

### Generate To Adapt: Aligning Domains using Generative Adversarial Networks

Swami Sankaranarayanan\* Yogesh Balaji\* Carlos D. Castillo Rama Chellappa UMIACS, University of Maryland, College Park CVPR 2018

### **Content:**

1. Introduction

2. Approach

3. Experiments and Results

4. Conclusion

## Introduciton:

Due to many factors (e.g., illumination, pose, and image quality), there is always a **distribution change or domain shift** between two domains that can degrade the performance.



### Introduction:

- 1. propose an approach that leverages unsupervised data to bring the source and target distributions closer in a learned joint feature space.
- 2. We accomplish this by inducing a symbiotic relationship between the learned embedding and a generative adversarial network.
- 3. by far the only GAN-based method that has been shown to work well across different datasets such as OFFICE and DIGITS.



 $x_g = [F(x), z, l],$ 

# Approach:

1. Given source images as input, D outputs two distributions D<sub>data</sub> and D<sub>cls</sub>

$$L_{data,src} + L_{cls,src} = \mathbf{E}_{x \sim S} \max_{D} \log(D_{data}(x)) + \log(1 - D_{data}(G(x_g))) + \log(D_{cls}(x)_y)$$

D data (x): the probability of the input being real; D cls (x): the class probability distribution of the input x; D cls (x)y : the probability assigned by the classifier mapping D cls for input x to class Y

2. G is updated using a combination of adversarial loss and classification loss .

$$L_G = \min_{G} \mathbf{E}_{x \sim S} - \log(D_{cls}(G(x_g))_y) + \log(1 - D_{data}(G(x_g))),$$

#### 3. F and C are updated .

$$L_C = \min_{C} \min_{F} \mathbf{E}_{x \sim S} - \log(C(F(x))_y),$$
$$L_{cls,src} = \min_{F} \mathbf{E}_{x \sim S} - \alpha \log(D_{cls}(G(x_g))_y))$$

4. The target embeddings output by F along with the random noise vector z and the fake label encoding I are input to G. The generated target images  $G(x_g)$  are then given as input to D.

$$L_{adv,tgt} = \max_{D} \mathbf{E}_{x \sim \mathcal{T}} \log(1 - D_{data}(G(x_g)))$$

Algorithm 1 Iterative training procedure of our approach

- 1: training iterations = N
- 2: for t in 1:N do
- 3: Sample k images with labels from source domain S:  $\{s_i, y_i\}_{i=1}^k$
- 4: Let  $f_i = F(s_i)$  be the embeddings computed for the source images.
- 5: Sample k images from target domain  $\mathcal{T}: \{t_i\}_{i=1}^k$
- 6: Let  $h_i = F(t_i)$  be the embeddings computed for the target images.
- 7: Sample k random noise samples  $\{z_i\}_{i=1}^k \sim \mathcal{N}(0, 1)$ .
- 8: Let  $f_{g_i}$  and  $h_{g_i}$  be the concatenated inputs to the generator.
- 9: Update discriminator using the following objectives:

$$L_D = L_{data,src} + L_{cls,src} + L_{adv,tgt} \tag{3}$$

• 
$$L_{data,src} = \max_{D} \frac{1}{k} \sum_{i=1}^{k} \log(D_{data}(s_i)) + \log(1 - D_{data}(G(f_{g_i}))))$$
  
•  $L_{cls,src} = \max_{D} \frac{1}{k} \sum_{i=1}^{k} \log(D_{cls}(s_i)_{y_i})$ 

• 
$$L_{adv,tgt} = \max_D \frac{1}{k} \sum_{i=1}^k \log(1 - D_{data}(G(h_{g_i})))$$

10: Update the generator, only for source data, through the discriminator gradients computed using real labels.

$$L_G = \min_G \frac{1}{k} \sum_{i=1}^k -\log(D_{cls}(G(f_{g_i}))_{y_i}) + \log(1 - D_{data}(G(f_{g_i}))))$$
(4)

11: Update the embedding F using a linear combination of the adversarial loss and classification loss. Update the classifier C for the source data using a cross entropy loss function.

$$L_F = L_C + \alpha \, L_{cls,src} + \beta \, L_{F_{adv}} \tag{5}$$

• 
$$L_C = \min_C \min_F \frac{1}{k} \sum_{i=1}^k -\log(C(f_i)_{y_i})$$
  
•  $L_{cls,src} = \min_F \frac{1}{k} \sum_{i=1}^k -\log(D_{cls}(G(f_{g_i}))_{y_i})$   
•  $L_{F_{adv}} = \min_F \frac{1}{k} \sum_{i=1}^k \log(1 - D_{data}(G(h_{g_i})))$ 

12: end for

### Experiments and Results:

#### 1. Digit Experiments

Table 1: Accuracy (mean  $\pm$  std%) values for cross-domain recognition tasks over five independent runs on the digits based datasets. The best numbers are indicated in **bold** and the second best are <u>underlined</u>. – denotes unreported results. MN: MNIST, US: USPS, SV: SVHN. MN $\rightarrow$ US (p) denotes the MN $\rightarrow$ US experiment run using the protocol established in [17], while MN $\rightarrow$ US (f) denotes the experiment run using the entire datasets. (Refer to Digits experiments section for more details)

Method	$MN \rightarrow US \ (p)$	$MN \rightarrow US \ (f)$	$\text{US} \rightarrow \text{MN}$	$SV \to MN$
Source only	$75.2 \pm 1.6$	$79.1\pm0.9$	$57.1 \pm 1.7$	$60.3\pm1.5$
RevGrad [4]	$77.1 \pm 1.8$	-	$73.0 \pm 2.0$	73.9
DRCN [5]	$91.8 \pm 0.09$	-	$73.7\pm0.04$	$82.0 \pm 0.16$
CoGAN [15]	$91.2\pm0.8$	-	$89.1\pm0.8$	-
ADDA [32]	$89.4 \pm 0.2$	-	$\underline{90.1} \pm 0.8$	$76.0 \pm 1.8$
PixelDA [1]	-	95.9	-	-
Ours	$\textbf{92.8} \pm \textbf{0.9}$	$95.3 \pm 0.7$	<b>90.8</b> ± 1.3	<b>92.4</b> ± 0.9



Figure 2: TSNE visualization of SVHN  $\rightarrow$  MNIST adaptation. In (a), the source data shown in *red* is classified well into distinct clusters but the target data is clustered poorly. On applying the proposed approach, as shown in (b), both the source and target distributions are brought closer in a class consistent manner.

### 2. OFFICE experiments

Table 2: Accuracy (mean  $\pm$  std%) values on the OFFICE dataset for the standard protocol for unsupervised domain adaptation [6]. Results are reported as an average over 5 independent runs. The best numbers are indicated in **bold** and the second best are <u>underlined</u>. – denotes unreported results. A: Amazon, W: Webcam, D: DSLR

Method	$\mathbf{A} \to \mathbf{W}$	$\mathrm{D} \to \mathrm{W}$	$W \to D$	$\mathbf{A} \to \mathbf{D}$	$\mathbf{D} \to \mathbf{A}$	$W \to A$	Average
ResNet - Source only [9]	$68.4 \pm 0.2$	$96.7\pm0.1$	$99.3\pm0.1$	$68.9\pm0.2$	$62.5\pm0.3$	$60.7\pm0.3$	76.1
TCA [23]	$72.7\pm0.0$	$96.7\pm0.0$	$\underline{99.6}\pm0.0$	$74.1\pm0.0$	$61.7\pm0.0$	$60.9\pm0.0$	77.6
GFK [6]	$72.8\pm0.0$	$95.0\pm0.0$	$98.2\pm0.0$	$74.5\pm0.0$	$63.4\pm0.0$	$61.0\pm0.0$	77.5
DDC [33]	$75.6\pm0.2$	$76.0\pm0.2$	$98.2\pm0.1$	$76.5\pm0.3$	$62.2\pm0.4$	$61.5\pm0.5$	78.3
DAN [16]	$80.5\pm0.4$	$97.1\pm0.2$	$\underline{99.6}\pm0.1$	$78.6\pm0.2$	$63.6\pm0.3$	$62.8\pm0.2$	80.4
RTN [18]	$84.5\pm0.2$	$96.8\pm0.1$	$99.4\pm0.1$	$77.5\pm0.3$	$66.2\pm0.2$	$64.8\pm0.3$	81.6
RevGrad [4]	$82.0\pm0.4$	$96.9\pm0.2$	$99.1\pm0.1$	$79.4\pm0.4$	$68.2\pm0.4$	$67.4\pm0.5$	82.2
JAN [19]	$\underline{85.4}\pm0.3$	$\underline{97.4}\pm0.2$	$\textbf{99.8}\pm0.2$	$\underline{84.7}\pm0.3$	$\underline{68.6}\pm0.3$	$\underline{70.0} \pm 0.4$	<u>84.3</u>
Ours	$\textbf{89.5}\pm0.5$	$\textbf{97.9}\pm0.3$	$\textbf{99.8}\pm0.4$	$\textbf{87.7}\pm0.5$	$\textbf{72.8}\pm0.3$	$\textbf{71.4}\pm0.4$	86.5

#### **3. Synthetic to Real experiments**

In this experiment, we use CAD synthetic dataset and a subset of PASCAL VOC dataset as our source and target sets respectively.

Table 3: Accuracy (mean  $\pm$  std%) values over five independent runs on the Synthetic to real setting. The best numbers are indicated in **bold**.

Method	$\mathrm{CAD} \to \mathrm{PASCAL}$		
VGGNet - Source only	$38.1 \pm 0.4$		
RevGrad [4]	$48.3\pm0.7$		
RTN [18]	$43.2\pm0.5$		
JAN [19]	$46.4\pm0.8$		
Ours	$50.4 \pm 0.6$		

#### **Ablation Study**

#### Table 5: Ablation study for OFFICE $A \rightarrow W$ setting

Setting	Accuracy(in %)		
Stream 1 - Source only	68.4		
Stream 1 + Stream 2 ( $C_1$ only)	80.5		
Stream 1 + Stream 2 $(C_1 + C_2)$	89.5		

# Conclusion:

- We proposed a joint adversarial discriminative approach that transfers the information of the target distribution to the learned embedding using a generator-discriminator pair.
- address the problem using experiments on three different tasks