

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 θ

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 ??

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

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

(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.

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

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.

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

class prototype.

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

(iv)

concentration estimation: ϕ (smaller ^(variance) 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