

A unifying MI view of metric learning cross-entropy us pairwise losses

Notations:

Dataset $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$

Input feab. X

Embedded features $Z \in \mathbb{R}^d$

Label space $Y \in \mathbb{R}^k$

Encoder $\phi_w: X \rightarrow Z$

Soft classifier $f_\theta: Z \rightarrow [0, 1]^k$

Information measurement:

$$\text{Entropy } H; H(Y) := \mathbb{E}_P[-\log P_Y(y)]$$

$$\text{Conditional Ent. } H(Y|Z) = \mathbb{E}_{P_{YZ}}[-\log P_{Y|Z}(y|z)]$$

$$\text{MI between } Y \text{ & } Z; I(Z; Y) = H(Y) - H(Y|Z)$$

Two views of MI

$$I(Z; Y) = \underbrace{H(Y)}_{\text{Red bracket}} - \underbrace{H(Y|Z)}_{\text{Green bracket}} = H(Z) - H(Z|Y)$$

Discriminative view
[label identification
loss]
CROSS Ent

Generative view
[feature sharing]
loss
PAIRWISE loss

Paige wise loss and generative view of MI

1. Contrastive loss

$$L_{\text{cent}} = \frac{1}{n} \sum_{i=1}^n \underbrace{\sum_{j: y_j = y_i} D_{ij}^2}_{\text{Contrast}} + \frac{1}{n} \sum_{i=1}^n \underbrace{\sum_{j: y_j \neq y_i} [m - D_{ij}]^2}_{\text{Contrast}}$$

Tightness term contrastive term

Hence,

$$\text{Contrast} = \underbrace{H(\hat{z}; \hat{z}|y)}_{\text{cross entropy}} = H(\hat{z}|y) + D_{KL}(\hat{z}||\hat{z}|y)$$

\hat{z} valid only when same label

$\geq H(\hat{z}|y)$

go to unique solution if not given other contrastive term.

we want them to go to unique point

we can show that Contrast is related to $H(\hat{z})$

Finally :

$$L_{\text{cent}} = \frac{1}{n} \sum_{i=1}^n \sum_{j: y_i \neq y_j} (D_{ij}^2 + \alpha m D_{ij}) - \frac{2m}{n} \sum_{i=1}^n \sum_{j=1}^n D_{ij} \alpha - I(\hat{z}; y)$$

$\alpha H(\hat{z}|y)$ $\alpha H(\hat{z})$

Metric learning: cross-entropy vs. pairwise losses

7

Table 2. Several well-known and/or recent DML losses broken into a *tightness* term and a *contrastive* term. Minimizing the cross-entropy corresponds to an approximate bound optimization of PCE.

Loss	Tightness part $\propto \mathcal{H}(\hat{Z} Y)$	Contrastive part $\propto \mathcal{H}(\hat{Z})$
Center [42]	$\frac{1}{2} \sum_{i=1}^n \ z_i - c_{y_i}\ ^2$	$-\frac{1}{n} \sum_{i=1}^n \log p_{iy_i}$
Contrast [7]	$\frac{1}{n} \sum_{i=1}^n \sum_{j:y_j=y_i} D_{ij}^2$	$\frac{1}{n} \sum_{i=1}^n \sum_{j:y_j \neq y_i} [m - D_{ij}]_+^2$
SNCA [43]	$-\frac{1}{n} \sum_{i=1}^n \log \left[\sum_{j:y_j=y_i} \exp \frac{D_{ij}^{\cos}}{\sigma} \right]$	$\frac{1}{n} \sum_{i=1}^n \log \left[\sum_{k \neq i} \exp \frac{D_{ik}^{\cos}}{\sigma} \right]$
MS [40]	$\frac{1}{n} \sum_{i=1}^n \frac{1}{\alpha} \log \left[1 + \sum_{j:y_j=y_i} e^{-\alpha(D_{ij}^{\cos}-m)} \right]$	$\frac{1}{n} \sum_{i=1}^n \frac{1}{\beta} \log \left[1 + \sum_{j:y_j \neq y_i} e^{\beta(D_{ij}^{\cos}-m)} \right]$
PCE Prop. 1	$-\frac{1}{2\lambda n^2} \sum_{i=1}^n \sum_{j:y_j=y_i} z_i^\top z_j$	$\frac{1}{n} \sum_{i=1}^n \log \left[\sum_{k=1}^K \exp \left[\frac{1}{\lambda n} \sum_{j=1}^n p_{jk} z_i^\top z_j \right] \right]$ $-\frac{1}{2K^2 \lambda^2} \sum_{k=1}^K \ c_k^*\ ^2$

Lemma 1. Let T_A denote the tightness part of the loss from method A. Assuming that features are ℓ_2 -normalized, and that classes are balanced, the following hold:

$$T_{SNCA} \stackrel{c}{\leq} T_{Center} \stackrel{c}{=} T_{Contrastive} \stackrel{c}{\leq} T_{MS} \quad (8)$$

Where $\stackrel{c}{\leq}$ stands for lower than, up to a multiplicative and an additive constant.

Notation:



means equal upto additive and
multiplicative constant.