

Non-local Neurons

① Non-local Neurons

Generic Non-local operation $y_i = \frac{1}{c(x)} \sum_{x_j} f(x_i, x_j) g(x_j)$

step 1 ↓
 Normalization ↓
 weighted sum output position

instantiation: Design choice of the f & g

$$g(x_j) = w_g x_j \quad // [1x1 convolution]$$

Gaussian: $f(x_i, x_j) = e^{x_i^T x_j}$

$$c(x) = \sum_{x_j} f(x_i, x_j) \quad \xrightarrow{w_\theta x_i} \xrightarrow{\theta(x_i)^T \phi(x_j)} w_\phi(x)$$

embedded Gaussian: $f(x_i, x_j) = e^{\theta(x_i)^T \phi(x_j)}$

* Generalization of non local

Dot product: $f(x_i, x_j) = \theta(x_i)^T \phi(x_j)$

concentration: $f(x_i, x_j) = \text{ReLU}(w_f^T [\theta(x_i), \phi(x_j)])$

* Wrapping everything in Non local box

$$z_i = w_z(y_i) + x_i \quad // \text{step 2}$$

① Slowfast video Network

Slowfast video Network

	<u>Slow Path</u>	<u>Fast path</u>
①	CNN (3D ResNet)	Another CNN
②	Large temporal stride. $\underline{\alpha}$ (T frames) process	Temporal stride \rightarrow small $\frac{\alpha}{\beta} ; \alpha > \beta$ (αT frames) process
③	Process one out of c frames.	No temporal downsampling
④	Has higher channel high computation	Low no of channels low computation. $B = \frac{1}{\beta} C$ (weakened spatial ability) ($\alpha T, s^2, \beta C$)
⑤	(T, s^2, C)	

Lateral Connection

- 3 ways
- ① Time to channel $\{ \alpha T, s^2, BC \} \xrightarrow{\text{if}} (T, s^2, \alpha BC)$
 - ② Time Strided coupling (T, s^2, BC)
 - ③ Time strided conv. (filter width $5 \times 1^2, 2BC$)

① MultiSiam

MultiSiam

Preliminary: BYOL:

$x \rightarrow \text{image}$

$$v \sim \mathcal{P}_n$$

$$v' \sim \mathcal{P}'_n$$

Online Net: Backbone with GELU, MLP projection, MLP predictor

Target Net: similar to Online
But No MLP predictor.

BYOL \rightarrow stop gradient prevents collapse. (on Target Net)

$$\ell_{ID-img} \triangleq -\cos(q, z') = \frac{\langle q, z' \rangle}{\|q\|_2 \|z'\|_2}$$

BYOL target parameters update

$$\theta \leftarrow \tau \theta + (1-\tau) \theta_t$$

$$\tau \in [0, 1]$$

Positive Samples in multi-instance Data:

Clustering with DL

① Clustering with DL

$$\text{Reconstruction loss : } L = \frac{1}{n} \sum_{i=1}^n \|x_i - f(x_i)\|^2$$

Auto en andere: Non-clusters bres.

cluster losses:

① k-means loss:

(Yang et al. 2006)

$$L(\theta) = \sum_{i=1}^N \sum_{k=1}^K s_{ik} \|z_i - \mu_k\|^2$$

↑ no of cluster

cluster center

↓ embedded Data point

boolean value for cluster assignment
(any soft class option) ??

⑥ Cluster Assignment hardening:

Master Assignment hardening:

$$\sum_{j'} \left(1 + \frac{\|z_i - \mu_j\|^2}{\sigma} \right)^{-\frac{v+1}{2}}$$

constant

softmax value.

for all instances ~~in~~ of class

$$P_{ij} = \frac{q_{ij}^2 / \sum_i q_{ij}}{\sum_j (q_{ij} / \sum_i q_{ij})}$$

A target
Distribution P

(ii) Finally, loss function (for Network to minimize)

$$L = KL(P || Q) = \sum_j \sum_{ij} p_{ij} \log \frac{p_{ij}}{q_{ij}}$$

(iii) Balanced Assignment loss:

(Dizaji et. al 2017)

→ uniform Distribution

$$L_{ba} = KL(\alpha || \mu)$$

where, $q_k = P(y=k) = \frac{1}{N} \sum_i q_{ik}$

(iv) Locality preserving loss: (Huang, 2019)

similarity measure between x_i, x_j

$$L_{lp} = \sum_i \sum_{j \in N_k(i)} s(x_i, x_j) \|z_i - z_j\|^2$$

↓
set of k nearest neighbor of x_i

(v) Group Sparsity loss: (Ng... 2002)

$$L_{gs} = \sum_{i=1}^N \sum_{g=1}^{G_i} \gamma_g \|\phi^*(x_i)\|_1$$

weight of sparsity $\gamma_g = \lambda \sqrt{n_g}$
 ↓
 constant ↓
 group size.

Combined losses: $L(\theta) = \underbrace{\alpha L_c(\theta)}_{\text{clustering loss}} + (1-\alpha) \underbrace{L_n(\theta)}_{\text{Non clustering loss}}$

① Theoretical understanding of CL

Theoretical Understanding of CL

Framework:

$$\left. \begin{array}{l} \mathcal{X} \rightarrow \text{dataset} \\ D_{\text{sim}} \sim (x, x^+) \\ D_{\text{neg}} \sim x_1^-, x_2^-, \dots, x_k^- \\ \text{Encoder family } f \\ \therefore f: \mathcal{X} \rightarrow \mathbb{R}^d \text{ such that, } \|f(\cdot)\| \leq R; R > 0 \end{array} \right\}$$

Latent class: $(x, x^+) \rightarrow$ similar pair

$$\left. \begin{array}{l} \text{family of latent class } C \\ c \in C \\ D_c \text{ over the } \mathcal{X} \\ D_c(x) : \text{How relevant } x \text{ to } c.? \\ \rho : \text{how class occurs naturally.} \end{array} \right\}$$

Semantic Similarity:

$$D_{\text{sim}}(x, x^+) = \underset{c \in P}{\text{E}} D_c(x) D_c(x^+)$$

$$D_{\text{neg}}(x^-) = \underset{c \in P}{\text{E}} D_c(x^-)$$

Supervised Tasks: for a labeled pair (x, c) ; $c \in \{q_1, \dots, q_{C+1}\}$

$$D_{\mathcal{F}}(x, c) = D_c(x) D_c(c)$$

①

④ Theoretical understanding of CL

Framework:

$\mathcal{X} \rightarrow \text{dataset}$

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

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

Encoder family f

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

Latent class: $(x, x^+) \rightarrow$ similar pair

most setup

family of latent class C

$c \in C$

D_c over the \mathcal{X}

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

ρ : how class occurs naturally.

Semantic Similarity:

$$D_{\text{sim}}(x, x^+) = \underset{c \in P}{E} D_c(x) D_c(x^+)$$

$$D_{\text{neg}}(x^-) = \underset{c \in P}{E} D_c(x^-)$$

Supervised Tasks: for a labeled pair (x, c) ; $c \in \{q, \dots, q_{C+1}\}$

$$D_C(x, c) = D_c(x) D_c(c)$$

Evaluation metric for representation:

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

function $g: X \rightarrow \mathbb{R}^{k+1}$ // 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 high]
If k dim vector of difference in coordinate

By considering standard hinge loss:

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

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

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

For linear classification: $g(x) = w^T f(x)$ finally $(k+1)$

$$\text{Further, } L_{\text{sup}}(\mathcal{T}, f) = \inf_{w \in \mathbb{R}^{(k+1) \times d}} L_{\text{sup}}(\mathcal{T}, wf)$$

TP: $w \not\models$ mean for each class representation:

mean classifier: w^k

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

$$\text{Now, } L_{\text{sup}}(\mathcal{T}, f) = L_{\text{sup}}(\mathcal{T}, w^k + t)$$

(11)

④ Theoretical understanding of CL

Aug Supervised loss:

$$L_{\text{sup}}(f) := E_{\substack{\{c_i\}_{i=1}^{k+1} \sim p^{k+1}}} \left[L_{\text{sup}}^{\text{M}}(\{c_i\}_{i=1}^{k+1}, f) \mid c_i \neq g \right]$$

for mean class

$$L_{\text{sup}}^{\text{M}}(f) := E_{\substack{\{c_i\}_{i=1}^{k+1} \sim p^{k+1}}} \left[L_{\text{sup}}^{\text{M}}(\{c_i\}_{i=1}^{k+1}, f) \mid c_i \neq g \right]$$

CL Algorithm:

unsupervised loss: Population loss: Neg number

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

Empirical counterparts: $(x_j^+, x_j^-, \bar{x}_j^+, \dots \bar{x}_j^-) \in \mathbb{D}_{\text{un}}^m \times \mathbb{D}_{\text{neg}}^n$

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

Now,

$$L_{\text{un}}(f) \underset{\substack{C^+, C^- \\ \sim p^{k+1}}}{=} E_{\substack{n, x_i^+ \sim D_{C^+}^2 \\ x_i^- \sim D_{C^-}}} \left[L \left(\{f(x)^T f(x^+) - f(x^-)^T f(x)\} \right) \right]$$

(iv)

Results and theorems:

Th. 1.

$$L_{\text{sup}}(\hat{f}) \leq \alpha L_{\text{un}}(f) + \gamma \underline{\text{Gen}}_m + \delta \quad \forall f \in \mathcal{F}$$

upper bound?? generalization error.

 $M \rightarrow \infty, \underline{\text{Gen}}_m \rightarrow 0$ $\alpha, \gamma \rightarrow 1, \delta \rightarrow 0$

If c is large, $L_{\text{un}}(f)$ can be small

$$L_{\text{sup}}(\hat{f}) \leq L_{\text{un}}^*(f) + B_S(f) + \gamma \underline{\text{Gen}}_m + \delta \quad \forall f \in \mathcal{F}$$

\Rightarrow dependent.

$\rho \rightarrow$ uniform, $|e| \rightarrow \infty$ then $\beta \rightarrow 0, \gamma \rightarrow 1$

Ideal result should be:

$$\boxed{L_{\text{sup}}(\hat{f}) \leq \alpha L_{\text{sup}}(f) + \gamma \underline{\text{Gen}}_m} \quad \forall f \in \mathcal{F}$$

However not true?!