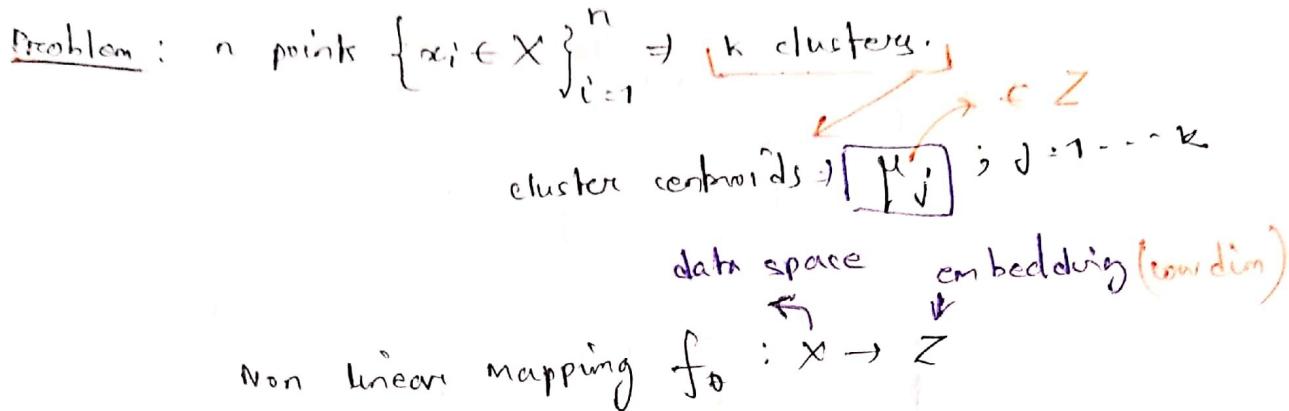


④

## Unsupervised deep embedding for clustering analysis (DEC)



Solution: Clustering with KL divergence:

① Mapping & initialization  $\mu_j$

to

$$(1 + \|z_i - \mu_j\|^2/\alpha)$$

$$f_\theta(x_i)$$

$$\text{inverse Hyperparam} = \frac{\alpha+1}{2}$$

Soft Assignment:

$$q_{ij} = \frac{1}{\sum (1 + \|z_i - \mu_j\|^2/\alpha)^{-(\frac{\alpha+1}{2})}}$$

KL Divergence Minimization: choice option for  $P$  ??

$$\text{Loss}, \quad L = \text{KL}(P || Q) = \sum_i \sum_j p_{ij} \log \frac{p_{ij}}{q_{ij}}$$

$$\Rightarrow p_{ij} = \frac{q_{ij}/f_j}{\sum_i q_{ij}/f_j} \quad // \quad f_j = \sum_i q_{ij}$$

⑤ DEC cont.

Optimization:

$$\frac{\partial L}{\partial z_i} = \frac{\alpha+1}{\alpha} \sum_j \left( 1 + \frac{\|z_i - \mu_j\|^2}{\alpha} \right)^{-1} \times (r_{ij} - q_{ij}) (z_i - \mu_j)$$

$$\frac{\partial L}{\partial \mu_j} = -\frac{\alpha+1}{\alpha} \sum_i \left( 1 + \frac{\|z_i - \mu_j\|^2}{\alpha} \right)^{-1} \times (r_{ij} - q_{ij}) (z_i - \mu_j)$$

Accuracy Metric (Classification Setting)

$$ACC = \max_m \frac{\sum_{i=1}^n \mathbb{I}\{l_i = m(c_i)\}}{n} // \text{Latent class}$$

$l_i \rightarrow$  true label

$m \Rightarrow$  mapping of cluster  $c_i$  to label  $\underline{l_i}$

DINO: distillation with no labels  
 DINO : Distillation with No Labels.

Approach: SSL with TD:  $\rightarrow$  k-dimensional output.

Student Network  $g_{\theta_S}$ ; Teacher Net  $g_{\theta_T}$

Softmax conversion of student output.

$$p_s(x)^i = \frac{\exp \{ g_{\theta_S}(x)^i / \tau_S \}}{\sum_{j=1}^k \exp \{ g_{\theta_S}(x)^j / \tau_S \}} \text{ // student}$$

$$p_t(x)^i = \frac{\exp \{ g_{\theta_T}(x)^i / \tau_T \}}{\sum_{j=1}^k \exp \{ g_{\theta_T}(x)^j / \tau_T \}} \text{ // teacher.}$$

optimization goal:  $\min_{\theta_S} H(p_t(x), p_s(x)) \text{ // } H(a, b) = -a \log b$

Actually! Local-to-Global view

$$\min_{\theta_S} \sum_{x \in \{x_1^g, x_2^g\}} \sum_{\substack{x' \in V \\ x' \neq x}} H(p_t(x), p_s(x'))$$

$V \rightarrow$  set of images with different view of  $x$ .

$$\left\{ x_1^g, x_2^g, x_1^L, \dots, x_n^L \right\}$$

$\downarrow$  global views      Local views { smaller dimension }  
 { passed through { student networks } }  
 { teacher only }

(11)

Teacher Network Update:

$$\theta_t \leftarrow \underline{\lambda \theta_t + (1-\lambda) \theta_S} // \text{Polyak-Ruppert Avg}$$

Interesting update.

Avoid collapse:

centering:  $g_t(x) \leftarrow g_t(x) + c$

$$c \leftarrow mc + (1-m) \cdot \frac{1}{B} \sum_{i=1}^B g_{t+i}(x)$$

Exp. Avg

①

## ④ Divide Mix

### Divide Mix

Training Data  $D = (\mathcal{X}, \mathcal{Y}) = \{(x_i, y_i)\}_{i=1}^N$

unlabeled  $\mathcal{U}$

one hot

$$\text{Cross entropy loss} ; \quad L(\theta) = \{\mathbf{w}\}_{c=1}^C = \left\{ \sum_{i=1}^N y_i^c \log \left( \frac{P_{\text{model}}(x_i; \theta)}{d} \right) \right\}_{c=1}^C$$

output for  $c$

Penalize the entropy term:

$$f_b = - \sum_c P_{\text{model}}(x_i; \theta) \log \left( P_{\text{model}}(x_i; \theta) \right)$$

Avoid near 0 normalized w's.

minibatch labeled data / unlabeled data

$$\{(x_b, y_b, w_b) ; b \in 1, \dots, B\}$$

$$\{w_b ; b \in 1, \dots, B\}$$

gaussian component (smaller loss)

Data clean probability  $w_i = P(g | l_i)$

↓

loss value.

compared with threshold  $\tau$

$$\text{co-refinement: } \bar{y}_b = w_b y_b + (1-w_b) p_b$$

✓

provided  
label

→ clean probability  
↓  
Avg across different  
Augmentation (prediction by net)

Applying Sharpening:

$$\hat{y}_b = \text{sharpen}(\bar{y}_b, T) = \frac{\bar{y}_b^{c^{\frac{1}{T}}}}{\sum_{c=1}^C \bar{y}_b^{c^{\frac{1}{T}}}} ; \text{ for all } T \in \mathbb{G}$$

Pair of sample  $(x_1, x_2)$  corresponding label  $(p_1, p_2)$   
 mixed  $\downarrow$   
 $(x', p')$

$$\begin{aligned} \lambda &\sim \text{Beta}(\alpha, \alpha) \\ x' &= \max(\lambda, 1-\lambda) \\ x' &= \lambda x_1 + (1-\lambda)x_2 \\ p' &= \lambda p_1 + (1-\lambda)p_2 \end{aligned} \quad \left. \begin{array}{l} \text{Kind of Augmentation.} \\ \rightarrow \text{closer to } x_1 \text{ than } x_2 \end{array} \right\}$$

Loss function:

$$\left\{ \begin{array}{l} \text{Cross entropy loss} \\ \text{loss} = \frac{1}{|x'|} \sum_{x, p \in x'} \sum_c p_c \log(p_{\text{model}}(x; \theta)) \end{array} \right.$$

$$\left\{ \begin{array}{l} \text{consistency loss} \\ \text{loss} = \frac{1}{|U'|} \sum_{x, p \in U'} \|p - p_{\text{model}}(x, \theta)\|_2^2 \\ \rightarrow \frac{1}{C} \text{ uniform} \end{array} \right.$$

$$\left\{ \begin{array}{l} \text{regularization loss} \\ \text{loss} = \sum_c \gamma_c \log \left( \frac{1}{1 + \sum_{x \in x' \cup U'} p_{\text{model}}^c(x, \theta)} \right) \end{array} \right.$$

① MixMatch

MixMatch

consistency regularization

$$\| P_{\text{model}}(y | \text{Augment}(x); \theta) - P_{\text{model}}(y | \text{Augment}(u); \theta) \|_2^2$$

Notation

Labeled data  $\mathcal{X}$ ; unlabeled data  $\mathcal{U}$

$$x' \in \mathcal{U} = \text{mixMatch}(\mathcal{X}, \mathcal{U}, T, K, \alpha)$$

$$L_x = \frac{1}{|\mathcal{X}'|} \sum_{x \in \mathcal{X}'} H(P, P_{\text{model}}(y | x; \theta))$$

$$L_u = \frac{1}{|\mathcal{U}'|} \sum_{u, q \in \mathcal{U}'} \| q - P_{\text{model}}(y | u; \theta) \|_2^2$$

Augmentation:

$$\begin{cases} \hat{x}_b = \text{Augment}(x_b) \\ \hat{u}_{b,k} = \text{Augment}(u_b) ; k \in 1 \dots K \end{cases}$$

Label guessing:

$$\hat{q}_b = \frac{1}{K} \sum_{k=1}^K P_{\text{model}}(y | \hat{u}_{b,k}; \theta)$$

Sharpening ( $P, T$ ):

$$P_i^T / \sum_j P_j^T$$

MixUp: Same as DivideMix.

① Temporal Cycle-consistency learning

Notation: Frame Sequence  $S = \{s_1, s_2, \dots, s_N\}$

Embedding  $u_i = \phi(s_i; \theta)$

[image based??]

Two video sequences:  $s, t$   $\checkmark_N$   $\downarrow$   $m$  length.

Embedding  $U = \{u_1, u_2, \dots, u_m\}, V = \{v_1, v_2, \dots, v_m\}$

Cycle Consistency:  $u_i \in U$

$$\text{find } v_i = \arg \min_{v_j \in V} \|u_i - v_j\| \quad \begin{matrix} \text{Nearest} \\ \text{neighbor} \\ \text{of } u_i \end{matrix}$$

$$\text{Repeating } u_k = \arg \min_{v_l \in V} \|v_l - u_k\| \quad \begin{matrix} \text{Nearest} \\ \text{neighbor} \\ \text{of } v_i \end{matrix}$$

[The points are cycle consistent iff  $i = k$ ]

Cycle back classification:

$$\text{Soft Nearest Neighbor: } \tilde{v} = \sum_{j=1}^M \alpha_j v_j \quad \begin{matrix} \text{input point} \\ \text{track it to } U \text{ set.} \end{matrix}$$

$$\text{Softmax weight: } \alpha_j = \frac{e^{-\|u_i - v_j\|^2}}{\sum_{k=1}^M e^{-\|u_i - v_k\|^2}}$$

$$N \text{ class classification problem: Logit } x_k = \frac{-\|\tilde{v} - u_k\|_2^2}{\sqrt{N}}$$

$$\hat{y} = \text{softmax}(x_k)$$

{ the smaller the better }

(ii)

Cross-Entropy loss: *only the  $y_i \in \pi$*

$$L_{CE} = -\sum_j^N y_i \log(\hat{y}_j)$$

Cycle back-propagation:

Proximity similarity vectors

$$\beta_k = \frac{\|\tilde{v} - u_k\|^2}{\sum_{j=1}^K \|\tilde{v} - u_j\|^2}$$

Variance

regularization

peaky around  $i$

$$L_{CBP} = \frac{\|i - \mu\|^2}{\sigma^2} + \lambda \log \sigma$$

where,

$$\left\{ \begin{array}{l} \mu = \frac{1}{K} \sum_{k=1}^K \beta_k \\ \sigma^2 = \frac{1}{K} \sum_{k=1}^K \beta_k \cdot (k - \mu)^2 \end{array} \right.$$

The losses are differentiable

Backpropagation

①

## ④ The whitening loss The whitening loss

$$\text{Embedding } z = f(x, \theta)$$

Goal,

$$\begin{cases} \min_{\theta} E \left[ \text{dist}(z_i, z_j) \right] \\ \text{s.t. } \text{cov}(z_i, z_i) = \text{cov}(z_j, z_j) = I \end{cases}$$

$$\begin{aligned} \text{dist}(z_i, z_j) &= \left\| \frac{z_i}{\|z_j\|_2} - \frac{z_j}{\|z_i\|_2} \right\|^2 \\ &= 2 - 2 \frac{\langle z_i, z_j \rangle}{\|z_i\|_2 \|z_j\|_2} \end{aligned}$$

Original Image No. N

$$\text{Batch } B = \{x_1, x_2, \dots, x_K\}$$

$K = \frac{N}{d}$   $\rightarrow$  No of positives.

$$V = \{v_1, v_2, \dots, v_K\}$$

Proposed W-MSE :

whitening ( $v$ )

$$L_{W\text{-MSE}}(v) = \frac{2}{N^2 d(d-1)} \sum \text{dist}(z_i, z_j); (i, j) \in \text{positive}$$

$$\text{whitening } (v) = w_v^\top (v - \mu_v)$$

$$\mu_v = \frac{1}{K} \sum_k v_k$$

$$w_v^\top w_v = \Sigma_v^{-1} \quad \text{Inverse of covariance Matrix.}$$

$$\Sigma_v = \frac{1}{K-1} \sum_k (v_k - \mu_v)(v_k - \mu_v)^\top$$

(P)

Memory bank

Memory Bank

①

Parametric Classifier:n images  $x_1, \dots, x_n$  in n classes.feature  $v_j = \frac{f_j(x_j)}{\|f_j(x_j)\|}$ 

$$P(i|v) = \frac{\exp(v^T w_i)}{\sum_{j=1}^n \exp(v^T w_j)}$$

weight vector for class i

Non Parametric classifier:

the instances itself ??

$$P(i|v) = \frac{\exp(v^T v_i)}{\sum_{j=1}^n \exp(v^T v_j)}$$

Learning Objective:

Augment version to itself ??

$$J(\theta) = -\sum_{i=1}^n \log P(i|f_\theta(x_i)) \quad // \text{minimize}$$

Equivalent to:  $\prod_{i=1}^n P_\theta(i|f(x_i)) \quad // \text{maximization}$ 

$$\text{NCE: } P(i|v) = \frac{\exp(v^T f_i / \sigma)}{Z_i}$$

$$Z_i = \sum_{j=1}^N \exp(v_j^T f_i / \sigma)$$

$$h(i, v) := P(D=1 | i, v) := \frac{P(i|v)}{P(i|v) + m P_d(i)}; P_d = \frac{1}{n}$$

↓  
m times more  
noise sample

$$\mathcal{J}_{\text{NCE}}(v) = -E_{p_i} [\log h(i, v)]$$

$$= -m E_{p_i} [\log (1 - h(i, v'))]$$

other image feature

$$z = z^* = n E_j \left[ \exp \left( v_j^T f_i / c \right) \right] = \frac{n}{m} \sum_{k=1}^m \exp \left( v_k^T f_i / c \right)$$

↑ Added for regularization

Proximal Regularization:

Loss function:  $-\log h(i, v_i^{(t-1)}) + \lambda \|v_i^{(t)} - v_i^{(t-1)}\|_2^2$

①

## PJRL

PIRL (pretext invariant representation learning)

Notation: Dataset  $\mathcal{D} = \{I_1, I_2, \dots, I_D\}$

$$I_n \in \mathbb{R}^{W \times H \times 3}$$

Transformation set,  $\mathcal{T}$

Network  $\phi_\theta(\cdot)$

$$v_I = \phi_\theta(I)$$

target optimization:

$$L_{\text{inv}}(\theta, \mathcal{D}) = \mathbb{E}_{t \sim P(T)} \left[ \frac{1}{|\mathcal{D}|} \sum_{I \in \mathcal{D}} L(v_i, v_{I+t}) \right]$$

contrastive target:

$$L_{\text{co}}(\theta, \mathcal{D}) = \mathbb{E}_{t \sim P(T)} \left[ \frac{1}{|\mathcal{D}|} \sum_{I \in \mathcal{D}} L_{\text{co}}(v_i, z(t)) \right]$$

measure properties  
of transformation +

Keep semantic  
irrelevant info.

Loss function:

$$h(v_i, v_{I+t}) = \frac{\exp\left(\frac{s(v_i, v_{I+t})}{c}\right)}{\exp\left(\frac{s(v_i, v_{I+t})}{c}\right) + \sum_{I' \in \mathcal{D}_N} \exp\left(\frac{s(v_i, v_{I'})}{c}\right)}$$

$$\begin{aligned} L_{\text{nce}}(I, I^t) &= -\log \left[ h(f(v_i), g(v_{I^t})) \right] \\ &\quad - \sum_{I' \in \mathcal{D}_N} \log \left[ 1 - h(g(v_{I'}), f(v_i)) \right] \end{aligned}$$

final loss function:

$$L(x, x^t) \rightarrow L_{\text{CNS}}(m_i, g(v_{it})) \quad \text{step 2.} \\ + (1-\alpha) L_{\text{ENSE}}(m_i, f(v_i)) \quad \text{step 1.} \quad \left. \right\} \text{kinda two step}$$

$\downarrow$

①  $f(v)$  similar to  $m_i$  // Dampening  
②  $f(v_i), f(v_i')$  dissimilar.

## ① EqCo EqCo (Equivalent Contrastive)

$$L_{\text{EqCo}} = \mathbb{E}_{q \sim D, k_+ \sim D^+, k_- \sim D^-} \left[ -\log \frac{e^{(q^T k_+ - m)/\tau}}{e^{(q^T k_+ - m)/\tau} + \sum_{i=1}^n e^{q^T k_i/\tau}} \right]$$

↑ Margin  
↑ Temperature

$q \rightarrow$  Query Sample Representation.

$k \rightarrow$  keys.  $0 \rightarrow$  positive;  $k_- \rightarrow$  negative.

$D \rightarrow$  Data Distribution.

Equivalent rule:  $m = \tau \log \frac{\alpha}{K}$

↑ constant !!  
No of Negatives.  
Temperature.  
Margin

EqCo: Batch size: Number of queries in 'N' per batch.

Negative Samples: K per query

Query Embedding ' $q$ '

Key Embedding  $x = \{x_i\}_{i=0, \dots, K}$

conditional  $P(x_i | q)$

independent dist  $P(x_i)$

$$P_{H_i}(x) = P(x_i | q) \prod_{j \neq i} P(x_j)$$

$\hookrightarrow \pi_{(k+1)}$  candidate distribution for  $x$ .

$$\begin{aligned}
 \Pr & \left[ x \sim H_0 \left| q_V, x \right. \right] = \frac{\Pr_{H_0}(x)}{\Pr_{H_0}(x) + \Pr \sum_{i=1}^k \Pr_{H_i}(x)} \\
 & = \frac{\frac{\Pr}{\Pr_{H_0}} \frac{\Pr(q_V|x)}{\Pr(q_V)}}{\frac{\Pr}{\Pr_{H_0}} \frac{\Pr(q_V|x)}{\Pr(q_V)} + \sum_{i=1}^k \frac{\Pr(q_V|x)}{\Pr(q_V)}}
 \end{aligned}$$

Let,  $\frac{\Pr(q_V|x)}{\Pr(q_V)} \propto e^{q_V^T k_i / c} \quad (i=0, \dots, k)$

$$\frac{\Pr}{\Pr_{H_0}} = e^{m/c}$$

then,  $L_{opt} \triangleq \mathbb{E}_{q_V, x} - \log \Pr \left[ x \sim H_0 \left| q_V, x \right. \right]$

$$\geq \log \left[ 1 + ke^{m/c} \right] - I(k_0, q_V)$$

$$I(k_0, q_V) \geq f_{\text{bound}}(m, k)$$

$$\begin{aligned}
 & \triangleq \log \left( 1 + ke^{m/c} \right) - L_{opt} \\
 & \approx \log \left( 1 + ke^{m/c} \right) - \mathbb{E}_{q_{WD}, k_0 \sim D(W)} \log \left( 1 + ke^{m/c} \frac{\Pr(k_0)}{\Pr(q_V | k_0)} \right)
 \end{aligned}$$

# Pseudo 3D CNN

①

## ① Pseudo 3D CNN

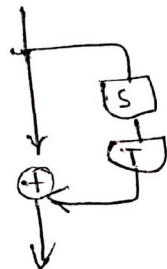
$$\text{Residual units: } \underbrace{x_{t+1}}_{\text{output of } t\text{th block.}} = h(x_t) + F(x_t)$$

↑ identity map      ↗ output of  $L$

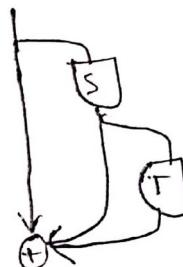
P3D Block design:  $S \rightarrow$  2D spatial filter

$T \rightarrow$  1d temporal filter.

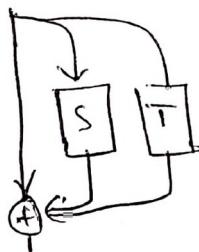
i) P3D - A



ii) P3D C



iii) P3D - B



①

(P) Partially view Alignment RL with noise robust  
 Partially View Alignment RL with Noise robust

$$\{x_i^i\}_{i=1}^N = \{x_1^i, x_2^i, \dots, x_N^i\}_{i=1}^N$$

Partially view aligned data.

$$\{x_i^i\}_{i=1}^N = \{x_i^i, v_i^i\}_{i=1}^N$$

aligned      unaligned

Target Align: utilizing A

Loss function:

$$L = \frac{1}{2m} \sum_{i=1}^N \left\{ p L_i^{pos} + (1-p) L_i^{neg} \right\}$$

$$L_i^{pos} = d(\hat{a}_i^1, a_c^2) \triangleq \|f_1(\hat{a}_i^1) - f_2(a_c^2)\|$$

$$L_i^{neg} = \frac{1}{m} \max \left( m d^2(\hat{a}_i^1, g_j^2) - d^2(\hat{a}_i^1, g_j^2), 0 \right)$$

$$m = \frac{i}{N_p} \cdot \sum d(\hat{a}_i^1, a_i^2) + \frac{1}{N_n} \sum d(\hat{a}_i^1, g_j^2)$$

modified

$$\text{Originally, } L_i^{neg} = \max \left( m - d(\hat{a}_i^1, g_j^2), 0 \right)^2$$