

Learning without limits

(Inspired by "Lifelong Domain Adaptation via Consolidated Internal Distribution" by Mohammad Rostami)

Group 69 : DCCLXXI

Indian Institute of Technology, Kanpur

November 26, 2024

Introduction

- Feature engineering using deep neural networks struggles with generalization to unseen data, especially when domain shifts occur between training and testing datasets.
- **Continual learning** (CL) aims to create models that can adapt to changing data distributions; learning new information without forgetting previous knowledge.
- **Unsupervised Domain Adaptation** (UDA) trains models for a target domain with unlabeled data by transferring knowledge from a related source domain with labeled data.
- **Catastrophic forgetting** refers to the loss of knowledge from previous tasks when learning a new task.

Unsupervised Domain Adaptation (UDA) :

- Existing domain alignment methods use probability distance to measure distribution discrepancies and train an encoder to minimize cross-domain distance.
- Wasserstein Distance (WD) minimizes cross-domain distance due to its non-vanishing gradients, essential for gradient-based optimization. The Sliced Wasserstein Distance (SWD) improves computational efficiency with a closed-form solution based on empirical samples.

Continual Learning (CL) :

- Model Regularization** consolidates crucial network weights to preserve past knowledge during updates.
- Experience Replay** stores key data points in a memory buffer, replaying them alongside new task data.

What to solve?

- Standard classification problem on source domain S with a labeled training dataset $D_S = (X_0, Y_0)$, where X_0 is the feature matrix, and Y_0 is the corresponding label matrix involves training a deep neural network $f_\theta : \mathcal{X} \rightarrow \mathcal{Y}$ by minimizing the loss function $L(\cdot)$, with optimal parameters: $\hat{\theta}_0 = \arg \min_{\theta} \sum_{i=1}^N L(f_\theta(x_i^0), y_i^0)$
- In this problem, we are presented with a sequence of target domains \mathcal{T}_t , where $t = 1, \dots, T$, each having an unlabeled dataset $D_{\mathcal{T}}^t = (X_t)$. Each domain \mathcal{T}_t follows its own unique input distribution $p_t(x)$, and these distributions differ: $\forall t_1, t_2 : p_{t_1} \neq p_{t_2}$.

Also need to tackle:

- Failure to apply traditional optimization techniques
- Non-stationary Distributions
- Catastrophic Forgetting
- Sequential Access Constraint

Solution Presented

The presented solution addresses domain shift and catastrophic forgetting by splitting the model f_θ into two components:

- **Deep Encoder** $\phi_v(\cdot) : X \rightarrow Z$: Maps input data to a latent space Z where source domain classes are well-separated.
- **Classifier** $h_w(\cdot) : Z \rightarrow Y$: Maps the latent embeddings Z to output predictions.

By stabilizing the distribution in Z , the solution minimizes the distance between source embeddings $\phi(p_0(x^0))$ and target embeddings $\phi(p_t(x^t))$. The encoder is trained to map the data into a multi-modal distribution $p_0^J(z)$ in the embedding space, modeled by a Gaussian Mixture Model (GMM) represented as: $p_0^J(z) = \sum_{j=1}^k \alpha_0^j N(z | \mu_0^j, \Sigma_0^j)$. Hence The MAP estimates are :

$$\hat{\alpha}_0^j = \frac{|S_0^j|}{N}, \quad \hat{\mu}_j = \frac{1}{|S_0^j|} \sum_{(x_0^j, y_0^j) \in S_0^j} \phi_v(x_0^j), \quad \hat{\Sigma}_0^j = \frac{1}{|S_0^j|} \sum_{(x_0^j, y_0^j) \in S_0^j} (\phi_v(x_0^j) - \hat{\mu}_j)^T (\phi_v(x_0^j) - \hat{\mu}_j) \quad (1)$$

Solution Presented

The continual learning process aligns the target domain's distribution with the learned distribution in the embedding space by combining Pseudo-Generated Data and Filtered Input Data to align target distributions with the embedding space.

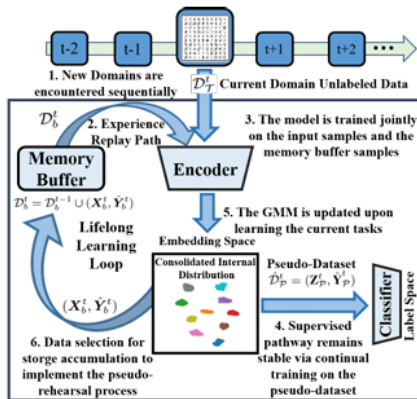
The optimization problem at time t is:

$$\min_{v,w} \sum_{i=1}^N L(h_w(z_i^p), \hat{y}_{p,t}^i) + \lambda D(\phi_v(p_t(X_t)), \hat{p}_t^J(Z_t^P))$$

To prevent catastrophic forgetting, experience replay is introduced, storing and replaying a subset of past data. The updated objective function incorporates buffer samples:

$$\min_{v,w} \left(\sum_{i=1}^N L(h_w(z_i^p), \hat{y}_{p,t}^i) + \sum_{i=1}^{N_b} L(h_w(\phi_v(x_i^b)), \hat{y}_i^b) \right. \\ \left. + \lambda D(\phi_v(p_t(X_t)), \hat{p}_t^J(Z_t^P)) + \lambda D(\phi_v(p_t(X_t^b)), \hat{p}_t^J(Z_t^P)) \right) \quad (2)$$

Model & Algorithm



Algorithm 1 LDAuCID (λ, τ, N_b)

- 1: **Source Training:**
- 2: **Input:** source labeled dataset $\mathcal{D}_S = (\mathbf{X}_0, \mathbf{Y}_0)$
- 3: $\hat{\theta}_0 = (\hat{\mathbf{w}}_0, \hat{\mathbf{v}}_0) = \arg \min_{\theta} \sum_i \mathcal{L}(f_{\theta}(\mathbf{x}_i^0), \mathbf{y}_i^0)$
- 4: **Internal Distribution Estimation:**
- 5: Use Eq. (2) and estimate α_j^0, μ_j^0 , and Σ_j^0
- 6: **Memory Buffer Initialization**
- 7: $\mathcal{D}_b^0 = (\mathbf{X}_b^0, \hat{\mathbf{Y}}_b^0)$
Pick the N_b/k samples with the least $d_{j,l}^t = \|\mu_j^t - \phi(\mathbf{x}_l^t)\|_2^2$, $\hat{\mathbf{y}}_b^{0,i} = \arg \max f_{\hat{\theta}_t}(\mathbf{x}_b^{0,N_b})$
- 8: **Continual Unsupervised Domain Adaptation:**
- 9: **for** $t = 1, \dots, T$ **do**
- 10: **Input:** target unlabeled dataset $\mathcal{D}_T^t = (\mathbf{X}_t)$
- 11: **Pseudo-Dataset Generation:**
- 12: $\hat{\mathcal{D}}_P^t = (\mathbf{Z}_P^t, \hat{\mathbf{Y}}_P^t) =$
- 13: $(\{\mathbf{z}_1^{p,t}, \dots, \mathbf{z}_N^{p,t}\}, \{\hat{\mathbf{y}}_1^{p,t}, \dots, \hat{\mathbf{y}}_N^{p,t}\})$, where:
 $\mathbf{z}_i^{p,t} \sim \hat{p}_J^{t-1}(\mathbf{z})$, $1 \leq i \leq N_p$ and
 $\hat{\mathbf{y}}_i^{p,t} = \arg \max_j \{h_{\hat{\mathbf{w}}_t}(\mathbf{z}_i^{p,t})\}$ if with confidence τ : $\max_j \{h_{\hat{\mathbf{w}}_t}(\mathbf{z}_i^{p,t})\} > \tau$
- 14: **for** $itr = 1, \dots, ITR$ **do**
- 15: draw data batches from \mathcal{D}_T^t and $\hat{\mathcal{D}}_P^t$
- 16: Update the model by solving Eq. (4)
- 17: **end for**
- 18: **Internal Distribution Estimate Update:**
- 19: Use Eq. (2) similar to step 5 above.
- 20: **Memory Buffer Update**
- 21: $\mathcal{D}_b^t = \mathcal{D}_b^{t-1} \cup (\mathbf{X}_b^t, \hat{\mathbf{Y}}_b^t)$, where $(\mathbf{X}_b^t, \hat{\mathbf{Y}}_b^t)$ is computed similar to step 7 above.
- 22: **end for**

Algorithm: LDAuCID

1 Source Training:

Input the source labeled dataset $\mathcal{DS} = (X^0, Y^0)$. Initialize the model parameters: $\hat{\theta}^0 = (\hat{w}^0, \hat{v}^0) = \arg \min_{\theta} \sum_i L(f_{\theta}(x_i^0), y_i^0)$.

2 Internal Distribution Estimation:

Use Eq. (1) to estimate α_j^0 , μ_j^0 , and Σ_j^0 .

3 Memory Buffer Initialization:

Initialize the memory buffer $\mathcal{D}_b^0 = (X_b^0, Y_b^0)$ by selecting N_b/k samples with the least distances:

$$d_{j,l}^t = |\mu_j^t - \phi(x_l^t)|_2^2$$

where $\hat{y}_{i,b}^0 = \arg \max_{\theta_t} f_{\theta_t}(x_{b,N_b}^0)$.

Algorithm: LDAuCID (contd.)

4 Continual Unsupervised Domain Adaptation

For each time step $t = 1, \dots, T$:

- Input the target unlabeled dataset $\mathcal{D}_T^t = (X^t)$.

- **Pseudo-Dataset Generation:**

Generate the pseudo-dataset $\hat{\mathcal{D}}_P^t = (Z_P^t, \hat{Y}_P^t)$ as follows:

$$Z_P^t = [z_1^{p,t}, \dots, z_N^{p,t}], \quad Y_P^t = [\hat{y}_1^{p,t}, \dots, \hat{y}_N^{p,t}]$$

where $z_i^{p,t} \sim \hat{p}_J^{t-1}(z)$, $1 \leq i \leq N_p$

and $\hat{y}^{p,t}_i = \arg \max_j h_{\hat{w}_t}(z^{p,t}_i)$ if the confidence threshold τ is satisfied:
 $\max_j h_{\hat{w}_t}(z^{p,t}_i) > \tau$.

Algorithm: LDAuCID (contd.)

④ Continual Unsupervised Domain Adaptation

For each time step $t = 1, \dots, T$:

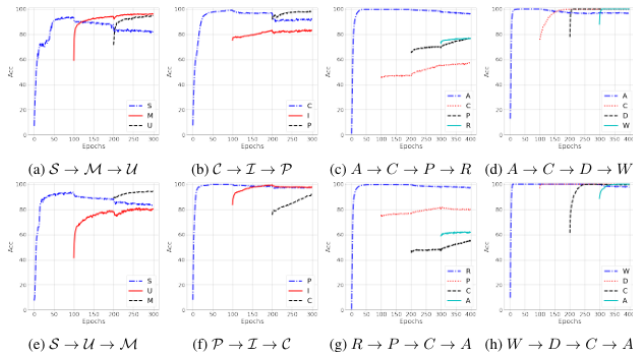
- **For** each iteration $\text{itr} = 1, \dots, \text{ITR}$:
 - Draw data batches from \mathcal{D}_T^t and $\hat{\mathcal{D}}_P^t$.
 - Update the model by solving Eq. (2).
- **Internal Distribution Estimate Update:** Update internal distribution parameters using Eq. (1) as in Step 2.
- **Memory Buffer Update:**
 - Update the memory buffer: $\mathcal{D}_b^t = \mathcal{D}_b^{t-1} \cup (X_b^t, \hat{Y}_b^t)$ where X_b^t and \hat{Y}_b^t are selected as in Step 2.

Results & Inference

Datasets: ImageCLEF-DA, Office-Home, Office-Caltech, Digit Recognition. **Methodology:** VGG16 (digit tasks), ResNet-50 (ImageCLEF-DA, Office-Home), Decaf6 (Office-Caltech).

Key Findings:

- **Learning Curves:** Performance drops initially, improves as model adapts.



Source: "Lifelong Domain Adaptation via Consolidated Internal Distribution" (Rostami, 2021)

Key Findings (cont.):

- **Catastrophic Forgetting:** Mitigated by LDAuCID, though more pronounced in challenging tasks (e.g., SVHN).
- **Dataset Performance:**
 - ImageCLEF-DA: Significant improvement due to balanced data.
 - Office-Home: Strong performance; CDAN outperforms.
- **Buffer/Hyperparameters:** Larger memory buffer (Nb) improves performance; $\tau \approx 1$ reduces label pollution.
- **Imbalanced Data:** LDAuCID performs well, though with reduced performance on imbalanced data.

- **LDAuCID** integrates **UDA** and **CL** to address domain shift and catastrophic forgetting by aligning internal data distributions in the embedding space.
- **Outperforms traditional UDA methods**, particularly in moderate domain shifts, and mitigates catastrophic forgetting effectively.
- **Performance** can be improved with techniques like **class-conditional alignment** for large domain gaps (e.g., Office-Home dataset).
- Uses **experience replay** and a **memory buffer** to retain acquired knowledge, enhancing continual learning, especially with imbalanced data or domain shifts.