

Modeling Sequential Domain Shift through Estimation of Optimal Sub-spaces for Categorization

Suranjana Samanta
 ssamanta@cse.iitm.ac.in
 Tirumarai A Selvan
 tirumarai.selvan@gmail.com
 Sukhendu Das
 http://www.cse.iitm.ac.in/~sdas/

Visualization and Perception Lab
 Dept. of CS&E
 Indian Institute of Technology Madras
 Chennai, India

Domain adaptation (DA) is the process in which labeled training samples available from one domain is used to improve the performance of statistical tasks performed on test samples drawn from a different domain. The domain from which the training samples are obtained is termed as the source domain, and the counterpart consisting of the test samples is termed as the target domain. Few unlabeled training samples are also taken from the target domain in order to approximate its distribution.

In this paper, we propose a new method of unsupervised DA, where a set of domain invariant sub-spaces are estimated using the geometrical and statistical properties of the source and target domains. This is a modification of the work done by Gopalan *et al.* [2], where the geodesic path from the principal components of the source to that of the target is considered in the Grassmann manifold, and the intermediary points are sampled to represent the incremental change in the geometric properties of the data in source and target domains. Instead of the geodesic path, we consider an alternate path of shortest length between the principal components of source and target, with the property that the intermediary sample points on the path form domain invariant sub-spaces using the concept of Maximum Mean Discrepancy (MMD) [3]. Thus we model the change in the geometric properties of data in both the domains sequentially, in a manner such that the distributions of projected data from both the domains always remain similar along the path. The entire formulation is done in the kernel space which makes it more robust to non-linear transformations.

Let X and Y be the source and target domains having n_X and n_Y number of instances respectively. If $\Phi(\cdot)$ is a universal kernel function, then in kernel space the source and target domains are $\Phi(X) \in \mathbb{R}^{n_X \times d}$ and $\Phi(Y) \in \mathbb{R}^{n_Y \times d}$ respectively. Let K_{XX} and K_{YY} be the kernel gram matrices of $\Phi(X)$ and $\Phi(Y)$ respectively. Let $D = [X; Y]$ denote the combined source and target domain data, and the corresponding data in kernel space is given as $\Phi(D)$. The kernel gram matrix formed using D is given by

$$K = \begin{bmatrix} K_{XX} & K_{XY} \\ K_{XY}^T & K_{YY} \end{bmatrix}, \text{ where } K_{XY} = \Phi(X)\Phi(Y)^T.$$

Let $\Phi(\tilde{X})$ and $\Phi(\tilde{Y})$ represent the projections of $\Phi(X)$ and $\Phi(Y)$ respectively onto a subspace $W_i \in \mathbb{R}^{d \times p}$, which is a point on the Grassmann manifold $G_{d,p}$. Here, d is the dimension of both source and target domains in RKHS and p is the dimension of the optimal sub-spaces. Then, the square of the distance between the means of two domains is given as:

$$\delta_\mu^2 = \text{tr} \left(W_i^T \Phi(D)^T \begin{bmatrix} I_1 & -I_2 \\ -I_2 & I_3 \end{bmatrix} \Phi(D) W_i \right) = \text{tr} \left(Z_i^T \Gamma Z_i \right) \quad (1)$$

where, $W_i = \Phi(D)^T Z_i$, $Z_i \in \mathbb{R}^{(n_X+n_Y) \times p}$, $\Gamma = \left(K \begin{bmatrix} I_1 & -I_2 \\ -I_2 & I_3 \end{bmatrix} K \right)$ and $[I_1]_{n_X \times n_X}$, $[I_2]_{n_Y \times n_X}$ and $[I_3]_{n_Y \times n_Y}$ are matrices containing all elements as $1/n_X^2$, $1/n_X n_Y$ and $1/n_Y^2$ respectively and Z_i is the unknown variable to be estimated.

If U_X^Φ and U_Y^Φ are the principal components of $\Phi(X)$ and $\Phi(Y)$ respectively, it can be proved that the principal components of $\Phi(X)$ and $\Phi(Y)U_Y^\Phi U_X^{\Phi T}$ are the same. Hence, the starting point of the path P^W is the principal components of $\Phi(D_s) = [\Phi(X); \Phi(Y)U_Y^\Phi U_X^{\Phi T}]$ and the end point of P^W can be obtained by the principal components of $\Phi(D_t) = [\Phi(X)U_X^\Phi U_Y^{\Phi T}; \Phi(Y)]$. Let, U_s^Φ and U_t^Φ be the principal components of $\Phi(D_s)$ and $\Phi(D_t)$ respectively. Also, V_X^Φ and V_Y^Φ be the eigen-vectors of K_{XX} and K_{YY} respectively. Similarly, let V_s^Φ and V_t^Φ be the eigen-vectors of K_s and K_t respectively, where K_s and K_t are the kernel gram matrices built on $\Phi(D_s)$ and $\Phi(D_t)$ respectively.

Let, G_i denote the i^{th} sampled point on the geodesic path P^G and the i^{th} sample point on P^W represent the sub-space W_i . The start and the end points of P^W are given by $W_1 = V_s^\Phi$ and $W_{N'} = V_t^\Phi$ respectively, while the intermediate points are denoted by W_i , $i = 1, \dots, N' - 1$. Now, P^W is

the path of shortest length if the sampled points from P^W is closest to the corresponding sampled points from P^G , i.e. $d_{proj}(G_i, W_i)$ is minimum, $\forall i = 2, \dots, (N' - 1)$. The square of the distance between two sub-spaces, P_i^G and P_i^W in the kernel space, is given as:

$$\delta_{proj}^2(W_i, G_i) = p - \text{tr}(Z_i^T \hat{K}_i V_i^\Phi V_i^{\Phi T} \hat{K}_i^T Z_i) = p - \text{tr}(Z_i^T \Pi_i Z_i) \quad (2)$$

where, $\Pi_i = \hat{K}_i V_i^\Phi V_i^{\Phi T} \hat{K}_i^T$. $\Phi(\hat{D}_i)$ is an appropriate projection of $\Phi(D)$. V_i^Φ is the i^{th} intermediary point sampled on the geodesic path from V_s^Φ to V_t^Φ and \hat{K}_i is the kernel gram matrix (for i^{th} sub-space in the sequence) given as $K V_i^\Phi V_i^{\Phi T} K$.

For an optimal value of Z_i , δ_{mu}^2 and $\delta_{proj}^2(G_i, W_i)$ given in Eqns. 1 and 2 should be minimum. The optimization framework to estimate Z_i is:

$$\underset{Z_i}{\text{maximize}} \quad \text{tr}(Z_i^T \Pi_i \Gamma^{-1} Z_i) \quad (3)$$

$$\text{subject to} \quad Z_i^T Z_i = I \quad (4)$$

After obtaining the set of optimal Z_i s, the projections of the data onto W_i s are given as $\Phi(D)W_i = KZ_i$, $\forall i = 2, \dots, (N' - 1)$. The projection of the data points onto the first and last (or initial and final) points of the path P^W i.e. on U_s^Φ and U_t^Φ are:

$$\Phi(D)U_s^\Phi = \Phi(D)\Phi(D_s)^T V_s^\Phi = \begin{bmatrix} K_{XX} & K_{XX} V_X^\Phi V_Y^{\Phi T} K_{YY} \\ K_{XY}^T & K_{XY}^T V_X^\Phi V_Y^{\Phi T} K_{YY} \end{bmatrix} V_s^\Phi \quad (5)$$

$$\Phi(D)U_t^\Phi = \Phi(D)\Phi(D_t)^T V_t^\Phi = \begin{bmatrix} K_{XY} V_Y^\Phi V_X^{\Phi T} K_{XX} & K_{XY} \\ K_{YY} V_Y^\Phi V_X^{\Phi T} K_{XX} & K_{YY} \end{bmatrix} V_t^\Phi \quad (6)$$

After obtaining the optimal sub-spaces, the projections of the source and target domains onto the intermediary sub-spaces are obtained and concatenated together, as done in [2], for training the KNN classifier.

We evaluate the performance of the proposed method of DA for improving the results of object categorization using Office + Caltech datasets [1]. The dataset contains four domains: Amazon (A), Caltech (C), Dslr (D) and Webcam (W), with 10 classes of objects in each of the domains. Table 1 shows the classification accuracies for 12 different pairs of source and target domains, using a 25-fold cross validation.

Table 1: Classification accuracies (in %-age) on Office+Caltech dataset [1], using different techniques of unsupervised domain adaptation.

Method	C→A	D→A	W→A	A→C	D→C	W→C
GFS [2]	36.9	32	27.5	35.3	29.4	21.7
GFK [1]	36.9	32.5	31.1	35.6	29.8	27.2
Proposed	42.63	44.16	44.65	34.40	41.56	43.26
Method	A→D	C→D	W→D	A→W	C→W	D→W
GFS [2]	30.7	32.6	54.3	31.0	30.6	66.0
GFK [1]	35.2	35.2	70.6	34.4	33.7	74.9
Proposed	38.82	43.64	80.57	39.31	42.27	78.03

The proposed method of unsupervised domain adaptation handles non-linear transformation of data as well as estimates intermediate domain invariant sub-spaces, making it more efficient.

- [1] Boqing Gong, Yuan Shi, Fei Sha, and Kristen Grauman. Geodesic flow kernel for unsupervised domain adaptation. In *CVPR*, 2012.
- [2] R. Gopalan, Ruonan Li, and Rama Chellappa. Domain adaptation for object recognition: An unsupervised approach. In *ICCV*, 2011.
- [3] Arthur Gretton, Karsten M. Borgwardt, Malte J. Rasch, Bernhard Schölkopf, and Alexander Smola. A kernel two-sample test. *Journal of Machine Learning*, 13:723–773, 2012.