# ⑧ Demystifying CL

Contrastive loss.

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

## Measuring Invariance:

transformation $t$

Invariant function $h$ iff. $\boxed{h(x) = h(t(x))}$

$\quad\longrightarrow$ label of image $t(x)$ $(x)$

formal Notion  iff  $y(x) = y(t(x))$

the  $\boxed{h^*(x) = h(t(x))}$  ; where, $t : x \to x$

$\Updownarrow$

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

## Definition of firing unit

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

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

$\curvearrowleft$ ith neuron

Global firing rate,  $G(i) = E\{f_i(x)\}$ // $t_i$ dependency.

$t_i$ chosen such that  $\oint G(i) = \dfrac{1}{|y|} \longrightarrow$ no of class.

wewant $\Downarrow$

is numbers of firing unit

one class $\to$ one section firing,

equal parts

Ⓟ Demystifying CL

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

Local firing rate is defined as below

ith neuron

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

↓
fraction
of time
i neuron fires.

Avg — All image in $X_y$ | thy | measuring local firing for $x$ & their transformation.

Target conditioned invariance $J_y(i) = \dfrac{L_y(i)}{G(i)}$  //

$\Big\}$ find (top-k) neurons.

Representation Invariance Score (RIS):

commonalities in top k neurons for each classes.