**Ⓕ Theoretical understanding of CL**

# Theoretical Understanding of CL

**framework:**

setup $\begin{cases} \mathcal{X} \to \text{Dataset} \\[6pt] \mathcal{D}_{sim} \sim (x, x^+) \\[6pt] \mathcal{D}_{neg} \sim x_1^-, x_2^-, \ldots x_k^- \\[6pt] \text{Encoder family } f \\[6pt] \therefore f: \mathcal{X} \to \mathbb{R}^d \quad \forall \text{such that,} \quad \|f(\cdot)\| \le R \; ; \; R > 0 \end{cases}$

**latent class:** $(x, x^+) \to$ similar pair

more setup $\begin{cases} \text{family of latent class } \mathcal{C} \\[6pt] \qquad c \in \mathcal{C} \\[6pt] D_c \quad \text{over the } \mathcal{X} \\[6pt] D_c(x) : \text{How relevant } x \text{ to } c \cdot ?? \\[6pt] \rho : \text{how class occurs naturally.} \end{cases}$

**Semantic Similarity:**

$$\mathcal{D}_{sim}(x, x^+) = \underset{c \sim \rho}{E} \; D_c(x) \, D_c(x^+)$$

$$\mathcal{D}_{neg}(x^-) = \underset{c \sim \rho}{E} \; D_c(x^-)$$

**Supervised Tasks:** for a labeled pair $(x, c)$ ; $c \in \{9 \cdots 9_{c+1}\}$

$$\mathcal{D}_{\mathcal{T}}(x, c) = D_c(x) \, D_{\mathcal{T}}(c)$$

**⊕ Theoretical understanding of CL**

**framework:**

$\mathcal{X} \to$ Dataset

**setup** $\left\{ \begin{array}{l} \end{array} \right.$

$D_{sim} \sim (x, x^+)$

$D_{neg} \sim x_1^-, x_2^-, \ldots x_k^-$

Encoder family $f$

$\therefore f: \mathcal{X} \to \mathbb{R}^d$ ∨such that, $\|f(\cdot)\| \leq R$ ; $R > 0$

**Latent class:**

$(x, x^+) \to$ /similar pair

**more setup** $\left\{ \begin{array}{l} \end{array} \right.$

family of latent class $C$

$c \in C$

$D_c$ over the $\mathcal{X}$

$D_c(x)$ : How relevant $x$ to $c.??$

$\rho$ : how class occurs naturally.

**Semantic Similarity:**

$$D_{sim}(x, x^+) = \underset{c \sim \rho}{E} \; D_c(x) \, D_c(x^+)$$

$$D_{neg}(x^-) = \underset{c \sim \rho}{E} \; D_c(x^-)$$

**Supervised 'Tasks'.** for a labeled pair $(x, c)$ ; $c \in \{q \ldots q_{c+1}\}$

$$D_{\vec{\varsigma}}(x, c) = D_c(x) \, D_{\vec{\varsigma}}(c)$$

**Evaluation metric for representation :**

Task $\mathcal{T} = \{c_1 \cdots c_{k+1}\}$

function $g: X \to \mathbb{R}^{k+1}$ // linear classifier.

point $(x, y) \in X \times \mathcal{T}$

loss $\triangleq L(\{g(x)_y - g(x)_{y'}\}_{y \neq y'})$ [Different from true class should be highest]

$f$ (k dim vector of Differences in Coordinate)

By considering standard **hinge loss** :

$$\ell(v) = \max\{0, 1 + \max_i \{-v_i\}\}$$

logistic loss $\ell(v) = \log_2\left(1 + \sum_i \exp(-v_i)\right)$ ; $v \in \mathbb{R}^k$

$$L_{sup}(\mathcal{T}, g) := \mathbb{E}_{(x,c) \sim D_c}\left[\ell\{g(x)_c - g(x)_{c'}\}_{c \neq c'}\right]$$

For linear Classifier: $g(x) = W f(x)$        formally (k+1)

Further, $L_{sup}(\mathcal{T}, f) = \inf\limits_{W \in \mathbb{R}^{(k+1) \times d}} L_{sup}(\mathcal{T}, Wf)$

TP: $W \Rightarrow$ mean for each class representation :

**Mean classifier :** $W^\mu$

c.th row $\Rightarrow \mu_c : \mathbb{E}_{x \sim D_c}[f(x)]$

Now, $L^\mu_{sup}(\mathcal{T}, f) = L_{sup}(\mathcal{T}, W^\mu f)$

④ Theore. understanding of CL

**Avg Supervised loss:**

$$L_{sup}(f) := \mathop{E}_{\{c_i\}_{i=1}^{k+1} \sim \rho^{k+1}} \left[ L_{sup}\left(\{c_i\}_{i=1}^{k+1} , f\right) \mid c_i \neq c_j \right]$$

for mean class

$$L_{sup}^{\mu}(f) := \mathop{E}_{\{c_i\}_{i=1}^{k+1} \sim \rho^{k+1}} \left[ L_{sup}^{\mu}\left(\{c_i\}_{i=1}^{k+1}, f\right) \mid c_i \neq c_j \right]$$

**CL Algorithm:**

unsupervised loss: Population loss:

Neg number
↑
?

$$L_{un}(f) := E \left[ L\left( \{f(x)^T (f(x^+) - f(x_i^-)) \}_{i=1}^k \right) \right]$$

Empirical Counterparts: $\left( x_j, x_j^+, \bar{x}_{j_1}, \dots x_{j_k}^- \right) \in D_{sim} \times D_{neg}^k$  $\quad j=1$

$$\hat{L}_{un}(f) = \frac{1}{m} \sum_{j=1}^{M} L\left( \{f(x_j)^T (f(x_j^+) - f(x_{j_i}^-)) \}_{i=1}^k \right)$$

Now,

$$L_{un}(f) := \mathop{E}_{\substack{c^+, c^- \\ \sim \rho^{k+1}}} \mathop{E}_{\substack{x, x^+ \sim D_{c^+}^2 \\ x_i^- \sim D_{c^-}}} \left[ L\left( \{f(x)^T (f(x^+) - f(x_i^-)) \} \right) \right]$$

# Results and theorems:

**Th. 1.**

$$L_{sup}(\hat{f}) \leq \alpha L_{un}(f) + \zeta \, Gen_M + \delta \quad \forall f \in f_0$$

upper bound?? generalization error.

$$M \to \alpha, \quad Gen_M \to 0$$

$$\alpha, \zeta \to 1, \quad \delta \to 0$$

[if c is large, $L_{un}(f)$ can be small]

$$L_{sup}(\hat{f}) \leq L_{un}^{\mathcal{E}}(f) + \beta s(f) + \zeta \, Gen_M \quad \forall f \in f_0$$

$\rho$ dependent.

$$\rho \to uniform, \; |\mathcal{E}| \to \alpha \text{ then } \beta \to 0, \zeta \to 1$$

Ideal result should be :

$$\boxed{L_{sup}(\hat{f}) \leq \alpha L_{sup}(f) + \zeta \, Gen_M} \quad \forall f \in f_0$$

However not true??