

①

Unsupervised Data Augmentation (UDA)

② Unsupervised Data Augmentation (UDA)

UDA: Target model $P_{\theta}(y|x)$

$P_L(x)$ // Labeled data dist.

$P_U(x)$ // Unlabeled data dist.

perfect model f^*

Supervised Augmentation: $\hat{x} \sim q(\hat{x}|x)$

ADA: input x : $P_{\theta}(y|x) \xrightarrow{\text{minimize}} D(P_{\theta}(y|x) || P_{\theta}(y|x, \epsilon))$

Divergence

$P_{\theta}(y|x, \epsilon)$

noise

quality??.

$$\hat{x} = q(x, \epsilon)$$

supervised

Total Objective

$$\begin{aligned} \min_{\theta} J(\theta) = & \left\{ E_{x_i \sim P_L(x)} \left[-\log P_{\theta}(f^*(y) | x_i) \right] \right. \\ & + \lambda E_{x_c \sim P_U(x)} E_{\hat{x} \sim q(\hat{x}|x_c)} \left[CE \left(P_{\theta}(y|x_c) || P_{\theta}(y|\hat{x}) \right) \right] \end{aligned}$$

fixed copy (no grad) updated net only

↓

unsupervised

Sharpening Indicator: indicator

$$\frac{1}{|B|} \sum_{x \in B} I(\max_y P_{\theta}(y|x) > \beta) CE \left(P_{\theta}^{\text{sharp}}(y|x) || P_{\theta}(y|\hat{x}) \right)$$

$$P_{\theta}^{\text{sharp}}(y|x) = \frac{\exp(\frac{-y}{\epsilon})}{\sum_{y'} \exp(\frac{-y'}{\epsilon})}$$

logit label for y

Theory:

In-domain: $P_U(\hat{x}) > 0$ for $\hat{x} \sim q(\hat{u}|x)$, $x \sim P_U(x)$

Label preserving: $f^* \circ q = f^*(\hat{x})$ for $q(\hat{u}|x) ; x \sim P_U(x)$

Reversible: if $q(\hat{x}|x) > 0$; then $q(x|\hat{x}) > 0$

Theorem: Under UDA, $\text{Pr}_U(A)$: Algo. can't infer the label

of new test example from labeled example from P_L

→ geometric (succes after m try)

$$\text{Pr}_U(A) = \sum_i P_U \left(1 - p_i\right)^m \quad // \text{Prob bound.}$$

$$\downarrow = \sum_{x \in C_i} p_L(x)$$

// labeled component.

$p_i \Rightarrow$ observed example fall in i -th component

component numbers

Further, if $m = O\left(\frac{1}{\epsilon}\right) \Rightarrow \text{Pr}_U(A) = O(\epsilon)$

Error Rate..

① SAM SAM

Training Dataset $S = \bigcup_{i=1}^n \{(x_i, y_i)\}$ from \mathcal{D} i.i.d

model parameters, $w \in W \subseteq \mathbb{R}^d$

Data point loss function $L: w \times \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}_+$

$$\text{Training loss: } L_S(w) \triangleq \frac{1}{n} \sum_{i=1}^n L(w, x_i, y_i)$$

$$\text{Population loss: } L_D(w) \triangleq \mathbb{E}_{(x, y) \sim \mathcal{D}} [L(w, x, y)]$$

Theorem: $\rho > 0$, w.h.p. over training set S , from \mathcal{D}

$$L_D(w) \leq \max_{\|\epsilon\|_2 \leq \rho} L_S(w + \epsilon) + h\left(\frac{\|w\|_2^2}{\rho^2}\right)$$

regularization
||w||₂
can somehow
do it ??

strictly
increasing function ($R_+ \rightarrow R_+$)

Rewriting conditioned on L_D

$$\left[\max_{\|\epsilon\|_2 \leq \rho} L_S(w + \epsilon) - L_S(w) \right] + L_S(w) + h\left(\frac{\|w\|_2^2}{\rho^2}\right)$$

sharpness more connected
+ the bound

So SAM optimization problem.

$$\min_w L_S^{\text{SAM}}(w) + \lambda \|w\|_2^2 \quad // \quad L_S^{\text{SAM}}(w) = \max_{\|\epsilon\|_p \leq \rho} L_S(w + \epsilon)$$

λ norm $\in [1, \infty]$

finding L_S^{SAM} via optimization

$$\epsilon^*(w) \triangleq \arg \max_{\|\epsilon\|_p \leq \rho} L_S(w + \epsilon) \approx \arg \max_{\|\epsilon\|_p \leq \rho} L_S(w) + \epsilon^T \nabla_w L_S(w)$$

$$= \arg \max_{\|\epsilon\|_p \leq \rho} \epsilon^T \nabla_w L_S(w)$$

Solution involves math.

(11)

Elementwise operation

$$\hat{\epsilon}(w) = \text{sign}(\nabla_w L_s(w)) \left[\nabla_w L_s(w) \right]^{q-1} \left(\frac{\|\nabla_w L_s(w)\|_q^2}{\text{norm}} \right)^{\frac{1}{q}}$$

where $\frac{1}{p} + \frac{1}{q} = 1$

Now,

$$\nabla_w L_s^{\text{STM}}(w) \approx \nabla_w L_s(w + \hat{\epsilon}(w)) \quad \text{at the point } w + \hat{\epsilon}(w)$$
$$= \frac{d(w + \hat{\epsilon}(w))}{dw} \nabla_w L_s(w) \Big|_{w + \hat{\epsilon}(w)}$$

$$= \nabla_w L_s(w) \Big|_{w + \hat{\epsilon}(w)} + \underbrace{\frac{d\hat{\epsilon}(w)}{dw} \nabla_w L_s(w) \Big|_{w + \hat{\epsilon}(w)}}_{\text{Dropped it!}}$$

for Acceleration : Dropping the Second term.

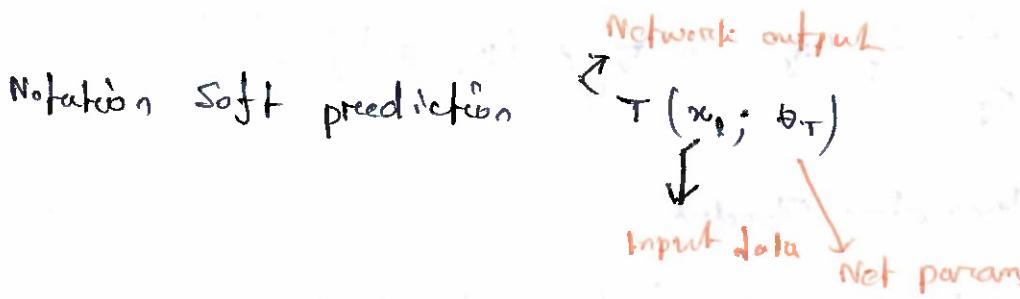
$$\nabla_w L_s^{\text{STM}}(w) \approx \nabla_w L_s(w) \Big|_{w + \hat{\epsilon}(w)}$$

①

① meta pseudo labels: Meta pseudo-labels

Teacher, $T \rightarrow \theta_T$ data labeled (x_L, y_L)

Student, $S \rightarrow \theta_S$



Pseudo label optimization :- (review)

$$\theta_S^{PL} = \arg \min_{\theta_S} E_{x_u} \left[CE \left(T(x_u; \theta_T), S(x_u; \theta_S) \right) \right] \\ := L_u(\theta_T, \theta_S)$$

$$E_{x_L, y_L} \left[CE(y_L, S(x_L; \theta_S^{PL})) \right] := L_L(\theta_S^{PL}) \quad // \text{should be low}$$

\downarrow

is a function of (θ_T)

further,

$$\min_{\theta_T} L_L(\theta_S^{PL}(\theta_T))$$

$$\text{where } \theta_S^{PL}(\theta_T) = \arg \min_{\theta_S} L_u(\theta_T, \theta_S)$$

pseudo label adjustment is possible \Rightarrow

New optimization problem unroll everything!!

⑩

Practical Approximation:

$$\theta_s^{PT}(\theta_T) \approx \theta_s - \eta_s \nabla_{\theta_s} L_u(\theta_T, \theta_s)$$

$$\min_{\theta_T} L_u(\theta_s - \eta_s \nabla_{\theta_s} L_u(\theta_T, \theta_s))$$

SGD optimization Objective:

$$\theta_s' = \theta_s - \eta_s \nabla_{\theta_s} L_u(\theta_T, \theta_s) \quad // \text{1st}$$

$$\theta_T' = \theta_T - \eta_s \nabla_{\theta_T} L_u \left(\theta_s - \nabla_{\theta_s} L_u(\theta_T, \theta_s) \right) \quad // \text{2nd}$$

labeled data unlabeled

① PAWS

PAWS

unlabeled dataset $D = (z_i)_{i \in [1, N]}$

support set $S = \{(z_i, y_i)\}_{i \in [1, m]} \quad m \ll N$

Leverage both $D \& S$

Pretraining D, S twice??
fine tune with S

$z_i \in D \rightarrow \begin{cases} \text{anchor} \\ \hat{z}_i \\ \text{positive} \end{cases}$ } minimize cross entropy between them.

SWAY, BOYL → positive only

Detailed: $x \in \mathbb{R}^{n \times (3H \times W)}$ ↗ view anchor
 $x_d \in \mathbb{R}^{n \times (3H \times W)}$ ↗ positive
 $x_s \in \mathbb{R}^{n \times (3H \times W)}$ } labeled [K total]
 $y_s \in \mathbb{R}^{n \times K}$ ↗ one hot

+ encoder: $\mathbb{R}^{3H \times W} \xrightarrow{f_\theta} \mathbb{R}^d$

$z \in \mathbb{R}^{n \times d}$

$z^+ \in \mathbb{R}^{n \times d}$

$z^- \in \mathbb{R}^{n \times d}$ ↗ label matrix.

j-th row

similarity classifier $\pi_d(z_i, z_j) = \sum_{(z_j, y_j) \in S} \left(\frac{d(z_i, z_j)}{\sum_{z_k \in S} d(z_i, z_k)} \right) y_j$

(ii)

similarity matrices: $d(a, b) = \exp\left(\frac{a^T b}{\|a\| \|b\|}\right)$

$$\pi_i : \pi(z_i; \underline{z}) = \text{softmax}_{\underline{z}}(z_i^T \underline{z}) \quad // \text{softmax prob.}$$

Sharpening function $\left[\rho(p_i)\right]_k := \frac{[p_i]^{1/T}}{\sum_{j=1}^K [p_j]^{1/T}}$; $k = 1, \dots, K$
temperature.
weights sharpening (p_i)

where, $p_i \in [0, 1]^K$

overall objective. for encoder. to minimize

$$\frac{1}{2n} \sum_{i=1}^n \left[H(\rho(p^+), p_i) + H(\rho(p_i), p^+) \right] - H(\bar{p})$$

cross entropy treble? Entropy
why?

Theoretical Bound:

Assumption: balanced class & target sharpening is not uniform.

Prop: Non-collapsing Representation: if rep collapse $\underline{z}_i = \underline{z} \forall i \in S$

then $\|\nabla_{p^+} H(p^+, \underline{z})\| > 0$; gradient is positive?

Proof: if $d(\underline{z}_i, \underline{z}) = d(\underline{z}_j, \underline{z})$

$$\Rightarrow p_i : \pi(\underline{z}, \underline{S}) = \frac{1}{|\underline{S}|} \sum_{i \in \underline{S}} y_i^T = \frac{1}{K} \underset{\text{not uniform}}{\bigcup} \underset{\text{yes}}{\bigcap} p_i \quad / \text{uniform}$$

p^+ not uniform then, $\|\nabla_{p^+} H(p^+, \underline{p})\| > 0$

not uniform uniform

① Maximum Entropy IRL Maximum Entropy IRL

Background: Agent behavior: S_i , trajectory

state: $s_i \rightarrow$ feature $f_{s_i} \in \mathbb{R}^k$

Action: a_i

Goal: Optimizing Some function f_{S_i} to reward value

$$f_S = \sum_{s_i \in S} f_{s_i} \rightarrow \text{reward for all the paths.}$$

$$\underline{\text{Reward}}(f_S) = \theta^T f_S = \sum_{s_i \in S} \theta^T f_{s_i}$$

feature expectation $\Rightarrow \sum_{\text{Path } S_i} P(S_i) f_{s_i} = \bar{f}$ // probabilistic problem.

Deterministic Path distribution: $P(S_i | \theta) = \frac{1}{Z(\theta)} e^{\theta^T f_{s_i}}$

Distribution

[Plans with higher reward is preferred]

$$\begin{aligned} \text{Non deterministic path: } P(S | \theta, T) &= \sum_{s \in S} P(s | \theta) \frac{e^{\theta^T f_s}}{Z(\theta, s)} \mathbb{I}_{S \in \theta} \\ &\approx \frac{e^{\theta^T f_S}}{Z(\theta, T)} \prod_{s_{t+1}, a_t, s_t \in S} P_t(s_{t+1} | a_t, s_t) \end{aligned}$$

$$\text{Stochastic policies: } \pi(\text{action} | \theta, T) \propto \sum_{S, a \in S} P(S | \theta, T)$$

Learning from Demonstration:

$$\theta^* = \arg \max_{\theta} L(\theta) = \arg \max_{\theta} \sum_{\text{examples}} \log P(\tilde{S} | \theta, T)$$

$$\nabla L(\theta) = \bar{f} - \sum_S P(S | \theta, T) f_S = \bar{f} - \sum_{s_i} P_{S_i} f_{s_i}$$

freq.
(Alg 1)

state visualisation
freq.
(Alg 1)

①

P NNCLR

NNCLR

$$\text{infonce: } f_i^{\text{infonce}} = -\log \frac{\exp(z_i \cdot z_i^+/c)}{\exp(z_i \cdot z_i^+/c) + \sum_{z_j \in N_i} \exp(z_i \cdot z_j^-/c)}$$

$$f_i^{\text{simCLR}} = -\log \sum_{k=1}^K \frac{\exp(z_i \cdot z_k^+/c)}{\exp(z_i \cdot z_k^+/c) + \sum_{j \neq i} \exp(z_i \cdot z_j^-/c)}$$

$$\text{NNCLR: } f_i^{\text{NNCLR}} = -\log \frac{\exp(\text{NN}(z_i; \theta) \cdot z_i^+/c)}{\sum_{k \neq i} \exp(\text{NN}(z_i; \theta) \cdot z_k^+/c)}$$

where, $\boxed{\text{NN}(z_i; \theta) = \arg \min_{q \in Q} \|z_i - q\|_2}$ key points.

①

① Perceptual rewards Perceptual Rewards

IRL

$s_t \rightarrow$ visual features activation at time $t \Rightarrow [s_{1t}, s_{2t}, \dots]$

$C \rightarrow \{s_1, \dots, s_T\}$ // sequence of trajectory.

$$P(C) = P(s_1, \dots, s_T) = \frac{1}{Z} \exp \left(\sum_{t=1}^T R_t(s_t) \right); \text{ maxEnt model}$$

↑
unknown rewards
(target)

↓
Boltzmann distribution.

Dynamic programming challenge: How to compute??

Now, Next state, $s_{t+1} = \begin{cases} f(a_t, s_t) \text{ deterministic} \\ \nu P(s_{t+1}|a_t, s_t) \text{ Probabilistic} \end{cases}$

Simplifying assumption,

$$P(C) = \prod_{t=1}^T \prod_{i=1}^N P(s_{it}) = \prod_{t=1}^T \prod_{i=1}^N \frac{1}{Z_{it}} \exp(R_i(s_{it}))$$

whence $\Rightarrow R_t(s_t) = \sum_{i=1}^N R_i(s_{it})$

Intermediate stage discovery:

$$P(C) = \prod_{t=1}^T \prod_{i=1}^N \frac{1}{Z_{it}} \exp(R_{ig}(s_{it}))$$

↓
index of goal/step
at time t

①

② improving MB

Improving MB

given images x_1, \dots, x_n & their encoding $f(x_i) = \phi(x_i, \theta) \in \mathbb{R}^d$

$$L_{CE} = \log \prod_{i=1}^n P(i|x_i) = \sum_{i=1}^n \log \frac{e^{f_i^\top f_i/c}}{\sum_{j=1}^n e^{f_j^\top f_i/c}} \quad || \text{Normal} \quad ①$$

mini large batch

K augmentation for each $x_m \rightarrow \{x_m^{(1)}, x_m^{(2)}, \dots, x_m^{(K)}\}$

$$\text{modified } L_{CE} = \sum_{i=1}^{|B|} \sum_{k=1}^K \log \frac{\exp(f_i^\top \tilde{f}_k/c)}{\sum_{j=1}^n \exp(\tilde{f}_j^\top \tilde{f}_i/c)} \quad || \text{New embed.} \quad ②$$

$$\tilde{f}_i = m \tilde{f}_i + (1-m) \sum_{k=1}^K \frac{1}{K} \tilde{f}_i^{(k)} \quad // \text{Aggregation with some moment.}$$

instance consistency ③

$$\text{consistency loss, } L_{cons} = \sum_{k=1}^K \sum_{i \neq j} KL(P(i|x_i^{(k)}) || P(i|x_j^{(k)})) \quad || \text{eq. ①} \quad || \text{diff same} \quad ④$$

the whole thing is contrastive

classified differently [different embedding looks different]

① Meta learning for Semi-Supervised FSL

① Meta learning for semi-supervised FSL

K-shot N-way episodes.

total K example from N classes.

i) support (training) set $S = \{(x_1, y_1), (x_2, y_2), \dots, (x_{N \times K}, y_{N \times K})\}$

ii) Query/Test set $Q = \{(x_1, y_1), (x_2, y_2), \dots, (x_T, y_T)\}$

Prototype, p_c of class c

$$p_c = \frac{h(x_i) z_{ic}}{\sum z_{ic}}, \quad z_{ic} = \mathbb{1}_{[y_i=c]}$$

$$P(c|x^*, \{p_c\}) = \frac{\exp(-\|h(x^*) - p_c\|_2^2)}{\sum_{c'} \exp(-\|h(x^*) - p_{c'}\|_2^2)},$$

Loss function: to minimize

$$-\frac{1}{T} \sum_i \text{wg}(y_i^* | x_i^*, \{p_c\}),$$

for test data (query)

Prototypical Net with soft K-means: extra terms [softmax]

$$\hat{p}_c = \frac{\sum_i h(x_i) z_{ic} + \sum_j h(\tilde{x}_j) \tilde{z}_{jc}}{\sum_i z_{ic} + \sum_j \tilde{z}_{jc}}, \quad \tilde{z}_{jc} = \frac{\exp(-\|h(\tilde{x}_j) - p_c\|^2)}{\sum_{c'} \exp(-\|h(\tilde{x}_j) - p_{c'}\|^2)}$$

labeled

unlabeled

Refinement prototypes: \tilde{p}_c

Protop. Net with k-means [A distractor class]

Assumption: $P_c = \begin{cases} \frac{\sum_i h(x_i) \cdot z_{ic}}{\sum_i z_{ic}}; & \text{for } c=1 \dots N \\ 0; & \text{for } c=N+1 \end{cases}$

$$\tilde{z}_{j,c} = \frac{\exp\left(-\frac{1}{\sigma_c^2} \| \tilde{x}_j - P_c \|^2 + A(\sigma_c)\right)}{\sum_c \exp\left(-\frac{1}{\sigma_c^2} \| \tilde{x}_j - P_c \|^2 + A(\sigma_c)\right)}$$

length scale
for
distractor
class.

$; A(\sigma) = \frac{1}{2} \log(k\pi) + \log(\sigma)$

Hence this paper, $\pi_1, \dots, \pi_N = 1$
 leave $\pi_{N+1} = ??$

PN soft k-means & masking: $\|h(x_j) - P_c\|^2$

Normalized Distance, $d_{j,c} = \frac{d_{j,c}}{\max_j d_{j,c}}$

soft threshold Id
 $[B_c, \delta_c] = \text{MLP}\left(\min_j(\hat{d}_{j,c}), \max_j(\hat{d}_{j,c}), \text{var}(\hat{d}_{j,c}), \text{skew}(\hat{d}_{j,c}), \text{kurt}(\hat{d}_{j,c})\right)$

$\tilde{P}_c = \frac{\sum_i h(x_i) \cdot z_{ic} + \sum_j h(\tilde{x}_j) \cdot \frac{m_{j,c}}{\delta_c}}{\sum_i z_{ic} + \sum_j \tilde{z}_{j,c} m_{j,c}}$, $m_{j,c} = \sigma \left(\gamma_c (\hat{d}_{j,c} - B_c) \right)$

$\tilde{z}_{j,c}$ a (modified term) sigmoid

modified cluster.

Minimum Entropy Regularization

① minimum Entropy regularization

original Data $L = \{x_i, y_i\}$ $x_i \in X$ data input
 $\downarrow \in \{w_1, \dots, w_K\}$

Some of the label is missing !!

Criterio de derivación:

New learning set $L_n = \{x_i, z_i\}$

$\rightarrow z_{ik} = 1$ if label is provided
 $z_{ik} = 0$ if label is not provided
 $\sum z_{ik} \leq 1$

$z_{ik} = 1$, $k \in 1, \dots, K$

when label is not provided,

\rightarrow means x_i is everybody !!

\rightarrow randomly distributed

for an unlabeled case

$$P(z|x_i, w_k) = P(z|x_i, w_k)$$

$$\neq (w_k, w_l)$$

$$\therefore \text{Now, } P(w_k|x, z) = \frac{\sum_k z_k P(w_k|x)}{\sum_{k=1}^K z_k P(w_k|x)}$$

conditional log likelihood

$$L(\theta, L_n) = \sum_{i=1}^n \log \left(\sum_{k=1}^K z_{ik} f_k(x_i; \theta) \right) + h(z_i),$$

model for $P(w_k|x)$
 \uparrow # parameter (concave)

independent of $P(x, y)$

①

when unlabeled data is informative:

conditional entropy unknown about $y|x, z$ is known

$$H(y|x, z) = - \mathbb{E}_{x,y,z} \log [P(y|x, z)]$$

maximum entropy in true!

$$\mathbb{E}_{\theta, \psi} [H(y|x, z)] = C$$

→ model params. → large enough

$$\therefore P(\theta; \psi) \propto \exp(-\lambda H(y|x, z)) // \text{prior on } \theta$$

$$H_{\text{emp}}(y|x, z; f_n) = -\frac{1}{n} \sum_{i=1}^n \sum_{k=1}^K P(w_k|x_i, z_i) \log P(w_k|x_i, z_i)$$

Entropy Regularization:

$$g_k(x, z; \theta) = \frac{z_k f_k(x; \theta)}{\sum_{k=1}^K z_k f_k(x; \theta)}$$

→ labeled case $g_k(x, z; \theta) \approx z_k$

→ model for $P(w_k|x, z)$

→ unlabeled case $g_k(x, z; \theta) = f_k(x; \theta)$

so the maximizer $\sim P(\theta; \psi) \log P(w_k|x, z) // \text{conditional}$.

$$c(\theta, \lambda; f_n) = L(\theta; f_n) - \lambda H_{\text{emp}}(y|x, z; f_n)$$

$$= \sum_{i=1}^n \log \left(\sum_{k=1}^K z_k f_k(x_i) \right) + \lambda \sum_{i=1}^n \sum_{k=1}^K g_k(x_i, z_i) \log g_k(x_i, z_i)$$

labeled data

unlabeled data

① Neural Turing Machine

Reading: memory $m_t \rightarrow \underset{\downarrow \text{time}}{\text{memory location}}, \underset{\text{matrix}}{M \times M} \text{ matrix}$, vector size at each location.

Attention weight: $0 \leq w_t(i) \leq 1 ; i \in \{1, \dots, m\}$ with constraint $\sum w_t(i) = 1$

read memory $\underset{(1 \times m) \text{ size}}{\sum_i w_t(i) m_t(i)}$ $\xrightarrow{\text{i-th row}}$ unmix combination, weight sum of results, differentiable.

Writing:

$$\begin{aligned} \tilde{m}_t(i) &\leftarrow m_{t-1}(i) \left[1 - w_t(i) e_t \right] \\ &\quad \xrightarrow{\text{pointwise multiply}} \text{erase vector } (1 \times m) \\ &\quad \downarrow \text{after erase add} \\ m_t(i) &\leftarrow \tilde{m}_t(i) + w_t(i) a_t \end{aligned}$$

both differentiable

! constructing weight vector:

focusing by content

$$w_t^c(i) \leftarrow \frac{\exp(\beta_t k [k_t, m_t(i)])}{\sum_j \exp(\beta_t k [k_t, m_t(j)])}$$

key strength, cosine sim, key vector ($k \times m$)

⑦

focusing by location:

↳ interpolation gate ($0, 1$)

$$\tilde{w}_t^g \leftarrow g_{t-1}^{w_t^c} + (1 - g_t) \tilde{w}_{t-1}$$

$$\tilde{w}_t^{(i)} \leftarrow \sum_{j=0}^{N-1} w_t^g(j) s(i-j)$$

shift weight

$$w_t^{(i)} \leftarrow \frac{\tilde{w}_t^{(i)} v_t}{\sum_i \tilde{w}_t^{(i)} v_t}$$

sharpening ↳
normalization

② Prototypical Network

Given N labeled examples $S_N = \{(x_1, y_1), \dots, (x_N, y_N)\}$

$$y_i \in \{1, \dots, K\}$$

$S_k \rightarrow$ Examples only from \underline{k} classes out of K classes.

$$\text{Prototype } c_k = \frac{1}{|S_k|} \sum_{(x_i, y_i) \in S_k} f_\phi(x_i) \quad \begin{matrix} \text{mean vectors.} \\ \text{Network parameters.} \end{matrix}$$

(uses later)

Distance function.

Now,

$$p_\phi(y=k|x) = \frac{\exp(-d(f_\phi(x), c_k))}{\sum_{k'} \exp(-d(f_\phi(x), c_{k'}))} \quad \begin{matrix} \text{probabilistic} \\ \text{approach} \end{matrix}$$

Target to minimize : $J = -\log p_\phi(y=k|x)$

Prototype as mixture density estimation:

$$\text{Bregman Divergence: } d_f(z, z') = f(z) - \left\{ f(z') + \langle \nabla f(z'), z - z' \rangle \right\}$$

(Definition) $= J$ $\begin{matrix} \text{(gradient)} \\ \text{derivative} \end{matrix}$

\downarrow not product

convex function. $\mathbb{R}^n \rightarrow \mathbb{R}$

Exponential family of distributions

$$p_\psi(z|\theta) = \exp\left(z^\top \theta - \psi(\theta) - g_\psi(z)\right)$$

output dist extra term

params. cumulant function

$$\begin{aligned} \text{Bregman Divergence} &= \exp\left(-d_\phi(z, \underline{f_\phi(\theta)})\right) \rightarrow \text{arg min} \\ &\quad - g_\phi(z) \end{aligned}$$

$f_\phi(z) = z$ Network embedding.

(4)

Exponential family with mixture model param: $\Gamma = \{\theta_k, \pi_k\}_{k=1}^K$

$$P(z|\Gamma) = \sum_{k=1}^K \pi_k P_\psi(z|\theta_k)$$

$$= \sum_{k=1}^K \pi_k \exp \left(-d_f(z, \underbrace{\mu(\theta_k)}_{f^T c_k}) - \underbrace{g_f(z)}_{\text{constant term}} \right)$$

proto type
 (low inductive bias)
 (just f(x))

gets canceled
 in prob equation.

$$\Rightarrow P(y=k|z) = \frac{\pi_k \exp(-d_f(z, \mu(\theta_k)))}{\sum_{k'} \pi_k \exp(-d_f(z, \mu(\theta_{k'})))}$$

Connection linear model: for $d_f(z, z') = \|z - z'\|^2$ Euclidean

$$-\|f_\phi(x) - c_k\|^2 = -f_\phi(x)^T f_\phi(x) + 2c_k^T f_\phi(x) - \underbrace{c_k^T c_k}_{w_k^T f_\phi(x) + b_k}$$

linear model term.

①

① ② Lifted structure. Lifted Structure

loss function:

$$J = \frac{1}{2|P|} \sum_{(i,j) \in P} \max(0, J_{i,j})^2$$

↳ positive pair

$$J_{i,j} = \max \left(\max_{\substack{(i,k) \in N \\ i \neq k}} \alpha - D_{i,k}, \max_{\substack{(j,l) \in N \\ j \neq l}} \alpha - D_{j,l} \right) + D_{i,j} + \text{rc}$$

↳ negative pair.

maximize
for what neg pairs.

embedding dimension $m \times c$
 Embedded feature vector $x \in \mathbb{R}^m \xrightarrow{\text{class no}}$

squared norm, $\|\tilde{x}\| = \left[\|f(x_1)\|_2^2, \|f(x_2)\|_2^2, \dots, \|f(x_m)\|_2^2 \right]$

pairwise density matrix $D^2 = \tilde{x} \tilde{x}^T + 1 \tilde{x}^T - 2 \tilde{x} \tilde{x}^T$
 which leads (Efficient computation)
 whose $D_{ij} = \|f(x_i) - f(x_j)\|_2^2$ simple compute.

using upper bound, the loss function.
 should be as close as possible to 0

$$\tilde{J}_{i,j} = \log \left(\underbrace{\sum_{(i,k) \in N} \exp \{ \alpha - D_{i,k} \}}_{(i,k) \in N} + \underbrace{\sum_{(j,l) \in N} \exp \{ \alpha - D_{j,l} \}}_{(j,l) \in N} \right) + D_{i,j}$$

$$J = \frac{1}{2|P|} \sum_{(i,j) \in P} \max(0, \tilde{J}_{i,j})$$

① Relation Network

①

Relation network

Problem definition: Episodic learning: Example

sample set: $S = \{(x_i, y_i)\}_{i=1}^m$ ($m = k \times c$)

query set: $Q = \{(x_i, y_i)\}_{i=1}^n$ (No. of class.)

one shot learning: Relation Network score.

$k \rightarrow$ for each k is k -shot learning.

$$r_{i,j} = g_\phi(C(f_\theta(x_i), f_\phi(x_j))) ; i=1, 2, \dots, c$$

Relation
Networks.

embedding
Networks

query

support

objective function:

$$f_\theta, g_\phi \leftarrow \arg \min_{\theta, \phi} \sum_{j=1}^n \sum_{i=1}^m \{r_{i,j} - \mathbb{1}(y_i = y_j)\}^2$$

① N-pair loss objective

Incorporate multiple negatives $\{x_i, x_i^+\}$, $x_1, \dots, x_{n-1}\}$

$$L(\{x, x^+\}, \{x_i\}_{i=1}^{n-1}; f) = \log \left(1 + \sum_{i=1}^{n-1} \exp(f^T f_i - f^T f^+) \right)$$

↑
pos
↑
neg
↓
proxy
↑
neg
↑
pos
↑
neg as possible

minimize it

$$= -\log \left[\frac{\exp(f^T f^+)}{\exp(f^T f^+) + \sum_{i=1}^{n-1} \exp(f^T f_i)} \right]$$

multiclass logistic loss !!

N-pair loss efficient deep metric learning

$$L_{N\text{-pair-mc}} \left(\{f(x_i), x_i^+\}_{i=1}^N; f \right) = \frac{1}{N} \sum_{i=1}^N \log \left(1 + \exp(f_i^T f_j^+ - f_i^T f_i^+) \right)$$

↑
maximizing / minimizing
what we have.

$$L_{N\text{-pair-ovo}} \left(\{x_i, x_i^+\}_{i=1}^N; f \right) = \frac{1}{N} \sum_{j=1}^N \sum_{i \neq j} \log \left(1 + \exp(f_i^T f_j^+ - f_i^T f_i^+) \right)$$

one-vs-one

①

Few shot meta-learning

few shot meta learning?

Simple view: Optimal model params

$$\theta^* = \arg \min_{\theta} \underset{D \sim P(D)}{E} [L_{\theta}(D)]$$

↓
Dataset

one dataset as one data sample??

$$D = \langle S, B \rangle \quad \begin{matrix} \xrightarrow{\text{Prediction}} \\ [k \text{ shot - N-class classification}] \end{matrix}$$

↓
Learning

Training in the same way as Testing!

$$D = \{(x_i, y_i)\} \in \mathcal{L}^{\text{label}}$$

classifier to

output probabilities for $y | x \rightarrow p_{\theta}(y|x)$

Optimal Parameters

$$\theta^* = \arg \max_{\theta} \underset{(x,y) \sim D}{E} [p_{\theta}(y|x)]$$

$$\theta^* = \arg \max_{\theta} \underset{B \subset D}{E} \left[\sum_{x,y \in B} p_{\theta}(y|x) \right]$$

// Expectation distribution changes.

Few Dataset \rightarrow small support set \rightarrow fake \rightarrow fast learning.

(i)

steps

① subset of labels $L \subset \mathcal{Y}^{\text{label}}$. taking few labels. [2 or 6 labels] may be

② sample support set $s^L \subset D$, training batch $B^L \subset D$

$$y \in L, \forall (x, y) \in s^L, B^L$$

③ support set \rightarrow part of model input.

④ optimization uses mini-batch θ_0

$$(s^L, B^L) \rightarrow \underline{\text{one data point?}}$$

Model trained \rightarrow generalized to other dataset.

$$\theta = \underset{\theta}{\operatorname{argmax}} \quad E_{L \subset \mathcal{Y}^{\text{label}}} \left[E_{s^L \subset D, B^L \subset D} \left[\sum_{(x, y) \in B^L} P_{\theta}(x, y | s^L) \right] \right]$$

[meta learning Parets.]

[optimization to be good at many]

Learner Vs Meta-learner:

Two stage update:

① $f_{\theta} \rightarrow$ learners model.

② update via support set $s^L; \theta' = g_{\theta}(s^L)$

$$E_{L \subset \mathcal{Y}^{\text{label}}} \left[E_{s^L \subset D, B^L \subset D} \left[\sum_{(x, y) \in B^L} P_{g_{\theta}(s^L)}(y | x) \right] \right]$$

① Evolving loss

method : multimodal learning.

$$L = \sum_m \sum_i \lambda_{m,i} l_{m,i} + \sum_d \lambda_d l_d$$

↑ task
↓ weighted [0, 1]
← modality [0, 1]

↑ weighted [0, 1]
constraint
↓ distillation

distillation : $l_d(L_i, M_i) = \|L_i - M_i\|_2$

↓
layer in main network ↓ layer in another network.

Evolving an unsupervised loss function?

① constraint

zipf's distribution matching.

feature $\overset{\in \mathbb{R}^P}{x_{\text{RGB}}} = \text{rgb}(I) \rightarrow \text{cluster into } k$.

$$P(x|c_i) = \frac{1}{\sqrt{2\sigma^2\pi}} \exp\left(-\frac{(x - c_i)^2}{2\sigma^2}\right) \quad // \text{centroid } c_i \in \mathbb{R}^D$$

\downarrow
1

$$P(c_i|x) = \frac{P(c_i) P(x|c_i)}{\sum_j P(x|c_j) P(c_j)} = \frac{\exp(-|x - c_i|^2)}{\sum_{j=1}^k \exp(-|x - c_j|^2)}$$

prior of $P(c_i) = \frac{1/c_i}{H_{k,IS}} \rightarrow \text{real constant}$??

$\hookrightarrow k$ th harmonic number ?? \downarrow law of total prob.

$$KL(P||q) = \sum_i P(c_i) \log \frac{P(c_i)}{q(c_i)} \quad P(c_i) = \frac{1}{n} \sum_{x \in V} P(c_i|x)$$

① what should be contrastive

①

What should be contrastive

view generations:

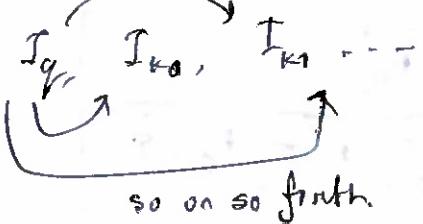
n atomic augmentation.

query view I_q

first key view $I_{k_0} \approx \{q, k_0\} = T\{x_1^{(q, k_0)}, x_2^{(q, k_0)}, \dots, x_n^{(q, k_0)}\}$

n views from reference images $I_{k_i}, \forall i \in \{1, \dots, n\}$

Provides with



contrastive embedding space:

$f(): \mathcal{X} \rightarrow \mathbb{V} \in \mathbb{R}^d$ $\xrightarrow{\text{head}} \mathbb{Z}$
int head each \mathbb{R}^d

MTC setup:

$v^q, v^{k_0}, \dots, v^{k_n} \in \mathbb{R}^d$

projected into $(n+1)$ normalized embedding.

$z_0, z_1, \dots, z_n \in \mathbb{R}^{d'} \text{ by head } h: \mathbb{V} \rightarrow \mathbb{Z}$

$x \rightarrow q$
 k_0
 \vdash
 \vdash
 \vdash

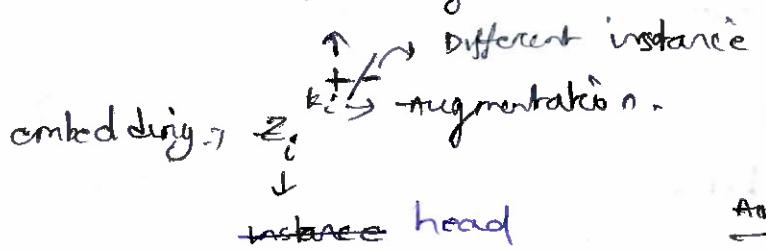
k_0
invariant
to all
(g)

dependent
on
1st transformation

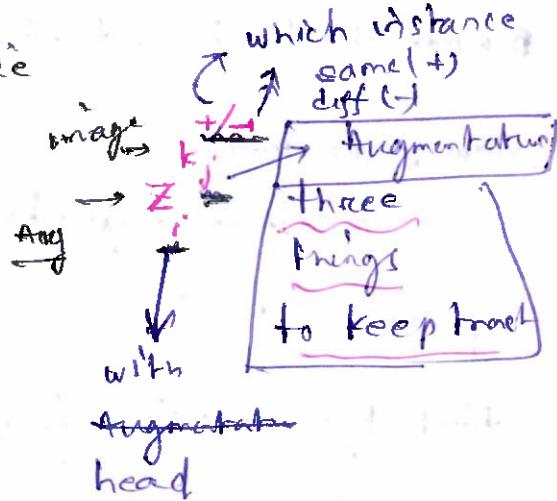
... so on so forth

⑨ what should be

same image ⑩



$$E_{ij}^{\{+, -\}} = \exp \left(\frac{z_i^q \cdot z_j^k f_i^+ - \beta}{\sigma} \right)$$



overall loss

$$L_q = -\frac{1}{n+1} \left[\log \frac{E_{0,0}^+}{E_{0,0}^+ + \sum_{k^-} E_{0,0}^-} \right]$$

$$\boxed{z_i^q}$$

$$-\frac{1}{n+1} \left[\log \frac{E_{i,i}^+}{\sum_{j=0}^n E_{i,j}^+ + \sum_{k^-} E_{i,i}^-} \right]$$

get loss

for each image

add all losses to get total loss

for each image

①

① intriguing Properties of CL loss

Intriguing properties of CL loss

Generalized loss function.

$$\text{L}_{\text{gen CL}} = \text{L}_{\text{align}} + \lambda \text{L}_{\text{distin}}$$

$$\text{L}_{\text{NT Xent}} = -\frac{1}{n} \sum_{i,j} \text{sim}(z_i, z_j) + \frac{\alpha}{n} \sum_{i,j} \log \sum_{k=1}^n \frac{\exp[\text{sim}(z_i, z_k)/\alpha]}{\sum_{k \neq j} \exp[\text{sim}(z_i, z_k)/\alpha]}$$

$$\text{L}_{\text{NT Xent}} = -\frac{1}{n} \sum_{i,j} \log \frac{\exp[\text{sim}(z_i, z_j)/\alpha]}{\sum_{k=1}^n \exp[\text{sim}(z_i, z_k)/\alpha]}$$

By some calculation, log breaking

$$\text{L}_{\text{NT Xent}} = -\frac{1}{n} \sum_{i,j} \text{sim}\left[\exp\left(\frac{\text{sim}(z_i, z_j)}{\alpha}\right)\right] + \frac{\alpha}{n} \sum_{i,j} \log \sum_{k=1}^n \frac{\exp\left(\frac{\text{sim}(z_i, z_k)}{\alpha}\right)}{\sum_{k \neq j} \exp\left(\frac{\text{sim}(z_i, z_k)}{\alpha}\right)}$$

remove it

Alignment

Distribution

[match hidden dist uniform
in hypersphere]

✓ [pairwise potential of gaussian kernel.]

✓ [minimized by perfect uniform encoder]

④ Introducing CL Loss.

(u)

sliced Wasserstein Distance (SWD) loss

Activation vectors. $H \in \mathbb{R}^{b \times d}$, a prior distribution S

Draw prior vector $P \in \mathbb{R}^{b \times d}$ using S

Generate random orthogonal vector matrix $W \in \mathbb{R}^{d \times d'}$ $\boxed{d > d'}$
??

Projection $H^{\perp} = HW \rightarrow P^{\perp} = PW$ // rotation.

Initialize $L = 0$ // SWD

for $j' \in \{1, 2, \dots, d'\}$ do

projection from all.

$$L = L + \left\| \text{sort}(H_{:, j'}^{\perp}) - \text{sort}(P_{:, j'}^{\perp}) \right\|^2$$

end for

\downarrow
column

return L/d'

where it fits ??

① measuring Invariance in DL

①

Measuring invariance in DL

-Measuring Invariance:

firing neuron, $s_i h_i(x) > t_i$

$$s_i \in \{-1, 1\} // \text{choose } s_i \text{ to maximize}$$

$$\text{firing, } f_i(x) = \mathbb{1}_{\{s_i h_i(x) > t_i\}}$$

Transformation function: $\tau(x, y)$

Local trajectory $\tau(x) \rightarrow$ semantically similar stimuli

$$\tau(x) = \{\tau(x, y) | y \in \Gamma\}$$

Global stimuli $G(i) = \mathbb{E}[f_i(x)] // \text{over all possible input}$

Local firing rate $L(i) = \frac{1}{|\mathcal{Z}|} \sum_{z \in \mathcal{Z}} \frac{1}{T(z)} \sum_{x \in T(z)} f_i(x) // \text{only over semantically similar input}$

Invariance score, $s(i) = \frac{L(i)}{G(i)}$