# Supplementary Material – PUnDA: Probabilistic Unsupervised Domain Adaptation for Knowledge Transfer Across Visual Categories

Behnam Gholami Department of Computer Science Rutgers University, Piscataway, NJ, USA Ognjen (Oggi) Rudovic MIT Media Lab, Cambridge, MA, USA orudovic@mit.edu

bb510@cs.rutgers.edu

Vladimir Pavlovic Department of Computer Science Rutgers University, Piscataway, NJ, USA

vladimir@cs.rutgers.edu

# **1. Introduction**

This Supplement contains additional materials related to the paper **PUnDA: Probabilistic Unsupervised Domain Adapta**tion for Knowledge Transfer Across Visual Categories. In particular, in Sec. 2 we present additional results that qualitatively demonstrate the advantages of our adaptation framework over competing approaches such as the ILS [2]. We contrast 2D embeddings of the features adapted by **PUnDA** to those of the pre-adapted fc6 layer and the ILS. We include detailed derivation of the Variational Bayes algorithm for **PUnDA** in Sec. 4. Finally, in Sec. 5, we provide the computational complexity analysis of our VB algorithm.

### 2. Visualization experiments

In this Section we provide additional results of experiments described in Sec. 4 of the main paper. Specifically, we show the 2D point clouds of the learned features of **PUnDA**, embedded using the t-SNE algorithm [3], and compare them to those of the ILS features and the standard, non-adapted features. Results will demonstrate that, as supported by the quantitative results in the main paper, joint optimization of the alignment and classification criteria, accomplished through **PUnDA**, leads to qualitatively superior domain matching **and** class separability, compared to competing approaches.

#### 2.1. Office + Caltech10 Results

We refer here to the setting of the Office + Caltech10 experiments introduced in the main paper in Sec. 4.2. Figs 1,2,3 depict the embedded features extracted from the Office+Caltech10 dataset for the cases of  $A \rightarrow C$ ,  $W \rightarrow C$  and  $D \rightarrow A$  adaptations, respectively. In the figures, blue and red colors indicate the source and target domains, respectively. Colors in the bottom rows correspond to different class instances.

The top row of the Figures illustrates how the features extracted using the MMD criteria (**PUnDA**, ILS) reduce the domain mismatch. Good features for domain adaptation should have a configuration where the red and blue colors are mixed. This effect can be seen in features extracted from the **PUnDA** (b) and ILS algorithms (c), which indicates that the domain mismatch is successfully reduced in the feature space, compared to the pre-aligned features in (a). Note that the domains obtained by **PUnDA** are more compact than those of the ILS, with increasingly matched source-target features.

In classification, good domain-adapted features should display large class separability. The bottom row highlights a major difference between **PUnDA** features and the original features and the ILS features in terms of the class separability: the **PUnDA** features are more clustered with respect to the classes than ILS features, with more prominent gaps among clusters. This is partly due to the fact that **PUnDA** exploits the unlabeled data to learn the classifier boundaries for the source domain adapted discriminatively to the target domain by minimizing the expected classification errors on the target domain, something that ILS fails to account for. As a consequence, **PUnDA** framework leads to superior cross-domain classification performance.



Figure 1. Feature visualization for embedding of  $A \rightarrow C$  (Office+Caltech10) data samples using t-sne algorithm. The top and bottom rows show the domains and classes respectively. (a,d) Original features. (b,e) **PUnDA** features. (c,f) ILS features.



Figure 2. Feature visualization. Embedding of  $W \rightarrow C$  (Office+Caltech10) data samples using t-sne algorithm. The top and bottom rows show the domains and classes respectively. (a,d) Original features. (b,e) **PUnDA** features. (c,f) ILS features.

# **3. Multi-PIE Results**

We refer here to the setting of the Multi-PIE experiments introduced in the main paper in Sec. 4.3. Figs 4 and 5 depict the embedded features extracted from the Multi-PIE dataset for the cases of  $C27 \rightarrow C05$ , and  $C27 \rightarrow C37$  adaptations, respectively. In the figures, blue and red colors indicate the source and target domains, respectively. Colors in the bottom rows correspond to different class instances.

The top row of the Figures illustrates how the features extracted using the MMD criteria (**PUnDA**, ILS) reduce the domain mismatch. Again, it can be seen that the domains obtained by **PUnDA** are more compact than those of the ILS, with increasingly matched source-target features.

The bottom row highlights a difference between **PUnDA** features and the original features and the ILS features in terms of the class separability: For the Multi-PIE dataset both the **PUnDA** features and **ILS** features are slightly more clustered with



Figure 3. Feature visualization for embedding of  $D \rightarrow A$  (Office+Caltech10) data samples using t-sne algorithm. The top and bottom rows show the domains and classes respective). (a,d) Original features. (b,e) **PUnDA** features. (c,f) ILS features.



Figure 4. Feature visualization for embedding of  $C27 \rightarrow C05$  (Multi-PIE) data samples using t-sne algorithm. The top and bottom rows show the domains and classes respective). (a,d) Original features. (b,e) **PUnDA** features. (c,f) ILS features.

respect to the classes than the original features.

# 4. VB Algorithm for PUnDA

In this section we describe, in detail, the Variational Bayes algorithm at the core of **PUnDA**. We first explicitly define the model and the objective function and then show how to optimize this function by tackling two key challenges of (i) non-conjugacy between the regularizer and the Gaussian distribution used in the model, and (ii) the non-conjugacy between the softmax function and the Gaussian distribution. The optimization approach, formulated as an EM algorithm, utilizes two bounds on the expectation of the softmax function described below.



Figure 5. Feature visualization for embedding of  $C27 \rightarrow C37$  (Multi-PIE) data samples using t-sne algorithm. The top and bottom rows show the domains and classes respective). (a,d) Original features. (b,e) **PUnDA** features. (c,f) ILS features.

The proposed model can be formally defined as

$$P(s_i) \sim \mathcal{N}(0, I_K), \quad i = 1, \dots, N \tag{1}$$

$$P(s'_j) \sim \mathcal{N}(0, I_K), \quad j = 1, ..., M$$
 (2)

$$P(\pi_k|a, b, K) \sim Beta(a/K, b(K-1)/K), \quad k = 1, ..., K$$
 (3)

$$P(\phi_k) \sim \mathcal{N}(0, I_d), \quad k = 1, \dots, K \tag{4}$$

$$P(\phi'_k) \sim \mathcal{N}(0, I_d), \quad k = 1, \dots, K \tag{5}$$

$$P(z_k|\pi_k) \sim Bernoulli(\pi_k), \quad k = 1, ..., K$$
(6)

$$P(w_c) \sim \mathcal{N}(0, I_K), \quad c = 1, ..., C \tag{7}$$

$$P(\gamma_s|c_1, d_1) \sim Gamma(c_1, d_1) \tag{8}$$

$$P(\gamma_t | c'_1, d'_1) \sim Gamma(c'_1, d'_1) \tag{9}$$

$$P(x_i|\Phi, Z, s_i, \gamma_s) \sim \mathcal{N}(\Phi^\top (Z \odot s_i), \gamma_s^{-1} I_d), \quad i = 1, ..., N$$
(10)

$$P(x'_j | \boldsymbol{\Phi}', Z, s'_i, \gamma_s) \sim \mathcal{N}(\boldsymbol{\Phi}'^\top (Z \odot s'_i), \gamma_t^{-1} I_d), \quad j = 1, ..., M$$
(11)

$$P(y_i = c | \boldsymbol{W}, Z, s_i) = \frac{e^{w_c^\top (Z \odot s_i)}}{\sum_{c'} e^{w_{c'}^\top (Z \odot s_i)}}, \quad i = 1, ..., N, \ c = 1, ..., C$$
(12)

For our framework to yield a computationally effective inference method, we use a factorized variational distribution with the following forms:

$$q(s_i) \sim \mathcal{N}(\mu_i^s, \beta_{s_i}^{-1}I), \quad i = 1, ..., N$$
 (13)

$$q(s'_j) \sim \mathcal{N}(\mu_j^{s'}, \beta_{s'_j}^{-1}I) \quad j = 1, ..., M$$
(14)

$$q(\pi_k) \sim Beta(a_k, b_k), \quad k = 1, \dots, K$$
(15)

$$q(\phi_k) \sim \mathcal{N}(\mu_k^{\phi}, \beta_{\phi_k}^{-1} I_d), \quad k = 1, \dots, K$$

$$(16)$$

$$q(\phi'_k) \sim \mathcal{N}(\mu_k^{\phi'}, \beta_{\phi'_k}^{-1} I_d), \quad k = 1, ..., K$$
 (17)

$$q(z_k) \sim Bernoulli(\rho_k), \quad k = 1, ..., K$$
 (18)

$$q(w_c) \sim \mathcal{N}(\mu_c^w, \beta_{w_c}^{-1} I_K), \quad c = 1, ..., C$$
 (19)

$$q(\gamma_t) \sim Gamma(c_t, d_t) \tag{20}$$

$$q(\gamma_s) \sim Gamma(c_s, d_s) \tag{21}$$

We need to solve the following optimization problem

$$\boldsymbol{\Delta}^{*} = \underset{\boldsymbol{\Delta}}{\arg\max} \mathbb{E}_{q} \Big[ \log(\boldsymbol{X}, \boldsymbol{Y}, \boldsymbol{X}', \boldsymbol{\Omega} | \boldsymbol{\Theta}) \Big] + \boldsymbol{H}[q(\boldsymbol{\Omega})] - \lambda \mathcal{L}(\boldsymbol{S}, \boldsymbol{S}') + \lambda' \mathcal{L}'(\boldsymbol{S}', \boldsymbol{W}, \boldsymbol{Z}),$$
(22)

where  $\mathbf{\Delta} = \{\{\mu_i^s, \beta_{s_i}\}, \{\mu_j^s, \beta_{s'_j}\}, \{a_k, b_k\}, \{\mu_k^{\phi}\}, \beta_{\phi}, \{\mu_k^{\phi'}\}, \beta_{\phi'}, \{\rho_k\}, \{\mu_c^w\}, \beta_w, c_t, d_t, c_s, d_s\}.$ 

The key terms in this objective are defined as follows.  $\mathbb{E}_q \big[ \log(X, Y, X', \Omega | \Theta) \big]$  can be expanded as

$$\mathbb{E}_{q}\left[\log(\boldsymbol{X}, \boldsymbol{Y}, \boldsymbol{X}', \boldsymbol{\Omega}|\boldsymbol{\Theta})\right] = -\sum_{i} \frac{1}{2} \mathbb{E}_{q}[\gamma_{s}||x_{i} - \boldsymbol{\Phi}^{\top}(\boldsymbol{Z} \odot s_{i})||^{2}] + \frac{dN}{2}\log(\gamma_{s}) + \frac{dM}{2}\log(\gamma_{t}) \\
-\sum_{j} \frac{1}{2} \mathbb{E}_{q}[\gamma_{s}||x_{j}' - \boldsymbol{\Phi}'^{\top}(\boldsymbol{Z} \odot s_{j}')||^{2}] + \sum_{k} \mathbb{E}_{q}[z_{k}\log\pi_{k}] + \mathbb{E}_{q}[(1 - z_{k})\log(1 - \pi_{k})] + \sum_{i}\sum_{c} \mathbf{1}[y_{i} = c]\mathbb{E}_{q}[w_{c}^{\top}(\boldsymbol{Z} \odot s_{i})] \\
-\sum_{i}\sum_{c} \mathbf{1}[y_{i} = c]\mathbb{E}_{q}\left[\log\left(\sum_{c'} e^{w_{c'}^{\top}(\boldsymbol{Z} \odot s_{i})}\right)\right] - \frac{1}{2}\sum_{i}\mathbb{E}_{q}||s_{i}||^{2} - \frac{1}{2}\sum_{j}\mathbb{E}_{q}||s_{j}'||^{2} - \frac{1}{2}\sum_{k}\mathbb{E}_{q}||\phi_{k}||^{2} - \frac{1}{2}\sum_{k}\mathbb{E}_{q}||\phi_{k}||^{2} \\
-\frac{1}{2}\sum_{c}\mathbb{E}_{q}||w_{c}||^{2} + \sum_{k}\left(\frac{a}{K} - 1\right)\mathbb{E}[\log\pi_{k}] + \left(\frac{b(K - 1)}{K} - 1\right)\mathbb{E}[\log(1 - \pi_{k})] + (c_{1} - 1)\mathbb{E}_{q}[\log\gamma_{s}] - d_{1}\mathbb{E}_{q}[\gamma_{s}] \\
+ (c_{2} - 1)\mathbb{E}_{q}[\log\gamma_{t}] - d_{2}\mathbb{E}_{q}[\gamma_{t}].$$
(23)

 $oldsymbol{H}[q(oldsymbol{\Omega})]$  is subsequently expanded as

$$\boldsymbol{H}[q(\boldsymbol{\Omega})] = -\frac{k}{2} \sum_{i} \log \beta_{s_i} - \frac{k}{2} \sum_{j} \log \beta_{s'_j} - \frac{dK}{2} \log \beta_{\phi} - \frac{dK}{2} \log \beta_{\phi'} - \frac{KC}{2} \log \beta_w + \log \Gamma(a_k) + \log \Gamma(b_k) - \log(a_k + b_k) - (a_k - 1)\psi(a_k) - (b_k - 1)\psi(b_k) + (a_k + b_k - 2)\psi(a_k + b_k) + c_s - \log d_s + \log \Gamma(c_s) + (1 - c_s)\psi(c_s) + c_t - \log d_t + \log \Gamma(c_t) + (1 - c_t)\psi(c_t) - \sum_k \left[ (1 - \rho_k) \log \rho_k + \rho_k \log(1 - \rho_k) \right],$$
(24)

where  $\psi$  denotes the **digamma** function defined as

$$\psi(x) = \frac{d}{dx} \log \Gamma(x), \tag{25}$$

and  $\Gamma(x)$  denotes the Gamma function defined as

$$\Gamma(x) = \int_0^\infty u^{x-1} e^{-u} du.$$
(26)

 $\mathcal{L}(\boldsymbol{S}, \boldsymbol{S}')$  is computed as

$$\mathcal{L}(\boldsymbol{S}, \boldsymbol{S}') = \sum_{i,i'} \frac{\mathcal{K}(q(s_i), q(s_j))}{N^2} - 2\sum_{i,j} \frac{\mathcal{K}(q(s_i), q(s'_j))}{NM} + \sum_{j,j'} \frac{\mathcal{K}(q(s'_j), q(s'_{j'}))}{M^2},$$
(27)

where  $\mathcal{K}(q(s),q(s'))$  is computed as

$$\mathcal{K}(q(s), q(s')) = \log \int_{\mathbb{R}^K} q(s)^{1/2} q(s')^{1/2} \, ds \, ds'.$$
(28)

Since q(s) and q(s') are Gaussian distributions, we can compute (28) in closed form as

$$\mathcal{K}(q(s), q(s')) = \frac{\beta_s + \beta_{s'}}{4\beta_s \beta_{s'}} \|\mu_s - \mu_t\|^2 - \frac{K}{4} \log \beta_s - \frac{K}{4} \log \beta_{s'} + \frac{K}{2} \log(\beta_s + \beta_{s'}).$$
(29)

 $\mathcal{L}'(\mathbf{S}', \mathbf{W}, Z)$  is computed as

$$\mathcal{L}'(\boldsymbol{W}, \boldsymbol{S}', Z) = \sum_{j} \sum_{c} \mathbb{E}_{P(y'_{j}|\boldsymbol{W}, s'_{j}, Z)} \log P(y'_{j} = c)$$

$$= \sum_{j} \sum_{c} \left[ \mathbb{E}_{q}[w_{c}^{\top}(Z \odot s'_{j})] - \mathbb{E}_{q}[\log \sum_{c'} exp(w_{c'}^{\top}(Z \odot s'_{j}))] \right] \left[ \frac{exp(w_{c}^{\top}(Z \odot s'_{j}))}{\sum_{c'} exp(w_{c'}^{\top}(Z \odot s'_{j}))]} \right].$$
(30)

There are two main issues which make deriving the VB update equations in closed-form intractable: (i) Non-conjugacy between the regularizer  $\mathcal{L}'(S', W, Z)$  and the Gaussian distribution, and (ii) Non-conjugacy between the softmax function and the Gaussian distribution.

To tackle the first issue, it is worth noting that  $\mathcal{L}'(\boldsymbol{W}, \boldsymbol{S}', Z)$  is equivalent to the sum of the negative Shanon Entropy of the *C*-dimensional probability vectors  $P(y'_j | \boldsymbol{W}, Z, s_j) = [p_j^1, ..., p_j^C]$ . Hence, we propose to use Renyi entropy [4] rather than Shanon entropy for defining the regularizer  $\mathcal{L}'(\boldsymbol{S}', \boldsymbol{W}, Z)$ . The Renyi entropy of order  $\alpha$ , where  $\alpha \ge 0$  and  $\alpha \ne 1$  is defined as

$$\boldsymbol{H}_{\alpha}[\boldsymbol{y}_{j}'] = \frac{1}{1-\alpha} \log \left( \sum_{c=1}^{C} (p_{c}^{j})^{\alpha} \right)$$
(31)

The intuition behind using the Renyi entropy is that the limiting value of  $H_{\alpha}[y'_j]$  as  $\alpha \to 1$  is the Shanon entropy [4]. We set  $\alpha = 0.99$  in the experiments to approximate the Shanon entropy with the Renyi entropy. By replacing the Shanon entropy with the Renyi entropy in (30), we have

$$\mathcal{L}'(\boldsymbol{W}, \boldsymbol{S}', \boldsymbol{Z}) = \frac{1}{1 - \alpha} \sum_{j} \mathbb{E}_{q} [\log \sum_{c} e^{\alpha w_{c}^{\top}(\boldsymbol{Z} \odot s_{j}')}] - \frac{\alpha}{1 - \alpha} \sum_{j} \mathbb{E}_{q} [\log \sum_{c} e^{w_{c}^{\top}(\boldsymbol{Z} \odot s_{j}')}].$$
(32)

To tackle the second non-conjugacy problem, we use a linear lower bound and a quadratic upper bound on the expectation of the log-sum function introduced in [1]:

$$\log \sum_{c=1}^{C} e^{w_c^{\top}(Z \odot s)} \ge \sum_{c=1}^{C} w_c^{\top}(Z \odot s)$$
(33)

$$\log \sum_{c=1}^{C} e^{w_{c}^{\top}(Z \odot s)} \leq \sum_{c=1}^{C} \left[ \eta(\xi_{c}) \left( (w_{c}^{\top}(Z \odot s))^{2} - \xi_{c}^{2} \right) - \log \sigma(\xi_{c}) \right] + \frac{1}{2} \sum_{c=1}^{C} w_{c}^{\top}(Z \odot s) + \xi_{c}, \tag{34}$$

where

$$\eta(\xi_c) = -\frac{1}{2\xi_c} \left( \frac{1}{1 + e^{-\xi_c}} - \frac{1}{2} \right)$$
(35)

and  $\{\xi_c \in [0, +\infty]\}$  denote the free variational parameters which are optimized to get the tightest possible bound. Hence, we replace  $\log \sum_c \exp(\alpha w_c^{\top}(Z \odot s'_j))$  with its lower bound and  $\log \sum_{c'} \exp(w_{c'}^{\top}(Z \odot s_i))$  and  $\log \sum_c \exp(w_c^{\top}(Z \odot s'_j))$  with their upper bound in the objective function, then we use the EM algorithm to optimize the factorized variational distribution and the free parameters which computes the variational posterior distribution in the E-step and maximizes the free parameters in the M-step, which goes as follows.

## E step:

In this step the free variational parameters  $\{\xi_c\}$  are fixed, and the variational distributions are updated by maximizing the objective function using a coordinate ascent algorithm.

#### Update for $\gamma_s$ :

One can show that the posterior parameters  $c_s, d_s$  can be updated as

$$c_s = c_1 + \frac{Nd}{2}, \quad d_s = d_1 + \frac{1}{2} \sum_{i=1}^d \sum_{j=1}^N \epsilon_{ij},$$
(36)

where

$$\epsilon_{ij} = x_{ij}^2 - 2x_{ij} \langle \phi_k \rangle^\top (\langle Z \rangle \odot \langle s_j \rangle) + \left( \langle \phi_k \rangle^\top (\langle Z \rangle \odot \langle s_j \rangle) \right)^2 + \sum_{k=1}^K \langle \phi_{ik}^2 \rangle (\langle z_k^2 \rangle \odot \langle s_{kj}^2 \rangle) - \sum_{k=1}^K \langle \phi_{ik} \rangle^2 (\langle z_k \rangle^2 \odot \langle s_{kj} \rangle^2)$$
(37)

and  $\langle . \rangle$  indicates the expectation operator. Update for  $\gamma_t$ :

One can show that the posterior parameters  $c_t, d_t$  can be updated as

$$c_t = c'_1 + \frac{Md}{2}, \quad d_s = d'_1 + \frac{1}{2} \sum_{i=1}^d \sum_{j=1}^M \epsilon_{ij},$$
(38)

where

$$\epsilon_{ij} = x_{ij}^{\prime}{}^2 - 2x_{ij}^{\prime}\langle\phi_k^{\prime}\rangle^{\top} (\langle Z \rangle \odot \langle s_j^{\prime}\rangle) + \left(\langle\phi_k^{\prime}\rangle^{\top} (\langle Z \rangle \odot \langle s_j^{\prime}\rangle)\right)^2 + \sum_{k=1}^K \langle\phi_{ik}^{\prime}{}^2\rangle (\langle z_k^2 \rangle \odot \langle s_{kj}^{\prime}{}^2\rangle) - \sum_{k=1}^K \langle\phi_{ik}^{\prime}\rangle^2 (\langle z_k \rangle^2 \odot \langle s_{kj}^{\prime}\rangle^2)$$
(39)

Update for  $\Pi = [\pi_1, ..., \pi_K]$ :

One can show that the posterior parameters  $a_k, b_k$  can be updated as

$$a_k = \frac{a}{K} + \langle z_k \rangle, \quad b_k = \frac{b(K-1)}{K} + 1 - \langle z_k \rangle. \tag{40}$$

Update for  $\mathbf{\Phi} = [\phi_1, ..., \phi_K]$ :

One can show that the posterior parameters  $\mu_k^{\phi}, \beta_{\phi_k}$  can be updated as

$$\beta_{\phi_k} = \langle \gamma_s \rangle \sum_{i=1}^N \langle z_k \rangle \langle s_{ki}^2 \rangle + 1/2, \qquad \mu_k^{\phi} = \frac{\langle \gamma_s \rangle}{\beta_{\phi_k}} \sum_{i=1}^N \langle z_k \rangle \langle s_{ki} \rangle \langle x_i^{-k} \rangle, \tag{41}$$

where  $x_i^{-k}$  is defined as

$$x_i^{-k} \equiv x_i - \Phi_{-k}^{\top}(z^{-k} \odot s_i^{-k})$$
(42)

and  $\Phi_{-k}, z^{-k}$  and  $s_i^{-k}$  are the matrix/vectors with the k-th column/element removed. Update for  $\Phi' = [\phi'_1, ..., \phi'_K]$ :

One can show that the posterior parameters  $\mu_k^{\phi'}, \beta_{\phi'_k}$  can be updated as

$$\beta_{\phi'_{k}} = \langle \gamma_{s} \rangle \sum_{j=1}^{M} \langle z_{k} \rangle \langle s'_{kj}{}^{2} \rangle + 1/2, \qquad \mu_{k}^{\phi'} = \frac{\langle \gamma_{t} \rangle}{\beta_{\phi'_{k}}} \sum_{j=1}^{M} \langle z_{k} \rangle \langle s_{kj} \rangle \langle x'_{j}{}^{-k} \rangle$$
(43)

where  $x'_{j}^{-k}$  is defined as

$$x'_{j}^{-k} \equiv x'_{j} - \Phi'_{-k}^{\top} (z^{-k} \odot s'_{j}^{-k})$$
(44)

**Update for**  $Z = [z_1, ..., z_K]$ :

One can show that the posterior parameter  $\rho_k$  can be updated as

$$\rho_k = \left(1 + e^{\langle \log \pi_k \rangle + \langle \log(1 - \pi_k) \rangle + \tau_s + \tau_t}\right)^{-1},\tag{45}$$

where

$$\tau_{s} = -\frac{\langle \gamma_{s} \rangle}{2} \left[ -2 \sum_{i=1}^{N} \langle s_{ki} \rangle \sum_{j=1}^{d} \left( \langle \phi_{jk} \rangle \sum_{k' \neq k}^{K} \langle \phi_{jk'} \rangle \langle s_{k'j} \rangle \langle z_{k'} \rangle \right) \right] + \sum_{i=1}^{N} \sum_{c=1}^{C} (1[y_{i} = c] - \frac{1}{2}) \langle w_{c} \rangle \langle s_{ki} \rangle \\ - \sum_{i=1}^{N} \sum_{c=1}^{C} \eta(\xi_{c}) \left[ \langle w_{c}^{2} \rangle \langle s_{ki}^{2} \rangle + 2 \langle w_{c} \rangle \langle s_{ki} \rangle \langle w_{c}^{-k} \rangle (\langle Z^{-k} \rangle \odot \langle s_{i}^{-k} \rangle) \right]$$

$$(46)$$

$$\tau_{t} = -\frac{\langle \gamma_{t} \rangle}{2} \bigg[ -2 \sum_{j=1}^{M} \langle s'_{kj} \rangle \sum_{i=1}^{d} \big( \langle \phi'_{ik} \rangle \sum_{k' \neq k}^{K} \langle \phi'_{ik'} \rangle \langle s'_{k'j} \rangle \langle z_{k'} \rangle \big) \bigg] + \sum_{j=1}^{M} \sum_{c=1}^{C} \frac{\lambda' \alpha}{2 - 2\alpha} \langle w_{c} \rangle \langle s_{ki} \rangle \\ - \frac{\lambda' \alpha}{1 - \alpha} \sum_{j=1}^{M} \sum_{c=1}^{C} \eta(\xi_{c}) \bigg[ \langle w_{c}^{2} \rangle \langle s'_{kj}^{2} \rangle + 2 \langle w_{c} \rangle \langle s'_{ki} \rangle \langle w_{c}^{-k} \rangle (\langle Z^{-k} \rangle \odot \langle s'_{j}^{-k} \rangle) \bigg]$$

$$(47)$$

and  $\langle \log \pi_k \rangle$  and  $\langle \log (1-\pi_k) \rangle$  can be computed as

$$\langle \log \pi_k \rangle = \psi(\frac{a}{K} + \langle z_k \rangle) - \psi(\frac{a + b(K - 1)}{K} + 1)$$
(48)

$$\langle \log(1-\pi_k) \rangle = \psi(\frac{b(K-1)}{K} + 1 - \langle z_k \rangle) - \psi(\frac{a+b(K-1)}{K} + 1).$$
(49)

Update for  $S = [s_1, ..., s_N]$ : One can show that the posterior parameter  $\mu_i^s, \beta_{s_i}$  can be updated as

$$\beta_{s_i} = \langle \gamma_s \rangle \langle z_k \rangle \sum_{j=1}^d \langle \phi_{jk}^2 \rangle + 2 \bigg( \sum_{c=1}^C \eta(\xi_c) (\langle w_c \rangle^\top \langle z_k \rangle) \bigg) + 1/2$$
(50)

$$\begin{aligned} (\mu_{i}^{s})_{k} &= \frac{\langle \gamma_{s} \rangle}{\beta_{s_{i}} + \lambda a} \sum_{j=1}^{d} \left( \langle \phi_{jk} \rangle \langle z_{k} \rangle \left( x_{ji} - \sum_{k' \neq k}^{K} \langle \phi_{jk'} \rangle \langle z_{k'}' \rangle \right) \right) + \frac{1}{\beta_{s_{i}}} \sum_{c=1}^{C} \mathbf{1}[y_{i} = c] \langle w_{kc} \rangle \langle z_{k} \rangle \\ &- \frac{2}{\beta_{s_{i}}} \sum_{c=1}^{C} \eta(\xi_{c}) \left( \langle w_{kc} \rangle \langle z_{k} \rangle \sum_{k' \neq k}^{K} \langle w_{k'c} \rangle \langle z_{k'}' \rangle \langle s_{k'i} \rangle \right) - \frac{1}{2\beta_{s_{i}}} \left( \langle z_{k} \rangle \sum_{c=1}^{C} \langle w_{kc} \rangle \right) \\ &+ \left( \frac{\lambda}{N^{2}} \sum_{i' \neq i}^{N} \frac{\beta_{s_{i}} + \beta_{s_{i'}}}{\beta_{s_{i}} \beta_{s_{i'}}} (\mu_{i'}^{s})_{k} \right) - \left( \frac{2\lambda}{MN} \sum_{j=1}^{M} \frac{\beta_{s_{i}} + \beta_{s'_{j}}}{\beta_{s_{i}} \beta_{s'_{j}}} (\mu_{j'}^{s'})_{k} \right) \end{aligned}$$
(51)

where  $(\mu_i^s)_k$  denotes the k-th entry of the vector  $(\mu_i^s)_k,$  and

$$a = \frac{1}{N^2} \sum_{i'=1}^{N} \frac{\beta_{s_i} + \beta_{s_{i'}}}{\beta_{s_i} \beta_{s_{i'}}} - \frac{2}{MN} \sum_{j=1}^{M} \frac{\beta_{s_i} + \beta_{s_j}}{\beta_{s_i} \beta_{s_j}}$$
(52)

Update for  $S' = [s'_1, ..., s'_N]$ : One can show that the posterior parameter  $\mu_i^{s'}, \beta_{s'_i}$  can be updated as

$$\beta_{s'_i} = \langle \gamma_t \rangle \langle z_k \rangle \sum_{j=1}^d \langle \phi'_{jk} \rangle^2 + \frac{2\lambda'\alpha}{1-\alpha} \left( \sum_{c=1}^C \eta(\xi_c) (\langle w_c \rangle^\top \langle z_k \rangle) \right) + 1/2$$
(53)

$$(\mu_{i}^{s'})_{k} = \frac{\langle \gamma_{t} \rangle}{\beta_{s'_{i}} + \lambda a} \sum_{j=1}^{d} \left( \langle \phi'_{jk} \rangle \langle z_{k} \rangle \left( x'_{ji} - \sum_{k' \neq k}^{K} \langle \phi'_{jk'} \rangle \langle z'_{k} \rangle \langle s'_{k'i} \rangle \right) \right) + \frac{\lambda' \alpha}{\beta_{s'_{i}}(1 - \alpha)} \sum_{c=1}^{C} \langle w_{kc} \rangle \langle z_{k} \rangle$$
$$- \frac{2\lambda' \alpha}{\beta_{s'_{i}}(1 - \alpha)} \sum_{c=1}^{C} \eta(\xi_{c}) \left( \langle w_{kc} \rangle \langle z_{k} \rangle \sum_{k' \neq k}^{K} \langle w_{k'c} \rangle \langle z'_{k} \rangle \langle s'_{k'i} \rangle \right) - \frac{\lambda' \alpha}{2\beta_{s'_{i}}(1 - \alpha)} \left( \langle z_{k} \rangle \sum_{c=1}^{C} \langle w_{kc} \rangle \right)$$
$$+ \left( \frac{\lambda}{M^{2}} \sum_{i' \neq i}^{M} \frac{\beta_{s'_{i}} + \beta_{s'_{i'}}}{\beta_{s'_{i}} \beta_{s'_{i'}}} (\mu_{i'}^{s'})_{k} \right) - \left( \frac{2\lambda}{MN} \sum_{j=1}^{N} \frac{\beta_{s'_{i}} + \beta_{s_{j}}}{\beta_{s'_{i}} \beta_{s_{j}}} (\mu_{j}^{s})_{k} \right) \right)$$
(54)

**Update for**  $W = [w_1, ..., w_C]$ :

One can show that the posterior parameter  $\mu_c^w, \beta_{w_c}$  can be updated as

$$\beta_{w_c} = 2\sum_{i=1}^N \lambda(\xi_c) trace(F_i) + \frac{2\lambda'\alpha}{1-\alpha} \sum_{j=1}^M \lambda(\xi_c) trace(F'_j) + \frac{1}{2}$$
(55)

where  $F_i = \langle (Z \odot s_i) (Z \odot s_i)^\top$  is a  $K \times K$  matrix which its elements are defined as

$$F_{i}(m,n) = \langle z_{m} \rangle \langle s_{mi}^{2} \rangle, \quad \text{if } m = n$$
  

$$F_{i}(m,n) = \langle z_{m} \rangle \langle z_{n} \rangle \langle s_{mi} \rangle, \quad \text{if } m \neq n$$
(56)

and where  $F'_j = \langle (Z \odot s'_j) (Z \odot s'_j)^\top$  is a  $K \times K$  matrix which its elements are defined as

$$F'_{j}(m,n) = \langle z_{m} \rangle \langle {s'}_{mj}^{2} \rangle, \text{ if } m = n$$
  

$$F'_{j}(m,n) = \langle z_{m} \rangle \langle z_{n} \rangle \langle {s'}_{mj} \rangle \langle {s'}_{nj} \rangle, \text{ if } m \neq n$$
(57)

$$\mu_c^w = \left[\operatorname{diag}(\beta_{w_c})\right]^{-1} \left(\sum_{i=1}^N (\mathbf{1}[y_i = c] - \frac{1}{2})(\langle Z \rangle \odot \langle s_i \rangle) + \frac{\lambda' \alpha}{2 - 2\alpha} \sum_{j=1}^M (\langle Z \rangle \odot \langle s'_j \rangle)\right)$$
(58)

where diag $(\beta_{w_c})$  is a  $K \times K$  diagonal matrix with entries taken from the elements  $\beta_{w_c}$ .

## M step:

In this step, the variational free parameters  $\{\xi_c\}$  are computed by maximizing the objective function while keeping the parameters of the posterior distribution  $Q(\Omega)$  fixed.

**Update for**  $\xi = [\xi_1, ..., \xi_C]$ :

it is easy to show that  $\xi_c$  can be computed as

$$\xi_c^2 = \operatorname{trace}\left[\sum_{i=1}^N F_i + \frac{\lambda'\alpha}{1-\alpha} \sum_{j=1}^M F_j \langle w_c w_c^\top \rangle\right]$$
(59)

### 5. Computational Complexity of our VB algorithm

Table 1. Computational Complexity of the PUnDA VB algorithm.

Step	Complexity
Update $\gamma_s$	O(NKd)
Update $\gamma_t$	O(MKd)
Update $\Pi$	O(K)
Update $\Phi$	$O(NK^3d))$
Update $\Phi'$	$O(MK^3d)$
Update $Z$	$O(K^2(N+M)d)$
Update $S$	O(NKd)
Update $S'$	O(MKd)
Update $W$	O(C(N+M)K)
Update <b>ξ</b>	$O(MK^2)$

In this section, we analyze the computational complexity of each iteration of the proposed VB algorithm. Each iteration in our VB algorithm, mainly includes matrix multiplications. The computational complexity of each parameter's updating is

summarized in Table 1. As can be seen from Table 1, the total complexity for VB algorithm is  $O(T \times (N + M)K^2d)$  and T is the total number of iterations. It is clear that the computational complexity of the proposed VB algorithm, for training, in one iteration is  $O((N + M)K^2d)$ , i.e., linear in the size of the source+target data N + M, the data dimensionality d, and quadratic in the dimensionality of the shared space  $K, K \ll d$ . For classifying a test data point the computational complexity is  $O(CK^2)$ , for C class instances.

## References

- G. Bouchard. Efficient bounds for the softmax function and applications to approximate inference in hybrid models. In NIPS 2007 workshop for approximate Bayesian inference in continuous/hybrid systems, 2007. 6
- [2] S. Herath, M. Harandi, and F. Porikli. Learning an invariant hilbert space for domain adaptation. *IEEE Conference on Computer Vision and Pattern Recognition*, 2017. 1
- [3] L. v. d. Maaten and G. Hinton. Visualizing data using t-sne. Journal of Machine Learning Research, 9(Nov):2579-2605, 2008. 1
- [4] R. Renner and S. Wolf. Smooth rényi entropy and applications. In *Information Theory*, 2004. ISIT 2004. Proceedings. International Symposium on, page 233. IEEE, 2004. 6