

Adversarial reciprocal Point learning

Problem :

n labeled examples $\mathcal{D}_L = \{(x_1, y_1), \dots, (x_n, y_n)\}$

N known classes: $y_i \in \{1, \dots, N\}$

test data: $\mathcal{D}_T = \{t_1, \dots, t_u\}$

t_i label $\in \{1, 2, \dots, N\} \cup \{N+1, \dots, N+U\}$

unknown classes

Deep embedding space for category $k \rightarrow S_k$

open space \mathcal{O}_k [not k]

$$\mathcal{O}_k = \mathcal{O}_k^{\text{pos}} \cup \mathcal{O}_k^{\text{neg}}$$

positive open
space of
other known
classes

infinite unknown
space

Here, $\mathcal{D}_L^k \in S^k$

$$\exists \mathcal{D}_k^{\text{pos}} \in \mathcal{D}_k^{\text{pos}}$$

$$\mathcal{D}_k \in \mathcal{D}_k^{\text{neg}} \Rightarrow \mathbb{R}^d$$

Binary prediction function ψ_k we have one of this for each class N

$$\psi_k \text{ for } k \text{ class: } 1 \text{ ; } x_i \in y_k$$

$$0 \text{ ; } x_i \notin y_k$$

objective: Learn ψ_k for all known N classes.

Expected error

$$\arg \min_{\psi_k} \left\{ \mathcal{R}_k \mid \mathcal{R}_k = \mathcal{R}_k(\psi_k, \mathcal{D}_k \cup \mathcal{D}_k^{\text{pos}}) + \alpha \mathcal{R}_0(\psi_k, \mathcal{D}_k^{\text{neg}}) \right\}$$

regularizer

open risk space

Empirical classification risk for known

further

$$\mathcal{R}_0(\psi_k, \mathcal{D}_k^{\text{neg}}) = \frac{\int_{\mathcal{D}_k^{\text{neg}}} \psi_k(x) dx}{1}$$

should be 0

$$\int_{S_k \cup \partial_k} \frac{\psi_k(x) dx}{1} \text{ only for } S_k \text{ region}$$

S_0 for multiclass prediction

$$f = \Theta(\psi_1, \psi_2, \dots, \psi_N)$$

↳ integration

final objective:

for all ψ_k

$$\arg \min_{f \in \mathcal{H}} \left\{ \mathcal{R}_L(f, \mathcal{D}_L) + \alpha \cdot \sum_{k=1}^N \mathcal{R}_0(f, \mathcal{D}_0) \right\}$$

$f: \mathbb{R}^d \mapsto \mathcal{N}$; measurable function.

Reciprocal point for classification

$\mathcal{P}^k \rightarrow$ latent representation of $\mathcal{D}_L^{\neq k} \cup \mathcal{D}_u$

↳ seem to sample from \mathcal{D}_k than to \mathcal{S}_k

$$\max (\mathcal{S}(\mathcal{D}_L^{\neq k} \cup \mathcal{D}_u, \mathcal{P}^k)) \leq d; \forall d \in \mathcal{S}(\mathcal{D}_L^k, \mathcal{P}^k)$$

↓

set of distance of all sample between
the reciprocal point & corresponding
known class

Distance between x & the reciprocal point P^k

$$d(C(x), P^k) = d_e(C(x), P^k) - d_d(C(x), P^k)$$

\swarrow Euclidean in \mathbb{R}^n (maximize) \swarrow Dot product (- minimize)

[opposite of their R_p]

Softmax function for class assignment:

$$P(y=k | x, C, P) = \frac{e^{\overbrace{d(C(x), P^k)}^{\text{max distance}}}}{\sum_{i=1}^N e^{d(C(x), P^i)}}$$

Learn Θ to minimize:

$$L_c(x; \theta, P) = -\log p(y=k | \underline{x}, \underline{C}, \underline{P})$$

↪ corresponding to $R_c(f, D_L)$

We also achieve: $\arg \max_{f \in H} \{ \mathcal{J}(D_L^k, P^k) \}$

Adversarial margin Constraint:

AMC \rightarrow minimize $R_o(f, D_u)$

Classwise openspace \rightarrow Global open space

$$D_G = \bigcap_{k=1}^N (D_k^{\text{pos}} \cup D_k^{\text{neg}}) = \bigcap_{k=1}^N D^k$$

further,

$$\max \left(\mathcal{J}(D_L^{\neq k} \cup D_u, P^k) \right) \leq R$$

↙ limiting the open space by R

by

$$L_o(x; \theta, P^k, R^k) = \max \left(\underbrace{d_c(C(x), P^k)}_{\downarrow} - R, 0 \right)$$

constraint on distance.

$$x \in S^k$$

∴ IDEA: Explained in th 1
Adverse to each other.

ℓ_0 , ℓ_c are minimize simultaneously if
there exists $\max(\mathcal{S}(\mathcal{D}_L^k, \mathcal{P}^k)) \leq R$

Bound for class k

$$\arg \min_{f \in \mathcal{H}} \left\{ \max \left(\left\{ \mathcal{S}(\mathcal{D}_L^k, \mathcal{P}^k) - R \right\} \cup \{0\} \right) \right\}$$

Learning open set network:

$$\mathcal{L}(x, y; \theta, \mathcal{P}, \mathcal{R}) = \mathcal{L}_c(x; \theta, \mathcal{P}) + \lambda \mathcal{L}_0(x; \theta, \mathcal{R}, \mathcal{P})$$

Adversarial Goal:

$$\max_G \underbrace{\frac{1}{n} \sum_{i=1}^n \left[-\frac{1}{N} \sum_{k=1}^N \overbrace{\mathcal{S}(z_i, \mathcal{P}^k)}^{\text{Softmax}(\mathcal{D}_k(C(G-z_i)), \mathcal{P}^k)} \cdot \log(\mathcal{S}(z_i, \mathcal{P}^k)) \right]}_{\text{Empirical Entropy!!}}$$

Combined with adversarial setting

GAN
Adversary

$$\max_D \log D(x_i) + \log D(1 - G(z_i))$$

$$\text{vs} \max_G \log D(G(z_i))$$

final goal:

$$\max_G \left[\underbrace{\log D(G(z_i))}_{\text{GAN}} + \beta \underbrace{H(z_i, \mathcal{P})}_{\text{Extra points.}} \right]$$

classifier C is optimized

$$\min_C \frac{1}{n} \sum_{i=1}^n \left[\underbrace{\mathcal{L}(x_i, y_i)}_{\text{ARPL loss}} - \beta H(z_i, \mathcal{P}) \right]$$

Reducing open space risks

discriminative.

