

①

# Barlow Twin

② Barlow twin

Batch of images  $X \begin{cases} y^A \text{ (Augmentation 1)} \\ y^B \text{ (Augmentation 2)} \end{cases}$

$$z^A = f_\theta(y^A) \quad ; \quad z^B = f_\theta(y^B) \quad // \text{mean centered.}$$

$$L_{BT} = \underbrace{\sum_i (1 - c_{ii})^2}_{\text{invariance term}} + \underbrace{\lambda \sum_i \sum_{j \neq i} c_{ij}}_{\text{redundancy reduction term}} \quad \begin{matrix} \text{cross} \\ \text{correlation} \\ \text{matrix} \end{matrix} \quad [\text{proxy entropy}]$$

where,  $c_{ij} = \frac{\sum_b z_{bi}^A z_{bj}^B}{\sqrt{\sum_b (z_{bi}^A)^2 \sum_b (z_{bj}^B)^2}}$    
 $\rightarrow b$  indexes batch sample   
 $\rightarrow i, j$  vector dimension of the network.

$\hookrightarrow$  summation over all the samples.

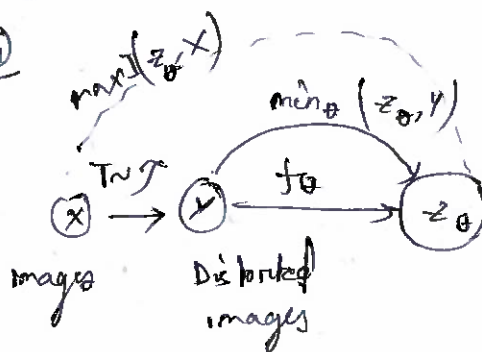
comparison with others:

$$d_{infoNCE} = \underbrace{- \sum_b \frac{\langle z_b^A z_b^B \rangle_i}{\|z_b^A\|_2 \|z_b^B\|_2}}_{\text{similarity term}} + \underbrace{\sum_b \log \left( \sum_{b' \neq b} \exp \frac{\langle z_b^A z_{b'}^B \rangle_i}{\|z_b^A\|_2 \|z_{b'}^B\|_2} \right)}_{\text{contrastive term.}}$$

non parametric estimation of entropy of reps.

④ Barlow twin (Appendix)

Information Bottleneck connection



$$\begin{aligned}
 IB_{\theta} &\triangleq I(z_{\theta}, y) - \beta I(z_{\theta}, x) \quad // \text{close means} \\
 &= \left[ H(z_{\theta}) - H(z_{\theta} | y) \right] - \beta \left[ H(z_{\theta}) - H(z_{\theta} | x) \right] \quad // \text{invariant to } \mathcal{Y} \text{ (Transform)} \\
 &\quad \text{As related to } x \text{ as possible.} \\
 &\quad \text{(can be negative) the smaller the better.} \\
 &\quad \text{A controllable}
 \end{aligned}$$

$$\approx H(z_{\theta} | x) + \frac{1-\beta}{\beta} H(z_{\theta}) \quad \text{if } \beta = 1 \quad // \text{always positive}$$

then <sup>+ve</sup> total positive (Entropy)

$z_{\theta} \sim$  gaussian constraints.

$$IB_{\theta} = E_x \log |C_{z_{\theta} | x}| + \frac{1-\beta}{\beta} \log |C_{z_{\theta}}|$$

Covariance function

connected to IB with some modification

constraints.

# Demystifying CL

## Demystifying Contrastive Learning

Contrastive loss.

$$\mathcal{L}(D, D^+) = - \sum_{(x, x^+) \in D^+} \frac{\exp(f(x)^T f(x^+)/\tau)}{\exp[f(x)^T f(x^+)/\tau] + \sum_{\substack{\tilde{x} \in D \\ x, \tilde{x} \notin D^+}} \exp(f(x)^T f(\tilde{x})/\tau)}$$

Measuring Invariance:

transformation  $t$

invariant function  $h$  iff:  $h(x) = h(t(x))$

formal notion iff  $y(x) = y(t(x))$  label of image  $t(x)$

where,  $t: \mathcal{X} \rightarrow \mathcal{X}$

$$\text{the } h^*(x) = h^*(t(x))$$

invariant for  $t(x)$  & label  $y$

definition of firing unit

$h(x) \in \mathbb{R}^n$  ; fire if  $s_i h_i(x) > t_i$  ;  $s_i \in \{-1, 1\}$

$f_i(x) = \mathbb{1}(s_i h_i(x) > t_i)$  ;  $f(x) \in \mathbb{R}^n$

global firing rate,  $G(i) = E[f_i(x)]$  //  $t_i$  dependency.

$t_i$  chosen such that  $G(i) = \frac{1}{|Y|}$  no of class.

we want  $\Downarrow$

↳ number of firing unit

one class  $\rightarrow$  one section firing;  
equal parts

## ② Demystifying CI

(11)

Local trajectory:  $T(x) = \{t(x, \gamma) \mid \forall \gamma\}$  // set of transformed version of  $x$  image.

Local firing rate is defined as below

$$L_y(i) = \frac{1}{|X_y|} \sum_{z \in X_y} \frac{1}{|T(z)|} \sum_{x \in T(z)} f_i(x) \quad \parallel \quad X_y = \{x \mid x \in X, y(x) = y\}$$

$L_y(i)$  (with arrow to  $i$ th neuron)   
 $|X_y|$  (Avg)   
 $\sum_{z \in X_y}$  (Average in  $X_y$ )   
 $\frac{1}{|T(z)|}$  (Avg)   
 $\sum_{x \in T(z)} f_i(x)$  (measuring local firing for  $x$  & their transformation.)

fraction of time  $i$  neuron fires.

Target conditioned invariance  $I_y(i) = \frac{L_y(i)}{G(i)}$  // find (top-k) neurons.

Representation Invariance Score (RIS):

commonalities in top  $k$  neurons for each classes.

# ① Prototypical Contrastive learning

## Prototypical Contrastive learning

### Preliminaries

$$X = \{x_1, \dots, x_n\} \quad n \text{ images.}$$

$f \rightarrow$  embedding function.

$$X \rightarrow V = \{v_1, \dots, v_n\}$$

$$v_i = f_{\theta}(x_i)$$

$$L_{\text{inference}} = \sum_{i=1}^n -\log \frac{\exp(v_i \cdot v'_i / \tau)}{\sum_{j=0}^r \exp(v_i \cdot v'_j / \tau)}$$

$\downarrow$   
 $r$  negs. includes the  $v'_i$

$$v'_i = f_{\theta'}(x_i)$$

$\theta' \rightarrow$  moving avg of  $\theta$

### PCL w EM:

$$\theta^* = \arg \max_{\theta} \sum_{i=1}^n \log p(x_i | \theta) \quad // \text{maximize log-likelihood.}$$

$$= \arg \max_{\theta} \sum_{i=1}^n \log \sum_{c_i \in C} p(x_i, c_i; \theta)$$

$\downarrow$   
 latent variable  
 law of total prob.

?? How to optimize this ??

$$\sum_{i=1}^n \log \sum_{c_i \in C} p(x_i, c_i; \theta) \geq \sum_{i=1}^n \sum_{c_i \in C} Q(c_i) \log \frac{p(x_i, c_i; \theta)}{Q(c_i)} \quad // \text{ELBO}$$

Focus  
 $\sum_{c_i} Q(c_i) = 1$

(11)

the equality holds 'u

$$Q(c_i) = P(c_i | x_i, \theta) = \frac{P(x_i, c_i | \theta)}{\sum_{c_i} P(x_i, c_i | \theta)}$$

E step:

estimate  $P(c_i | x_i, \theta)$

k means on feature  $v_i' = f_{\theta'}(x_i)$

↓  
momentum encoder.

prototype  $c_i \rightarrow$  centroid of the cluster.

compute  $P(c_i | x_i, \theta) = \mathbb{1}_{(x_i \in c_i)}$

↓  
sharper pdf.

M-step:

$$\sum_{i=1}^n \sum_{c_i \in C} Q(c_i) \log P(x_i, c_i | \theta) = \sum_{i=1}^n \sum_{c_i \in C_i} P(c_i | x_i, \theta) \log P(x_i, c_i | \theta)$$

$$= \sum_{i=1}^n \sum_{c_i \in C} \mathbb{1}_{(x_i \in c_i)} \log P(x_i, c_i | \theta)$$

$$P(x_i, c_i | \theta) = P(x_i | c_i, \theta) P(c_i | \theta) = \frac{1}{K} P(x_i | c_i, \theta)$$

↓  
uniformity assumption

(11)

# ⑦ Prototypical CL

assuming isotropic Gaussian.

$$P(x_i | c_i, \theta) = \exp\left(\frac{-(v_i - c_s)^2}{2\sigma_s^2}\right) \bigg/ \sum_{j=1}^K \exp\left(\frac{-(v_i - c_j)^2}{2\sigma_j^2}\right)$$

readable format ↗

By applying normalization of  $v$  &  $c$  we get.

$$P(x_i | c_i, \theta) = \exp\left(\frac{-(z - z_j^T c_s)^2}{2\sigma_s^2}\right) \bigg/ \sum_{j=1}^K \exp\left(\frac{-(z - z_j^T c_j)^2}{2\sigma_j^2}\right)$$

vary

so maximizing log likelihood falls into.

$$\theta^* = \arg \min_{\theta} \sum_{i=1}^n -\log \frac{\exp[v_i \cdot c_s / \phi_s]}{\sum_{j=1}^K \exp[v_i \cdot c_j / \phi_j]} ; \phi \propto \sigma^2$$

centroid of  $j$  cluster.

In practice the overall objective becomes.

$$\mathcal{L}_{\text{prototype}} = \sum_{i=1}^n \left[ \underbrace{\log \frac{\exp(v_i \cdot v_j^* / \sigma)}{\sum_{j=0}^n \exp(v_i \cdot v_j^* / \sigma)}}_{\text{NCE}} + \exp \frac{1}{m} \sum_{m=1}^m \log \frac{\exp(v_i \cdot c_s^m / \phi_s^m)}{\sum_{j=0}^n \exp(v_i \cdot c_j^m / \phi_j^m)} \right]$$

class prototype.

cluster  $m$  times !!  
with different number  
of cluster ??  
what if 1 is bad ??

(iv)

concentration estimation:  $\phi$  (smaller <sup>(variance)</sup> means high concentration)

$\phi \leftarrow$  momentum features  $\{v_z^1\}_{z=1}^Z$  of same cluster  $c$ .

$$\phi = \frac{\sum_{z=1}^Z \|v_z^1 - c\|_2}{Z \log(Z + \alpha)}$$

should be smaller

smooth params.

scaling factor for  $c_s^m$