# Adversarial Support Alignment

Shangyuan Tong[1]    Timur Garipov[1]    Yang Zhang[2]    Shiyu Chang[3]    Tommi Jaakkola[1]

[1]MIT CSAIL   [2]MIT-IBM Watson AI Lab   [3]UC Santa Barbara

## 1 Background: distribution alignment

**Given**   $\mathcal{P} = \{p^\theta \mid \theta \in \Theta\}$   $\mathcal{Q} = \{q^\theta \mid \theta \in \Theta\}$   $p^\theta(x)$   $q^\theta(x)$

**Find**   $\theta^* : p^{\theta^*} = q^{\theta^*}$

$p^{\theta^*} = q^{\theta^*}$

- **Generative Models** (GAN, Goodfellow et al. 2014)
  - alignment of generated and data distributions.
- **Domain Adaptation** (DANN, Ganin et al. 2016)
  - alignment of representations across domains.

### Domain adaptation

**Goal:** learn classifier on target domain using
- labeled dataset in source domain.
- unlabeled dataset in target domain.

Source $p(x)$   Labeled examples

car

plane    horse

Feature extractor $F^\theta : \mathcal{X} \to \mathcal{Z}$

Target $q(x)$   Unlabeled examples

**Distribution alignment approach:**
Feature extractor network $F^\theta : \mathcal{X} \to \mathcal{Z}$

Training:
- separate classes in source.
- align distributions of source/target features.

## 2 Motivation

- Distribution alignment is not always desired.
- Example: domain adaptation under **label distribution shift.**
- We propose to align only the **supports** of distributions.

No adaptation       Distribution alignment       **Support alignment** (ours)

SourceTarget

Target accuracy: 63%    Target accuracy: 75%    Target accuracy: 94%

Source class distribution       Target class distribution
[33%, 33%, 33%]                 [23%, 65%, 12%]

## 3 Support divergence

**Symmetric Support Difference (SSD)** divergence   $q(x)$   $p(x)$

$$\mathcal{D}_\triangle(p,q) = \mathop{\mathbb{E}}_{x^q \sim q}\Big[d(x^q, \mathrm{supp}(p))\Big] + \mathop{\mathbb{E}}_{x^p \sim p}\Big[d(x^p, \mathrm{supp}(q))\Big]$$

$$d(x^q, \mathrm{supp}(p)) = \inf_{x^p \in \mathrm{supp}(p)} d(x^q, x^p) \qquad d(x^p, \mathrm{supp}(q)) = \inf_{x^q \in \mathrm{supp}(q)} d(x^p, x^q)$$

1) $\mathcal{D}_\triangle(p,q) \geq 0 \ \forall p, q;$    2) $\mathcal{D}_\triangle(p,q) = 0 \iff \mathrm{supp}(p) = \mathrm{supp}(q)$

## 4 Support alignment via log-loss discriminator

$$\sup_{f:\mathcal{X} \to [0,1]} \mathop{\mathbb{E}}_{x \sim p}\big[\log f(x)\big] + \mathop{\mathbb{E}}_{y \sim q}\big[\log(1 - f(y))\big]$$

Optimal solution:   $f^*(x) = \dfrac{p(x)}{p(x) + q(x)}$

$q(x)$   $p(x)$

**Theorem**

The mapping $f^* : \mathcal{X} \to [0,1]$
realized by the optimal discriminator
preserves support discrepancy

$$\mathcal{D}_\triangle(p,q) = 0 \iff \mathcal{D}_\triangle(f^*_\sharp p, f^*_\sharp q) = 0$$

$f^*_\sharp p, \ f^*_\sharp q$ — pushforward distributions

$x \in \mathrm{supp}(p) \cap \mathrm{supp}(q)$

$x \in \mathrm{supp}(q) \setminus \mathrm{supp}(p)$    $x \in \mathrm{supp}(p) \setminus \mathrm{supp}(q)$

$f^*(x) = 0$    $f^*(x) \in (0,1)$    $f^*(x) = 1$

**Remark:** the result holds for $g : \mathcal{X} \to \mathbb{R}$
$g(x) : f(x) = \mathrm{sigmoid}(g(x))$

0   $f^*_\sharp q$    $f^*_\sharp p$   1

## 5 Method: Adversarial Support Alignment (ASA)

Discriminator's objective:   $\min\limits_{g:\mathcal{X} \to \mathbb{R}} \mathcal{L}_D(\theta, g)$

$$\mathcal{L}_D(\theta, g) = \mathop{\mathbb{E}}_{x \sim p^\theta}\big[\log\big(1 + e^{-g(x)}\big)\big] + \mathop{\mathbb{E}}_{x \sim q^\theta}\big[\log\big(1 + e^{g(x)}\big)\big]$$

Alignment objective:   $\min\limits_{\theta} \mathcal{L}_A(\theta, g)$

$$\mathcal{L}_A(\theta, g) = \mathop{\mathbb{E}}_{x \sim p^\theta}\big[d(g(x), \mathrm{supp}(g_\sharp q))\big] + \mathop{\mathbb{E}}_{x \sim q^\theta}\big[d(g(x), \mathrm{supp}(g_\sharp p))\big]$$

$[g_\sharp q](t)$       $[g_\sharp p](t)$

$L_q(t) = d(t, \mathrm{supp}(g_\sharp p))$     $L_p(t) = d(t, \mathrm{supp}(g_\sharp q))$
$= |t - \pi_p^*(t)|$                           $= |t - \pi_q^*(t)|$
$\pi_p^*(t) = \operatorname*{argmin}\limits_{t^p \in \mathrm{supp}(g_\sharp p)} |t - t^p|$     $\pi_q^*(t) = \operatorname*{argmin}\limits_{t^q \in \mathrm{supp}(g_\sharp q)} |t - t^q|$

$t_1^p$ $t_2^p$ $t_3^p$ $t_4^p$

$t_1^q$ $t_2^q$  $t_3^q$  $t_4^q$

## 6 Results: domain adaptation under label distribution shift

Average and minimum class accuracy (%) on USPS→MNIST, LeNet
across different levels of shifts in label distributions ($\alpha$).

| Algorithm | $\alpha = 0.0$ no shift | | $\alpha = 1.0$ | | $\alpha = 1.5$ | | $\alpha = 2.0$ severe shift | |
|---|---|---|---|---|---|---|---|---|
|  | average | min | average | min | average | min | average | min |
| No DA | 71.9 | 20.3 | 72.9 | 25.8 | 71.3 | 27.5 | 71.3 | 16.6 |
| DANN | 97.8 | 96.0 | 83.5 | 25.1 | 70.0 | 01.1 | 57.8 | 00.9 |
| VADA | **98.0** | 96.2 | 88.2 | 48.9 | 78.2 | 06.6 | 61.9 | 01.4 |
| IWDAN | 97.5 | 95.7 | 95.7 | 81.3 | 86.5 | 15.2 | 74.4 | 07.3 |
| IWCDAN | **98.0** | **96.6** | **96.7** | 85.1 | 91.3 | 66.5 | 77.5 | 22.2 |
| sDANN-4 | 87.4 | 05.6 | 94.9 | **85.7** | 86.8 | 21.6 | 81.5 | 39.3 |
| ASA-sq (ours) | 93.7 | 89.2 | 92.3 | 83.5 | 90.9 | 69.9 | 87.2 | 62.5 |
| ASA-abs (ours) | 94.1 | 88.9 | 92.8 | 78.9 | **92.5** | **82.4** | **90.4** | **68.4** |

STL→CIFAR, DeepCNN

| Algorithm | $\alpha = 0.0$ no shift | | $\alpha = 2.0$ severe shift | |
|---|---|---|---|---|
|  | average | min | average | min |
| No DA | 69.9 | 49.8 | 65.8 | 43.7 |
| DANN | 75.3 | 54.6 | 63.3 | 27.0 |
| VADA | **76.7** | **56.9** | 63.2 | 25.5 |
| IWDAN | 69.9 | 50.5 | 64.4 | 36.8 |
| IWCDAN | 70.1 | 47.8 | 64.5 | 37.0 |
| sDANN-4 | 71.8 | 52.1 | 66.4 | 39.0 |
| ASA-sq (ours) | 71.7 | 52.9 | **68.1** | **44.7** |
| ASA-abs (ours) | 71.6 | 49.0 | 67.8 | 40.9 |

VisDA-17, ResNet-50

| Algorithm | $\alpha = 0.0$ no shift | | $\alpha = 2.0$ severe shift | |
|---|---|---|---|---|
|  | average | min | average | min |
| No DA | 49.5 | 22.2 | 45.3 | 19.5 |
| DANN | **75.4** | 36.7 | 43.1 | 03.6 |
| VADA | 75.3 | 40.5 | 43.9 | 08.5 |
| IWDAN | 73.2 | 31.7 | 45.1 | 04.6 |
| IWCDAN | 71.6 | 27.6 | 38.3 | 00.6 |
| sDANN-4 | 72.4 | 37.8 | 50.7 | 18.6 |
| ASA-sq (ours) | 64.9 | 35.7 | 51.9 | 18.3 |
| ASA-abs (ours) | 64.8 | **40.6** | **52.5** | **19.7** |

## 7 Takeaways

- **Support alignment** as extreme relaxation of **distribution alignment**.
- **Adversarial Support Alignment (ASA):** method to align distribution supports.
- Evaluation on domain adaptation datasets under **label distribution shift**.
- Bonus: spectrum of **relaxed alignment** approaches based on **optimal transport**.

**Wasserstein distance**
$\mathcal{D}_W(p,q)$
$\|$
$\inf\limits_{\gamma} \mathop{\mathbb{E}}_{(x,y) \sim \gamma}[d(x,y)],$
$s.t. \int \gamma(x,y)dy = p(x)$
$\int \gamma(x,y)dx = q(y)$

$\mathcal{D}_W(p,q) = 0$
$\updownarrow$
$p = q$

**$\beta$-Wasserstein distance** (Wu et al. 2019)
$\mathcal{D}_W^{\beta,\beta}(p,q) = \mathcal{D}_W^\beta(p,q) + \mathcal{D}_W^\beta(q,p)$
$\|$
$\inf\limits_{\gamma} \mathop{\mathbb{E}}_{(x,y) \sim \gamma}[d(x,y)],$
$s.t. \int \gamma(x,y)dy = p(x)$
$\int \gamma(x,y)dx \leq (1+\beta)q(y)$

$\mathcal{D}_W^{\beta,\beta}(p,q) = 0$
$\updownarrow$
$\frac{1}{1+\beta} \leq \frac{p}{q} \leq 1 + \beta$

**SSD divergence**
$\mathcal{D}_\triangle(p,q) = \mathcal{D}_{\widetilde{W}}^\infty(p,q) + \mathcal{D}_{\widetilde{W}}^\infty(q,p)$
$\|$
$\mathop{\mathbb{E}}_{x \sim p(x)}\big[d(x, \mathrm{supp}(q))\big]$

$\mathcal{D}_\triangle(p,q) = 0$
$\updownarrow$
$\mathrm{supp}(p) = \mathrm{supp}(q)$