} z)$	72.9





Ablation Study on ImageNet-100 with ResNet18. Multiple samples are generated in various ways in a constant computation budget. Cropping small-sized image patches.

Drawing from external memory bank.

Processing large-sized patch via momentum-free mode.