

①

② unsupervised Data Augmentation (UDA)

UDA: Target model $P(y|x)$ $P_L(x)$ // Labeled data dist
 $P_U(x)$ // unlabeled data dist.

Perfect model f^*

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

UDA: input x : $P_\theta(y|x)$ $P_\theta(y|x, \epsilon)$ $\xrightarrow[\text{Divergence}]{\text{minimize}}$ $D(P_\theta(y|x) \| P_\theta(y|x, \epsilon))$
 noise

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

Supervised

Total Objective

$$\min_{\theta} \mathcal{J}(\theta) = \underbrace{E_{x_1 \sim P_L(x)} \left[-\log P_\theta(f^*(x) | x_1) \right]}_{\text{Supervised}} + \lambda \underbrace{E_{x \sim P_U(x)} E_{\hat{x} \sim q(\hat{x}|x)} \left[CE \left(P_{\theta}^{\text{fixed copy (no grad)}}(y|x) \| P_{\theta}^{\text{updated net copy}}(y|\hat{x}) \right) \right]}_{\text{Unsupervised}}$$

Sharpening Prediction: indicator

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

$$P_{\theta}^{\text{sharp}}(y|x) = \frac{\exp(z_y/c)}{\sum_{y'} \exp(z_{y'}/c)} \quad \text{logit label for } \underline{y}$$

Theory:

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

Label preserving: $f^*(x) = f^*(\hat{x})$ 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, $\text{Pr}(A)$: Algo. can't infer the label of new test example from P_L

$$\text{Pr}(A) = \sum_i P_i (1 - P_i)^m \quad \begin{array}{l} \text{geometric (succeed after } m \text{ try)} \\ \text{// Prob bound.} \end{array}$$

$$\downarrow = \sum_{x \in C_i // \text{labeled component.}} P_L(x)$$

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

component number

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

Error Rate..

①

② 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 $L: w \times \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}_+$

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

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

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

hyper
param.

$$L_D(w) \leq \max_{\|e\|_2 \leq \rho} L_S(w+e) + h(\|w\|_2^2 / \rho^2)$$

regularizer

$\|w\|_2$
can somehow
do it !!

strictly
increasing function $(\mathbb{R}_+ \rightarrow \mathbb{R}_+)$
conditioned on L_D

rewriting

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

sharpness

more
connected
to the bound

So SAM optimization problem.

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

finding L_S^{SAM} via optimization

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

(11)

solution involves math:

$$\hat{\epsilon}(w) = \rho \operatorname{sign}(\nabla_w L_S(w)) \left[\nabla_w L_S(w) \right]^{q-1} \left(\underbrace{\|\nabla_w L_S(w)\|_q^2}_{\text{norm}} \right)^{1/p}$$

where $1/p + 1/q = 1$

Elementwise operation

Now,

$$\nabla_w L_S^{\text{sam}}(w) \approx \nabla_w L_S(w + \hat{\epsilon}(w))$$

$$= \frac{d(w + \hat{\epsilon}(w))}{dw} \nabla_w L_S(w) \Big|_{w + \hat{\epsilon}(w)}$$

$$= \nabla_w L_S(w) \Big|_{w + \hat{\epsilon}(w)} + \underbrace{\frac{\partial \hat{\epsilon}(w)}{\partial w} \nabla_w L_S(w) \Big|_{w + \hat{\epsilon}(w)}}_{\text{Dropped !!}}$$

for acceleration : Dropping the Second term.

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

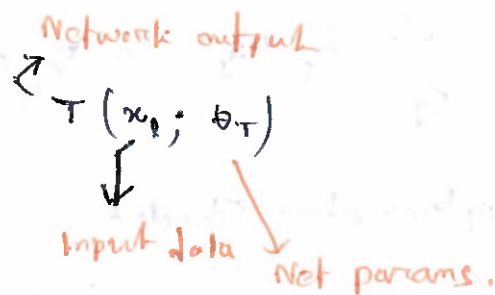
①

Ⓢ meta pseudo labels:

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

Student, $S \rightarrow \theta_S$

Notation Soft prediction



Pseudo label Optimization : (Review)

$$\theta_S^{PL} = \arg \min_{\theta_S} \mathbb{E}_{x_u} \left[\text{CE} \left(T(x_u; \theta_T), S(x_u; \theta_S) \right) \right]$$

$$:= \mathcal{L}_u(\theta_T, \theta_S)$$

$$\mathbb{E}_{x_L, y_L} \left[\text{CE} \left(y_L, S(x_L; \theta_S^{PL}) \right) \right] := \mathcal{L}_L(\theta_S^{PL}) \quad // \text{should be low}$$

↓
is a function of (θ_T)

further,

$$\min_{\theta_T} \mathcal{L}_L(\theta_S^{PL}(\theta_T))$$

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

pseudo label adjustment 'is possible in

new optimization problem without everything!!

②

Practical Approximation:

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

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

SGD optimization Objective:

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

$$\theta_T' = \theta_T - \eta_T \nabla_{\theta_T} L_L(\underbrace{\theta_s - \nabla_{\theta_s} L_u(\theta_T, \theta_s)}_{\text{labeled data}}) \quad // 2nd.$$

unlabeled

①

① PAWS

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

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

Leverage both \mathcal{D} & S pretraining \mathcal{D}, S fine tune with S twice??

swAV, BOYL \rightarrow positive only

$x_i \in \mathcal{D}$ \rightarrow \hat{w}_i (Anchor)
 \rightarrow \hat{w}_i^+ (positive)
 $\left. \begin{array}{l} \hat{w}_i \\ \hat{w}_i^+ \end{array} \right\}$ minimize cross entropy between them.

Detailed:

in hand

$x \in \mathbb{R}^{n \times (3 \times H \times W)}$

$x_d \in \mathbb{R}^{n \times (3 \times H \times W)}$

$x_s \in \mathbb{R}^{m \times (3 \times H \times W)}$

$y_s \in \mathbb{R}^{m \times K}$

\Rightarrow view Anchor

\Rightarrow positive

labeled $[K \text{ total}]$

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

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

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

$\underline{z} \in \mathbb{R}^{S \times d}$

label matrix.

j -th row

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

(1)

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

$$p_i: \pi(z_i, \underline{z}) = \sigma\left(\underline{z}_i^T \underline{z}_s\right) y_s \quad // \text{softmax prob.}$$

sharpening function $[p(p_i)]_k := \frac{[p_i]_k^{1/T}}{\sum_{j=1}^K [p_i]_j^{1/T}}$; $k = 1, \dots, K$

temperature.

weights sharpening $(p_i)_j$

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

overall objective, for encoder, to minimize

$$\frac{1}{2n} \sum_{i=1}^n \left[\underbrace{H(p(p_i^+), p_i)}_{\text{cross entropy}} + \underbrace{H(p(p_i), p^+)}_{\text{the heck? why?}} \right] - \underbrace{H(\bar{p})}_{\text{entropy}}$$

Theoretical Bound:

Assumption: Balanced class & target sharpening is not uniform.

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

then $\|\nabla_b H(p^+, p)\| > 0$; gradient is positive?

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

$$\Rightarrow p: \pi(\underline{z}, S) = \frac{1}{n} \sum_{i=1}^n y_i = \frac{1}{n} \sum_{i=1}^n y_i = \frac{1}{n} \sum_{i=1}^n y_i \quad // \text{uniform}$$

p^+ not uniform then, $\|\nabla H(p^+, p)\| > 0$

not uniform (uniform)

① Maximum Entropy RL

Background: Agent behavior: \bar{S} , / trajectory

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

Action: a_i

Goal: Optimizing some function $f_{s_i} \rightarrow$ Reward value

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

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

feature expectation $\Rightarrow \sum_{\text{path } \bar{S}_i} P(\bar{S}_i) f_{\bar{S}_i} = \bar{f}$ // probabilistic problem.

Deterministic Path distribution: $P(\bar{S}_i | \theta) = \frac{1}{Z(\theta)} e^{\theta^T f_{\bar{S}_i}}$
 Distribution Partition function

Plan with higher reward is preferred

Non deterministic path: $P(\bar{S} | \theta, T) = \sum_{o \in \mathcal{F}} P(o) \frac{e^{\theta^T f_{\bar{S}}}}{Z(\theta, o)} \mathbb{I}_{\bar{S} \in o}$
 $\approx \frac{e^{\theta^T f_{\bar{S}}}}{Z(\theta, T)} \prod_{\bar{s}_{t+1}, a_t, s_t \in \bar{S}} P(\bar{s}_{t+1} | a_t, s_t)$
 Identity green??

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

Learning from Demonstration:

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

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

state visualization freq. (avg 1)

①

① NNCLR

$$\text{infoNCE}; \quad \mathcal{L}_i^{\text{infoNCE}} = -\log \frac{\exp(z_i \cdot z_i^+ / \tau)}{\exp(z_i \cdot z_i^+ / \tau) + \sum_{z^- \in \mathcal{N}_i} \exp(z_i \cdot z^- / \tau)}$$

$$\mathcal{L}_i^{\text{simCLR}} = -\log \frac{\exp(z_i \cdot z_i^+ / \tau)}{\sum_{k=1}^K \exp(z_i \cdot z_k^+ / \tau)}$$

$$\text{NNCLR}: \quad \mathcal{L}_i^{\text{NNCLR}} = -\log \frac{\exp(\text{NN}(z_i; \theta) \cdot z_i^+ / \tau)}{\sum_{k=1}^K \exp(\text{NN}(z_i, z_k) \cdot z_k^+ / \tau)}$$

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

①

② Perceptual rewardsIRL $s_t \rightarrow$ visual feature activation at time $t \Rightarrow [s_{1t}, s_{2t}, \dots]$ $\tau \Rightarrow \{s_1, \dots, s_T\}$ // sequence of trajectory.

$$P(\tau) = P(s_1, \dots, s_T) = \frac{1}{Z} \exp\left(\sum_{i=1}^T R(s_t)\right); \text{ max Ent model}$$

↑ unknown rewards (target)

{ dynamic programming }

challenge: How to compute?

Boltzmann distribution.

$$\text{Now, next state, } s_{t+1} = \begin{cases} f(a_t, s_t) \text{ deterministic} \\ \text{or } P(s_{t+1} | a_t, s_t) \text{ probabilistic.} \end{cases}$$

Simplifying Assumption,

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

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

Intermediate stage discovery:

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

↓
index of goal/step
at time t

①

② Improving MB

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

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

mini large batch

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

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

two sum New embed.

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

Instance consistency (3)

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

eq ① diff
same

the whole thing is contrastive
classified differently [different embedding looks different]

①

① metalearning for semi-supervised FSL

k-shot many episodes.

total k example from
classes.

1) support (training) set $\mathcal{S} = \{(x_1, y_1), (x_2, y_2) \dots (x_{n_{\text{sup}}}, y_{n_{\text{sup}}})\}$

2) Query / test set $\mathcal{Q} = \{(x_1, y_1), (x_2, y_2) \dots (x_T, y_T)\}$

Prototype, p_c of class c

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

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

Loss function: to minimize

argmax
 $y_i = c$

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

for test data (query)

Prototypical Net with Soft K-means: Extra terms [softmax]

$$\hat{p}_c = \frac{\sum_i h(x_i) z_{ic} + \sum_j h(\tilde{x}_j) \tilde{z}_{j,c}}{\sum_i z_{ic} + \sum_j \tilde{z}_{j,c}}; \quad \tilde{z}_{j,c} = \frac{\exp(-\|h(\tilde{x}_j) - p_c\|_2^2)}{\sum_{c'} \exp(-\|h(\tilde{x}_j) - p_{c'}\|_2^2)}$$

labeled

unlabeled

Refinement prototypes: \tilde{p}_c

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

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

$$\tilde{z}_{j,c} = \frac{\exp \left(-\frac{1}{\sigma_c^2} \| \tilde{x}_j - p_c \|^2 - A(\sigma_c) \right)}{\sum_{c'} \exp \left(-\frac{1}{\sigma_c^2} \| \tilde{x}_j - p_{c'} \|^2 - A(\sigma_{c'}) \right)}$$

length scale for distractor class.

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

Here this paper, $\pi_1, \dots, \pi_N = 1$
learn $\sigma_{N+1} = ??$

PN soft k-means & masking:

Normalized Distance, $\tilde{d}_{j,c} = \frac{\|h(x_j) - p_c\|_2^2}{\frac{1}{n} \sum_j d_{j,c}}$

soft threshold \uparrow
 $[\beta_c, \delta_c] = \text{MLP} \left(\left[\min_j (\tilde{d}_{j,c}), \max_j (\tilde{d}_{j,c}), \text{var}_j (\tilde{d}_{j,c}), \text{skew}_j (\tilde{d}_{j,c}), \text{kurt}_j (\tilde{d}_{j,c}) \right] \right)$
 slope \downarrow

$$\tilde{p}_c = \frac{\sum_i h(x_i) z_{ic} + \sum_j h(\tilde{x}_j) \tilde{z}_{j,c} m_{j,c}}{\sum_i z_{ic} + \sum_j \tilde{z}_{j,c} m_{j,c}} ; m_{j,c} = \sigma \left(-\gamma_c (\tilde{d}_{j,c} - \beta_c) \right)$$

\downarrow
 modified cluster.

①

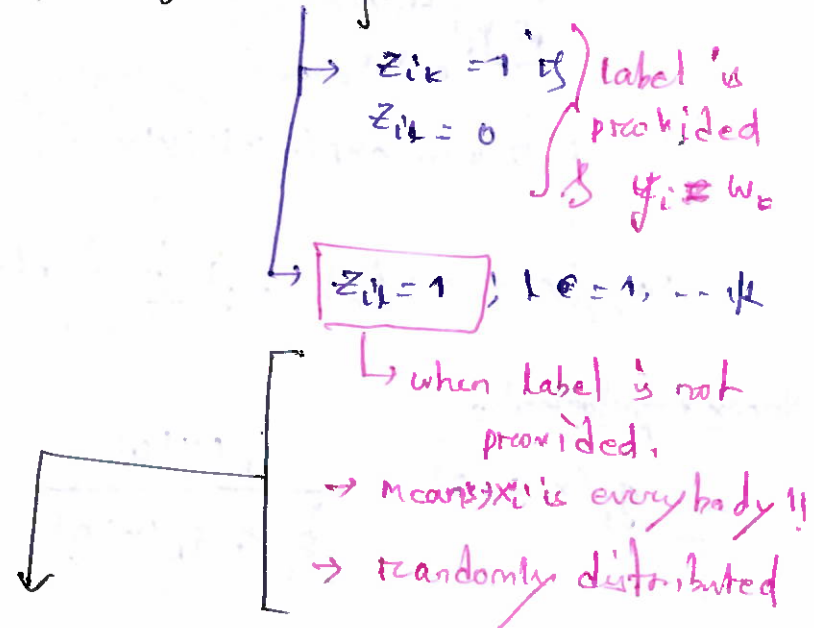
② minimum Entropy regularization

original Data $L_n = \{x_i, y_i\}$ $x_i \in X$ Data input
 $y_i \in \{w_1, \dots, w_K\}$ \rightarrow target output

Some of the label is missing !!

Criterion Derivation:

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



for an unlabeled case $P(z | x_i, w_k) = P(z | x_i, w_l)$
 $\forall (w_k, w_l)$

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

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 # parameters (concave)
 independent of $P(x, y)$

⑦

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 θ :

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

$\underbrace{\theta, \psi}_{\text{model params.}} \rightarrow \text{log-trick}$

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

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

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 $g_k(x, z; \theta) = z_k$

→ model for $p(w_k|x, z)$

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

So the maximizer $\sim p(\theta, \psi) \propto p(w_k|x, z) // \text{conditional.}$

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

$$= \underbrace{\sum_{i=1}^n \log \left(\sum_{k=1}^K z_{ik} f_k(x_i) \right)}_{\text{labeled data}} + \lambda \sum_{i=1}^n \sum_{k=1}^K \underbrace{g_k(x_i, z_i) \log g_k(x_i, z_i)}_{\text{unlabeled data}}$$

labeled data

unlabeled data

① Neural Turing Machine

①

Reading: memory $M_t \rightarrow$ $N \times M$ matrix
 ↓ time
 memory location.
 vector size at each location.

attention weight, $0 \leq w_t(i) \leq 1$; $i \sim \{1, \dots, N\}$ with constraint $\sum_i w_t(i) = 1$

read memory
 $(1 \times M)$ size
 $r_t \leftarrow \sum_i w_t(i) m_t(i)$
 i th row
 vector combination.
 weighted sum of rows.
 differentiable.

Writing:

$\tilde{m}_t(i) \leftarrow m_{t-1}(i) \left[1 - w_t(i) e_t \right]$
 erase vector $(1 \times M)$
 pointwise multiply
 after erase add
 $m_t(i) \leftarrow \tilde{m}_t(i) + w_t(i) a_t$
 both differentiable

! constructing weight vector:

focusing by content

key strength.
 cosine sim
 key vector $(k \times M)$
 $w_t^c(i) \leftarrow \frac{\exp(A_t k [k_t, m_t(i)])}{\sum_j \exp(A_t k [k_t, m_t(j)])}$

(17)

focusing by location:

interpolation gate (0, 1)

$$w_t^g \leftarrow g_t w_{t-1}^c + (1 - g_t) w_{t-1}^u$$

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

$$w_t(i) \leftarrow \frac{\tilde{w}_t(i)^{\gamma_t}}{\sum_i \tilde{w}_t(i)^{\gamma_t}} \left[\begin{array}{l} \text{sharpening} \\ \text{normalization} \end{array} \right]$$

① Prototypical Networks

Given N labelled examples $S_N = \{(x_1, y_1) \dots (x_n, y_n)\}$

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

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

Prototype $c_k = \frac{1}{|S_k|} \sum_{(x_i, y_i) \in S_k} f_\phi(x_i)$ // mean vectors.
 (uses later) Network parameters,

Now,

$$P_\phi(y = k | x) = \frac{\exp(-d(f_\phi(x), c_k))}{\sum_{k'} \exp(-d(f_\phi(x), c_{k'}))}$$

// probabilistic approach

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

Prototype as mixture density estimation:

Bregman Divergence: $d_f(z, z') = f(z) - \{f(z') + \underbrace{\langle \nabla f(z'), z - z' \rangle}_{\text{dot product}}\}$
 (Definition) convex function, $\mathbb{R}^n \rightarrow \mathbb{R}$

Exponential family of Distributions

$$P_\psi(z | \theta) = \exp(z^T \theta - \psi(\theta) - g_\psi(z)) = \exp(-d_\phi(z, \mu(\theta)))$$

output θ dot extra term Bregman Divergence
params. Params. cumulant function neg !!
 $f_\phi(x) = -x$ Network embedding.

(ii)

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

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

$$= \sum_{k=1}^K \pi_k \exp \left(-d_f(z, \mu(\theta_k)) - g_f(z) \right)$$

 $\int c_k$

prototype

(low inductive bias)
(just Avg)constant term
gets cancelled
in prob equation.

$$\Rightarrow p(y=k|z) = \frac{\pi_k \exp(-d_f(z, \mu(\theta_k)))}{\sum_{k'=1}^K \pi_{k'} \exp(-d_f(z, \mu(\theta_{k'})))}$$

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

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

↓
linear model/term.

①① Lifted structure.

loss function:

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

positive pairs

$$\tilde{J}_{i,j} = \max_{(i,k) \in \hat{N}} \left(\alpha - D_{i,k} \right) - \max_{(j,l) \in \hat{N}} \left(\alpha - D_{j,l} \right) + D_{i,j}$$

maximize for what neg pairs.

Distance between +ve

negative pairs.

Embedded feature vector $x \in \mathbb{R}^{m \times c}$ $\xrightarrow{\text{embedding dimension}}$ 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)

where $D_{i,j} = \|f(x_i) - f(x_j)\|_2^2$

simple compute.

using upper bound, the loss function.

should be as close as possible now

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

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

① Relation Network

Problem definition: Episodic learning: Example

sample set: $S = \{(x_i, y_i)\}_{i=1}^m$ ($m = k \times c$)
no. of class.

query set: $Q = \{(x_i, y_i)\}_{i=1}^n$

one shot learning: Relation Network score.

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

$$r_{ij} = g_{\phi} \left(C \left(f_{\psi}(x_i), f_{\phi}(x_j) \right) \right); \quad i=1, 2, \dots, c$$

Relation Networks | Embedding Networks

query support

Objective function:

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

⑦

① N-pair loss objective

incorporate multiple negatives $\{x, x^+, x_1, \dots, x_{n-1}\}$
 pos ↑
 ↓ quocry
 negs

$$L(\{x, x^+, \{x_i\}_{i=1}^{n-1}; f) = \log \left(1 + \sum_{i=1}^{n-1} \exp \left(\underbrace{f^T f_i - f^T f^+}_{\substack{\uparrow \text{neg} \quad \uparrow \text{pos} \\ \text{neg as possible}}} \right) \right)$$

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?

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

maximizing what we have. // multiclass

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

// one-vs-one

①

few shot meta learning?

simple view: Optimal model params

$$\theta^* = \arg \min_{\theta} \underbrace{E_{D \sim P(D)}}_{\text{Dataset}} [L_{\theta}(D)]$$

one dataset as one data sample??

$$D = \langle S, B \rangle \quad \begin{array}{l} \nearrow \text{Prediction} \\ \downarrow \text{Learning} \end{array} \quad [K \text{ shot - } N\text{-class classification}]$$

Training in the same way as Testing!

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

classifier f_{θ}

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

Optimal Parameters

$$\theta^* = \arg \max_{\theta} E_{\underbrace{(x,y) \sim D}} [P_{\theta}(y | x)]$$

$$\theta^* = \arg \max_{\theta} E_{\underbrace{B \sim D}} \left[\sum_{x,y \in B} P_{\theta}(y | x) \right]$$

// Expectation distribution changes.

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

(ii)

steps

(i) subset of labels $L \subset \mathcal{L}^{\text{label}}$. taking few labels. [2 out of 6 classes] maybe

(ii) sample support $s^L \subset \mathcal{D}$, training batch $B^L \subset \mathcal{D}$

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

(iii) support set \rightarrow part of model input.

(iv) optimization uses mini batch B_L

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

model trained to generalized to other dataset.

$$\theta = \underset{\theta}{\operatorname{argmax}} \mathbb{E}_{L \subset \mathcal{L}} \left[\mathbb{E}_{s^L \subset \mathcal{D}, B^L \subset \mathcal{D}} \left[\sum_{(x, y) \in B^L} P_{\theta}(x, y, s^L) \right] \right]$$

[meta learning Params.]

[optimization: to be good at many]

Learner Vs meta-learner:

Two stage updates:

(i) $f_{\theta} \rightarrow$ learners model.

(ii) update via support set s ; $\theta' = g_{\phi}(\theta, s)$

$$\mathbb{E}_{L \subset \mathcal{L}} \left[\mathbb{E}_{s^L \subset \mathcal{D}, B^L \subset \mathcal{D}} \left[\sum_{(x, y) \in B^L} P_{g_{\phi}(\theta, s^L)}(y | x) \right] \right]$$

⑧ Evolving loss

①

method: multimodal learning.

$$L = \sum_m \sum_t \lambda_{m,t} L_{m,t} + \sum_d \lambda_d L_d$$

task

weighted [0,1]

modality

distillation

constraint

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 $x_{P(A,B)} \in \mathbb{R}^D$ $= E_{\text{neg}}(I) \rightarrow$ cluster into k .

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

$\dots \{c_1 \dots c_k\}$

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

prior of $q(c_i) = \frac{1/c_i^s \rightarrow \text{real constant} \quad !! ?}{H_{k,s}}$

$\hookrightarrow k$ th harmonic number !! ?

law of total prob.

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

① what should be contrastive

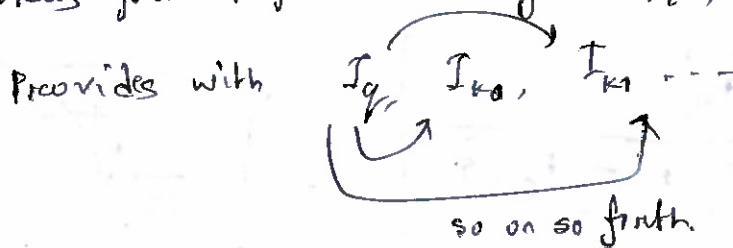
view generations:

n atomic augmentation.

query view I_q

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

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



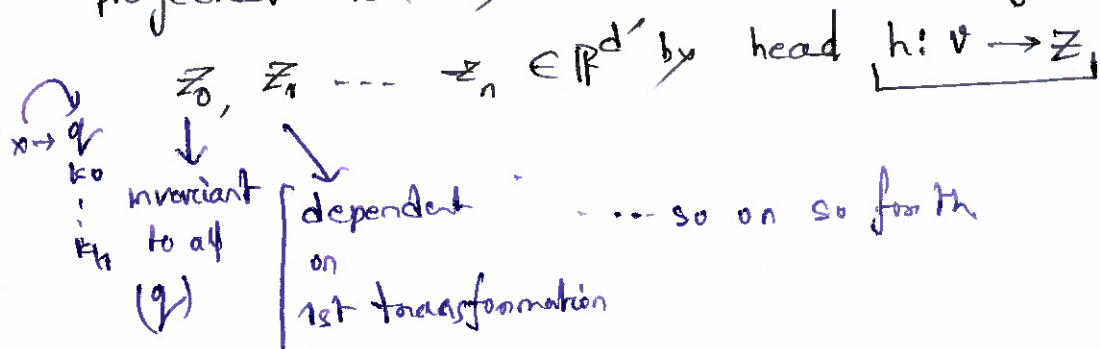
Contrastive Embedding Space:

$$f() : \mathcal{X} \rightarrow v \in \mathbb{R}^d \begin{matrix} \uparrow \downarrow \\ \text{at least each } \mathbb{R}^d \end{matrix}$$

MTE setup:

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

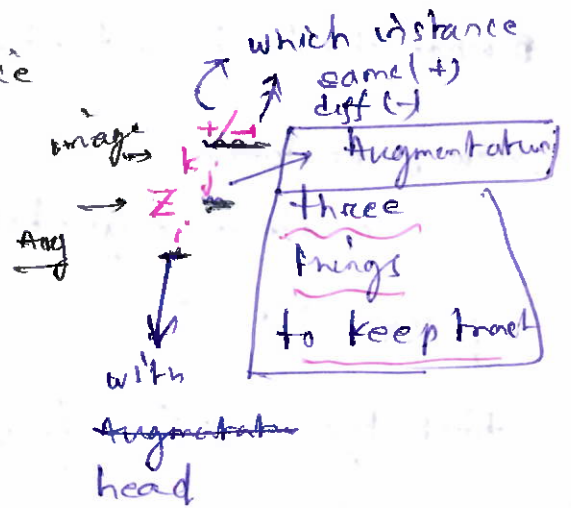
projected into $(n+1)$ normalized embedding.



⑧ what should be same image ⑪

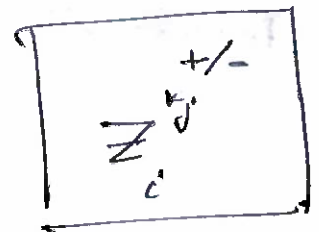
↑ different instance
 $k_i \rightarrow$ augmentation
 embedding $\rightarrow z_i$
 ↓
 instance head

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



overall loss

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



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



Generalized loss function.

$$L_{\text{gen CL}} = L_{\text{augn}} + \lambda L_{\text{dist'n}}$$

$$\tau_{\text{LNT cost}} = -\frac{1}{N} \sum_{i=1}^N \log \sum_{k=1}^N \frac{1}{[k \neq i]} \exp(\text{sim}(z_i, z_k))$$

$$f_{\text{NTXent}} = -\frac{1}{n} \sum_{i,j \in \mathcal{N}} \log \frac{\exp[\sin(z_i, z_j)/\tau]}{\sum_{k=1}^n \mathbb{1}_{k \neq j} \exp[\sin(z_i, z_k)/\tau]}$$

By some calculation, log breaking

by some calculation, log breaking

$$L^{NTXent} = -\frac{1}{n} \sum_{i,j} \text{sim}[\mathbf{z}_i, \mathbf{z}_j] + \frac{1}{n} \sum_{i,j} \log \sum_{k=1}^K \exp(\text{sim}(\mathbf{z}_i, \mathbf{z}_k))$$

remove it

Alignment

Distribution.

[match hidden dist uniform
in hypersphere]

- ✓ [pairwise potential of gaussian kernel.
- ✓ [minimized by perfect uniform encoder]

② 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 $H^T = HW$; $P^T = PW$ // rotation.

initialize $L = 0$ // SWD

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

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

end for

return $L/(d')$

where it fits ??

⑧ measuring Invariance in DL ①

- measuring Invariance:

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

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

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

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

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

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

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

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

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