

# Optimal Bayesian Transfer Learning

Xiaoning Qian,  
Alireza Karbalayghareh, and Edward R Dougherty

Texas A&M University



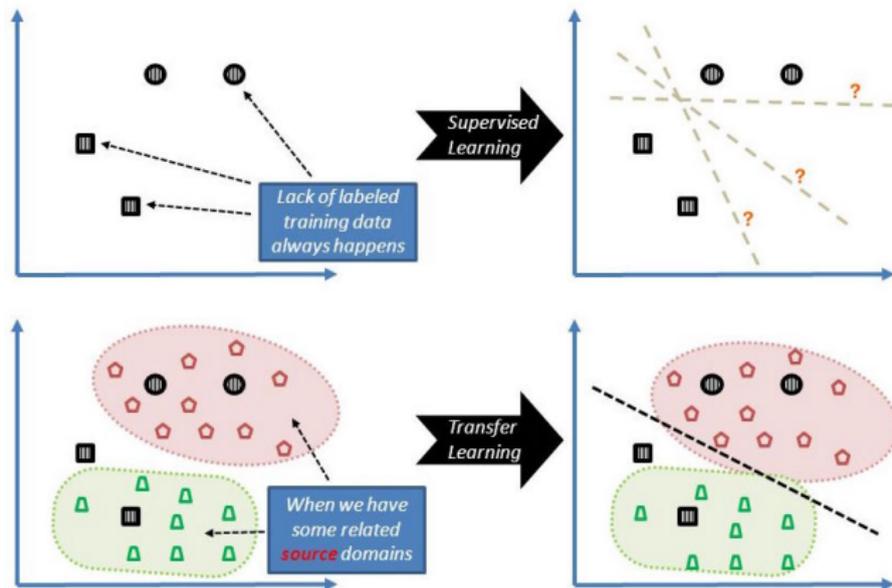
August 7, 2018

# Outline

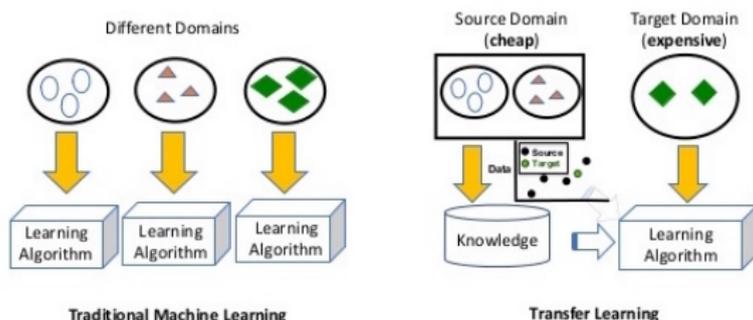
- 1 Transfer Learning
- 2 Optimal Bayesian Classifier
- 3 Optimal Bayesian Transfer Learning for Multivariate Gaussian Data
- 4 Optimal Bayesian Transfer Learning for Count Data
- 5 Conclusions
- 6 References

# Transfer Learning Basics

- Traditional machine learning vs. transfer learning



# Transfer Learning Basics



- Suppose we want to do a supervised learning but there is lack of labeled data in the domain of interest (**target domain**).
- Therefore, the classifier cannot be trained well and error rate would be high.
- At the same time, suppose we have plenty of labeled data in a different but **relevant** domain (**source domain**).
- The problem of transfer learning is to answer **when** and **how** to employ those source data in order to design a more accurate classifier in the target domain.

# Domain Adaptation

- Distributions of source and target data are different (not **i.i.d.** as in traditional machine learning).
- **Domain adaptation** [1] aims to find a common domain where both source and target data can be transformed to have similar distributions.
- Often, transformation is forced to source and target data but no theoretical guarantee that the prediction performance in the target domain will be enhanced.
- There is no rigorous reasoning for “transferability” and it does not answer if the two domains are actually relevant.
- More critically, is there a way to optimally transfer the relevant knowledge and data from source to target?

# Optimal Bayesian Classifier

- Feature-label distribution:  $p(\mathbf{x}, l|\theta) = p(\mathbf{x}|l; \theta)p(l|\theta)$
- Prior distribution:  $p(\theta)$
- Likelihood:  $p(\mathcal{D}|\theta) = \prod_n p(\mathbf{x}^n, l^n|\theta)$
- Posterior:  $p(\theta|\mathcal{D}) = \frac{p(\theta) \prod_n p(\mathbf{x}^n, l^n|\theta)}{p(\mathcal{D})}$
- Posterior predictive (effective class-conditional) distribution given a new feature vector  $\mathbf{x}^*$ :  $p(x^*|l; \mathcal{D}) \propto \int d\theta p(x^*|l; \theta)p(\theta|\mathcal{D})$
- Optimal Bayesian classifier:

$$\arg \max_{l \in \{1, \dots, L\}} p(\theta_l|\mathcal{D})p(x^*|l; \mathcal{D})$$

# Bayesian Transfer Learning

- We formulate a Bayesian transfer learning framework to transfer source domain knowledge and data for learning in target domain.
- Our Bayesian framework directly models the **feature-label distributions** in source and target domains.
- The “transferability” across domains can be characterized by a **joint prior distribution** on model parameters of feature-label distributions across domains.
- The relevance of source and target problems can be studied through the **joint posterior distribution** of model parameters.
- Under such a Bayesian framework, we show how to **optimally** transfer abundant source data to the target domain and define the Optimal Bayesian Transfer Learning (OBTL) classifier.

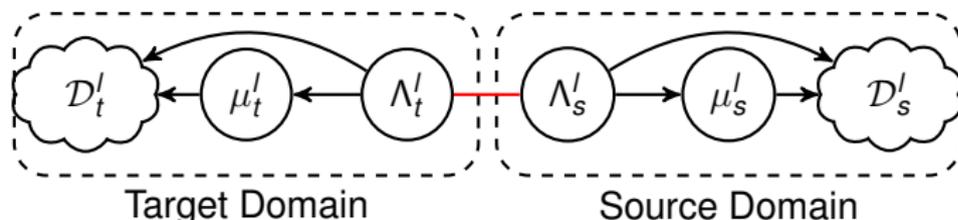
# Multivariate Gaussian Data

- Distributions of data in source and target domains:

$$\mathbf{x}_s^l \sim \mathcal{N}(\mu_s^l, (\Lambda_s^l)^{-1}), \quad \mathbf{x}_t^l \sim \mathcal{N}(\mu_t^l, (\Lambda_t^l)^{-1}), \quad l \in \{1, \dots, L\},$$

- Joint prior for the parameters of the two domains:

$$p(\mu_s^l, \mu_t^l, \Lambda_s^l, \Lambda_t^l) = p(\mu_s^l | \Lambda_s^l) p(\mu_t^l | \Lambda_t^l) p(\Lambda_s^l, \Lambda_t^l), \quad l \in \{1, \dots, L\},$$
$$\mu_s^l | \Lambda_s^l \sim \mathcal{N}(\mathbf{m}_s^l, (\kappa_s^l \Lambda_s^l)^{-1}), \quad \mu_t^l | \Lambda_t^l \sim \mathcal{N}(\mathbf{m}_t^l, (\kappa_t^l \Lambda_t^l)^{-1}),$$



- In the case of one domain, Wishart matrices are used for a conjugate prior for the distribution of precision matrices.
- The main question here is: **how to define a joint distribution between two Wishart matrices  $p(\Lambda_s^l, \Lambda_t^l)$ ?**

## Theorem ([2])

If  $\Lambda \sim W_d(\mathbf{M}, \nu)$ , and  $\mathbf{A}$  is an  $r \times d$  matrix of rank  $r$ , where  $r \leq d$ , then  $\mathbf{A}\Lambda\mathbf{A}' \sim W_r(\mathbf{A}\mathbf{M}\mathbf{A}', \nu)$ .

## Corollary

If  $\Lambda \sim W_d(\mathbf{M}, \nu)$  and  $\Lambda = \begin{pmatrix} \Lambda_{11} & \Lambda_{12} \\ \Lambda_{12}' & \Lambda_{22} \end{pmatrix}$ , where  $\Lambda_{11}$  and  $\Lambda_{22}$  are  $d_1 \times d_1$  and  $d_2 \times d_2$  submatrices, respectively, and if  $\mathbf{M} = \begin{pmatrix} \mathbf{M}_{11} & \mathbf{M}_{12} \\ \mathbf{M}_{12}' & \mathbf{M}_{22} \end{pmatrix}$  is the corresponding partition of  $\mathbf{M}$  with  $\mathbf{M}_{11}$  and  $\mathbf{M}_{22}$  being two  $d_1 \times d_1$  and  $d_2 \times d_2$  submatrices, respectively, then  $\Lambda_{11} \sim W_{d_1}(\mathbf{M}_{11}, \nu)$  and  $\Lambda_{22} \sim W_{d_2}(\mathbf{M}_{22}, \nu)$ .

## Theorem ([3])

Let  $\Lambda = \begin{pmatrix} \Lambda_{11} & \Lambda_{12} \\ \Lambda_{12}' & \Lambda_{22} \end{pmatrix}$  be a  $(d_1 + d_2) \times (d_1 + d_2)$  partitioned Wishart random matrix, where the diagonal partitions are of sizes  $d_1 \times d_1$  and  $d_2 \times d_2$ , respectively. The Wishart distribution of  $\Lambda$  has  $\nu \geq d_1 + d_2$  degrees of freedom and positive-definite scale matrix  $\mathbf{M} = \begin{pmatrix} \mathbf{M}_{11} & \mathbf{M}_{12} \\ \mathbf{M}_{12}' & \mathbf{M}_{22} \end{pmatrix}$  partitioned in the same way as  $\Lambda$ . The joint distribution of the two diagonal partitions  $\Lambda_{11}$  and  $\Lambda_{22}$  have the density function given by

$$\begin{aligned} p(\Lambda_{11}, \Lambda_{22}) &= K \operatorname{etr} \left( -\frac{1}{2} \left( \mathbf{M}_{11}^{-1} + \mathbf{F}' \mathbf{C}_2 \mathbf{F} \right) \Lambda_{11} \right) \operatorname{etr} \left( -\frac{1}{2} \mathbf{C}_2^{-1} \Lambda_{22} \right) \\ &\times |\Lambda_{11}|^{\frac{\nu-d_2-1}{2}} |\Lambda_{22}|^{\frac{\nu-d_1-1}{2}} {}_0F_1 \left( \frac{\nu}{2}; \frac{1}{4} \mathbf{G} \right), \end{aligned} \quad (1)$$

where  $\mathbf{C}_2 = \mathbf{M}_{22} - \mathbf{M}_{12}' \mathbf{M}_{11}^{-1} \mathbf{M}_{12}$ ,  $\mathbf{F} = \mathbf{C}_2^{-1} \mathbf{M}_{12}' \mathbf{M}_{11}^{-1}$ ,  $\mathbf{G} = \Lambda_{22}^{\frac{1}{2}} \mathbf{F} \Lambda_{11} \mathbf{F}' \Lambda_{22}^{\frac{1}{2}}$ ,  $K^{-1} = 2^{\frac{(d_1+d_2)\nu}{2}} \Gamma_{d_1} \left( \frac{\nu}{2} \right) \Gamma_{d_2} \left( \frac{\nu}{2} \right) |\mathbf{M}|^{\frac{\nu}{2}}$ , and  ${}_0F_1$  is the generalized matrix-variate hypergeometric function.

Multivariate gamma function given by  $\Gamma_d(\alpha) = \pi^{\frac{d(d-1)}{4}} \prod_{i=1}^d \Gamma \left( \alpha - \frac{i-1}{2} \right)$ .

# Hypergeometric functions of matrix arguments

## Definition ([4])

The generalized hypergeometric function of one matrix argument is defined by

$${}_pF_q(a_1, \dots, a_p; b_1, \dots, b_q; \mathbf{X}) = \sum_{k=0}^{\infty} \sum_{\kappa \vdash k} \frac{(a_1)_{\kappa} \cdots (a_p)_{\kappa}}{(b_1)_{\kappa} \cdots (b_q)_{\kappa}} \frac{C_{\kappa}(\mathbf{X})}{k!}, \quad (2)$$

where  $a_i, i = 1, \dots, p$ , and  $b_j, j = 1, \dots, q$ , are arbitrary complex (real in our case) numbers,  $C_{\kappa}(\mathbf{X})$  is the zonal polynomial of  $d \times d$  symmetric matrix  $\mathbf{X}$  corresponding to the ordered partition  $\kappa = (k_1, \dots, k_d)$ ,  $k_1 \geq \dots \geq k_d \geq 0$ ,  $k_1 + \dots + k_d = k$  and  $\sum_{\kappa \vdash k}$  denotes summation over all partitions  $\kappa$  of  $k$ . The generalized hypergeometric coefficient  $(a)_{\kappa}$  is defined by

$$(a)_{\kappa} = \prod_{i=1}^d \left( a - \frac{i-1}{2} \right)_{k_i}, \quad (3)$$

where  $(a)_r = a(a+1) \cdots (a+r-1)$ ,  $r = 1, 2, \dots$ , with  $(a)_0 = 1$ .

# Hypergeometric functions of matrix arguments

Most special cases are:

$$\begin{aligned} {}_0F_0(\mathbf{X}) &= \sum_{k=0}^{\infty} \sum_{\kappa \vdash k} \frac{C_{\kappa}(\mathbf{X})}{k!} = \sum_{k=0}^{\infty} \frac{(\text{tr}(\mathbf{X}))^k}{k!} = \text{etr}(\mathbf{X}), \\ {}_1F_0(a; \mathbf{X}) &= \sum_{k=0}^{\infty} \sum_{\kappa \vdash k} \frac{(a)_{\kappa} C_{\kappa}(\mathbf{X})}{k!} = |\mathbf{I}_m - \mathbf{X}|^{-a}, \quad \|\mathbf{X}\| < 1, \\ {}_0F_1(b; \mathbf{X}) &= \sum_{k=0}^{\infty} \sum_{\kappa \vdash k} \frac{C_{\kappa}(\mathbf{X})}{(b)_{\kappa} k!}, \quad (\text{Confluent hypergeometric limit function}) \\ {}_1F_1(a; b; \mathbf{X}) &= \sum_{k=0}^{\infty} \sum_{\kappa \vdash k} \frac{(a)_{\kappa} C_{\kappa}(\mathbf{X})}{(b)_{\kappa} k!}, \quad (\text{Confluent hypergeometric function of the first kind}) \\ {}_2F_1(a, b; c; \mathbf{X}) &= \sum_{k=0}^{\infty} \sum_{\kappa \vdash k} \frac{(a)_{\kappa} (b)_{\kappa} C_{\kappa}(\mathbf{X})}{(c)_{\kappa} k!}, \quad \|\mathbf{X}\| < 1, \quad (\text{Gauss hypergeometric function}) \end{aligned} \tag{4}$$

# Hypergeometric functions of matrix arguments

## Theorem ([2])

Let  $\mathbf{Z}$  be a complex symmetric matrix whose real part is positive-definite, and let  $\mathbf{X}$  be an arbitrary complex symmetric matrix. Then

$$\int_{\mathbf{R}>0} \text{etr}(-\mathbf{Z}\mathbf{R}) |\mathbf{R}|^{\alpha - \frac{d+1}{2}} C_{\kappa}(\mathbf{R}\mathbf{X}) d\mathbf{R} = \Gamma_d(\alpha) (\alpha)_{\kappa} |\mathbf{Z}|^{-\alpha} C_{\kappa}(\mathbf{X}\mathbf{Z}^{-1}), \quad (5)$$

the integration being over the space of positive-definite  $d \times d$  matrices, and valid for all complex numbers  $\alpha$  satisfying  $\text{Re}(\alpha) > \frac{d-1}{2}$ .  $\Gamma_d(\alpha)$  is the multivariate gamma function defined in (??).

## Theorem ([5])

If  $\mathbf{Z} > 0$  and  $\text{Re}(\alpha) > \frac{d-1}{2}$ , and  $\mathbf{X}$  is a  $d \times d$  symmetric matrix, we have

$$\begin{aligned} & \int_{\mathbf{R}>0} \text{etr}(-\mathbf{Z}\mathbf{R}) |\mathbf{R}|^{\alpha - \frac{d+1}{2}} {}_pF_q(a_1, \dots, a_p; b_1, \dots, b_q; \mathbf{R}\mathbf{X}) d\mathbf{R} \\ &= \int_{\mathbf{R}>0} \text{etr}(-\mathbf{Z}\mathbf{R}) |\mathbf{R}|^{\alpha - \frac{d+1}{2}} {}_pF_q(a_1, \dots, a_p; b_1, \dots, b_q; \mathbf{R}^{1/2} \mathbf{X} \mathbf{R}^{1/2}) d\mathbf{R} \\ &= \Gamma_d(\alpha) |\mathbf{Z}|^{-\alpha} {}_{p+1}F_q(a_1, \dots, a_p, \alpha; b_1, \dots, b_q; \mathbf{X}\mathbf{Z}^{-1}). \end{aligned}$$

# Joint prior for two precision matrices

- We define the following joint prior for the precision matrices of the source and target domains:

$$\begin{aligned} p(\Lambda'_t, \Lambda'_s) &= K^l \text{etr} \left( -\frac{1}{2} \left( (\mathbf{M}'_t)^{-1} + \mathbf{F}' \mathbf{C}' \mathbf{F}' \right) \sim'_t \right) \\ &\text{etr} \left( -\frac{1}{2} (\mathbf{C}')^{-1} \Lambda'_s \right) |\Lambda'_t|^{\frac{\nu' - d - 1}{2}} |\Lambda'_s|^{\frac{\nu' - d - 1}{2}} {}_0F_1 \left( \frac{\nu'}{2}; \frac{1}{4} \mathbf{G}' \right), \end{aligned} \quad (6)$$

where  $\mathbf{M}' = \begin{pmatrix} \mathbf{M}'_t & \mathbf{M}'_{ts} \\ \mathbf{M}'_{ts} & \mathbf{M}'_s \end{pmatrix}$  is an  $2d \times 2d$  positive definite scale matrix, and

$\nu' \geq 2d$  is degrees of freedom.  $\mathbf{C}' = \mathbf{M}'_s - \mathbf{M}'_{ts} \mathbf{M}'_t^{-1} \mathbf{M}'_{ts}$ ,  
 $\mathbf{F}' = (\mathbf{C}')^{-1} \mathbf{M}'_{ts} \mathbf{M}'_t^{-1}$ ,  $\mathbf{G}' = \Lambda'_s \mathbf{F}' \Lambda'_t \mathbf{F}' \Lambda'_s$ ,  $(K^l)^{-1} = 2^{d\nu'} \Gamma_d^2 \left( \frac{\nu'}{2} \right) |\mathbf{M}'|^{\frac{\nu'}{2}}$ .

- The marginal distributions are Wishart for each domain (we are interested in understanding **how source data may help better learn the marginal distribution in target domain**):

$$\Lambda'_z \sim W_d(\mathbf{M}'_z, \nu'), \quad l \in \{1, \dots, L\}, \quad z \in \{\mathbf{s}, \mathbf{t}\}. \quad (7)$$

# Posteriors

- Joint likelihood of source and target:

$$\begin{aligned}\rho(\mathcal{D}_t, \mathcal{D}_s | \mu_t, \mu_s, \Lambda_t, \Lambda_s) &= \rho(\mathcal{D}_t | \mu_t, \Lambda_t) \rho(\mathcal{D}_s | \mu_s, \Lambda_s) \\ &= \rho(\mathcal{D}_t^1, \dots, \mathcal{D}_t^L | \mu_t^1, \dots, \mu_t^L, \Lambda_t^1, \dots, \Lambda_t^L) \\ &\quad \times \rho(\mathcal{D}_s^1, \dots, \mathcal{D}_s^L | \mu_s^1, \dots, \mu_s^L, \Lambda_s^1, \dots, \Lambda_s^L) \\ &= \prod_{l=1}^L \rho(\mathcal{D}_t^l | \mu_t^l, \Lambda_t^l) \prod_{l=1}^L \rho(\mathcal{D}_s^l | \mu_s^l, \Lambda_s^l).\end{aligned}\tag{8}$$

- Joint posterior of source and target:

$$\begin{aligned}\rho(\mu_t, \mu_s, \Lambda_t, \Lambda_s | \mathcal{D}_t, \mathcal{D}_s) \\ &\propto \rho(\mathcal{D}_t, \mathcal{D}_s | \mu_t, \mu_s, \Lambda_t, \Lambda_s) \rho(\mu_t, \mu_s, \Lambda_t, \Lambda_s) \\ &\propto \prod_{l=1}^L \rho(\mathcal{D}_t^l | \mu_t^l, \Lambda_t^l) \prod_{l=1}^L \rho(\mathcal{D}_s^l | \mu_s^l, \Lambda_s^l) \prod_{l=1}^L \rho(\mu_t^l, \mu_s^l, \Lambda_t^l, \Lambda_s^l) \\ &\propto \prod_{l=1}^L \rho(\mathcal{D}_t^l | \mu_t^l, \Lambda_t^l) \rho(\mathcal{D}_s^l | \mu_s^l, \Lambda_s^l) \rho(\mu_s^l | \Lambda_s^l) \rho(\mu_t^l | \Lambda_t^l) \rho(\Lambda_s^l, \Lambda_t^l)\end{aligned}\tag{9}$$

# Posteriors of Target Parameters

- Posterior of target given both the source and target data:

$$\begin{aligned}\rho(\mu_t, \Lambda_t | \mathcal{D}_t, \mathcal{D}_s) &= \int_{\mu_s, \Lambda_s} \rho(\mu_t, \mu_s, \Lambda_t, \Lambda_s | \mathcal{D}_t, \mathcal{D}_s) d\mu_s d\Lambda_s \\ &= \prod_{l=1}^L \int_{\mu_s^l, \Lambda_s^l} \rho(\mu_t^l, \mu_s^l, \Lambda_t^l, \Lambda_s^l | \mathcal{D}_t^l, \mathcal{D}_s^l) d\mu_s^l d\Lambda_s^l \\ &= \prod_{l=1}^L \rho(\mu_t^l, \Lambda_t^l | \mathcal{D}_t^l, \mathcal{D}_s^l),\end{aligned}$$

where

$$\begin{aligned}\rho(\mu_t^l, \Lambda_t^l | \mathcal{D}_t^l, \mathcal{D}_s^l) &= \int_{\mu_s^l, \Lambda_s^l} \rho(\mu_t^l, \mu_s^l, \Lambda_t^l, \Lambda_s^l | \mathcal{D}_t^l, \mathcal{D}_s^l) d\mu_s^l d\Lambda_s^l \\ &\propto \rho(\mathcal{D}_t^l | \mu_t^l, \Lambda_t^l) \rho(\mu_t^l | \Lambda_t^l) \\ &\quad \times \int_{\mu_s^l, \Lambda_s^l} \rho(\mathcal{D}_s^l | \mu_s^l, \Lambda_s^l) \rho(\mu_s^l | \Lambda_s^l) \rho(\Lambda_s^l, \Lambda_t^l) d\mu_s^l d\Lambda_s^l.\end{aligned}\tag{10}$$

# Posteriors of Target Parameters

## Lemma

If  $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  where  $\mathbf{x}_i$  is a  $d \times 1$  vector and  $\mathbf{x}_i \sim \mathcal{N}(\mu, (\Lambda)^{-1})$ , for  $i = 1, \dots, n$ , and  $(\mu, \Lambda)$  has a Normal-Wishart prior, such that,  $\mu | \Lambda \sim \mathcal{N}(\mathbf{m}, (\kappa \Lambda)^{-1})$  and  $\Lambda \sim W_d(\mathbf{M}, \nu)$ , then the posterior of  $(\mu, \Lambda)$  upon observing  $\mathcal{D}$  is also a Normal-Wishart distribution:

$$\begin{aligned}\mu | \Lambda, \mathcal{D} &\sim \mathcal{N}(\mathbf{m}_n, (\kappa_n \Lambda)^{-1}), \\ \Lambda | \mathcal{D} &\sim W_d(\mathbf{M}_n, \nu_n),\end{aligned}\tag{11}$$

where

$$\begin{aligned}\kappa_n &= \kappa + n, \quad \nu_n = \nu + n, \quad \mathbf{m}_n = \frac{\kappa \mathbf{m} + n \bar{\mathbf{x}}}{\kappa + n}, \\ \mathbf{M}_n^{-1} &= \mathbf{M}^{-1} + \mathbf{S} + \frac{\kappa n}{\kappa + n} (\mathbf{m} - \bar{\mathbf{x}})(\mathbf{m} - \bar{\mathbf{x}})'\end{aligned}\tag{12}$$

depending on the sample mean and covariance matrix

$$\bar{\mathbf{x}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i, \quad \mathbf{S} = \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})'\tag{13}$$

# Posteriors of Target Parameters

- Using the previous lemma and theorems, we can find the closed-form posterior distribution of mean and precision matrix of the target domain, which is a function of matrix-variate Confluent hypergeometric function of first kind:

$$\begin{aligned} p(\mu_t^l, \Lambda_t^l | \mathcal{D}_t^l, \mathcal{D}_s^l) = & \\ A^l |\Lambda_t^l|^{\frac{1}{2}} \exp\left(-\frac{\kappa_{t,n}^l}{2} (\mu_t^l - \mathbf{m}_{t,n}^l)' \Lambda_t^l (\mu_t^l - \mathbf{m}_{t,n}^l)\right) & \\ \times |\Lambda_t^l|^{\frac{\nu^l + n_t^l - d - 1}{2}} \text{etr}\left(-\frac{1}{2} (\mathbf{T}_t^l)^{-1} \Lambda_t^l\right) & \\ \times {}_1F_1\left(\frac{\nu^l + n_s^l}{2}; \frac{\nu^l}{2}; \frac{1}{2} \mathbf{F}^l \Lambda_t^l \mathbf{F}^{l'} \mathbf{T}_s^l\right), & \end{aligned} \quad (14)$$

- We see that as opposed to one-domain posterior which is Normal-Wishart, here the posterior is Normal-Hypergeometric.

where

$$(A')^{-1} = \left( \frac{2\pi}{\kappa'_{t,n}} \right)^{\frac{d}{2}} 2^{-\frac{d(\nu^l + n'_t)}{2}} \Gamma_d \left( \frac{\nu^l + n'_t}{2} \right) \times \left| \mathbf{T}'_t \right|^{\frac{\nu^l + n'_t}{2}} {}_2F_1 \left( \frac{\nu^l + n'_t}{2}, \frac{\nu^l + n'_t}{2}; \frac{\nu^l}{2}; \mathbf{T}'_s \mathbf{F}' \mathbf{T}'_t \mathbf{F}' \right), \quad (15)$$

and

$$\begin{aligned} \kappa'_{t,n} &= \kappa'_t + n'_t, & \kappa'_{s,n} &= \kappa'_s + n'_s, \\ \mathbf{m}'_{t,n} &= \frac{\kappa'_t \mathbf{m}'_t + n'_t \bar{\mathbf{x}}'_t}{\kappa'_t + n'_t}, & \mathbf{m}'_{s,n} &= \frac{\kappa'_s \mathbf{m}'_s + n'_s \bar{\mathbf{x}}'_s}{\kappa'_s + n'_s}, \\ (\mathbf{T}'_t)^{-1} &= (\mathbf{M}'_t)^{-1} + \mathbf{F}' \mathbf{C}' \mathbf{F}' + \mathbf{S}'_t + \frac{\kappa'_t n'_t}{\kappa'_t + n'_t} (\mathbf{m}'_t - \bar{\mathbf{x}}'_t)(\mathbf{m}'_t - \bar{\mathbf{x}}'_t)', & (16) \\ (\mathbf{T}'_s)^{-1} &= (\mathbf{C}'_s)^{-1} + \mathbf{S}'_s + \frac{\kappa'_s n'_s}{\kappa'_s + n'_s} (\mathbf{m}'_s - \bar{\mathbf{x}}'_s)(\mathbf{m}'_s - \bar{\mathbf{x}}'_s)', \end{aligned}$$

depending on the corresponding sample mean vectors and sample covariance matrices as follows:

$$\begin{aligned} \bar{\mathbf{x}}'_t &= \frac{1}{n'_t} \sum_{i=1}^{n'_t} \mathbf{x}'_{t,i}, & \bar{\mathbf{x}}'_s &= \frac{1}{n'_s} \sum_{i=1}^{n'_s} \mathbf{x}'_{s,i}, \\ \mathbf{S}'_t &= \sum_{i=1}^{n'_t} (\mathbf{x}'_{t,i} - \bar{\mathbf{x}}'_t) (\mathbf{x}'_{t,i} - \bar{\mathbf{x}}'_t)', & (17) \\ \mathbf{S}'_s &= \sum_{i=1}^{n'_s} (\mathbf{x}'_{s,i} - \bar{\mathbf{x}}'_s) (\mathbf{x}'_{s,i} - \bar{\mathbf{x}}'_s)'. \end{aligned}$$

# Effective Class-Conditional Densities

- The effective class-conditional densities (thereafter posterior predictive):

$$p(\mathbf{x}|l) = \int_{\mu_t^l, \Lambda_t^l} p(\mathbf{x}|\mu_t^l, \Lambda_t^l) \pi^*(\mu_t^l, \Lambda_t^l) d\mu_t^l d\Lambda_t^l \quad (18)$$

$$\begin{aligned} O_{\text{BTTL}}(\mathbf{x}|l) &= p(\mathbf{x}|l) = \\ &\pi^{-\frac{d}{2}} \left( \frac{\kappa_{t,n}^l}{\kappa_{\mathbf{x}}^l} \right)^{\frac{d}{2}} \Gamma_d \left( \frac{\nu^l + n_t^l + 1}{2} \right) \Gamma_d^{-1} \left( \frac{\nu^l + n_t^l}{2} \right) \\ &\times |\mathbf{T}_{\mathbf{x}}^l|^{\frac{\nu^l + n_t^l + 1}{2}} {}_2F_1 \left( \frac{\nu^l + n_s^l}{2}, \frac{\nu^l + n_t^l + 1}{2}; \frac{\nu^l}{2}; \mathbf{T}_s^l \mathbf{F}^l \mathbf{T}_{\mathbf{x}}^l \mathbf{F}^{l'} \right) \\ &\times |\mathbf{T}_t^l|^{-\frac{\nu^l + n_t^l}{2}} {}_2F_1^{-1} \left( \frac{\nu^l + n_s^l}{2}, \frac{\nu^l + n_t^l}{2}; \frac{\nu^l}{2}; \mathbf{T}_s^l \mathbf{F}^l \mathbf{T}_t^l \mathbf{F}^{l'} \right), \end{aligned} \quad (19)$$

where

$$\begin{aligned} \kappa_{\mathbf{x}}^l &= \kappa_{t,n}^l + 1 = \kappa_t^l + n_t^l + 1, \\ \mathbf{m}_{\mathbf{x}}^l &= \frac{\kappa_{t,n}^l \mathbf{m}_{t,n}^l + \mathbf{x}}{\kappa_{t,n}^l + 1}, \\ (\mathbf{T}_{\mathbf{x}}^l)^{-1} &= (\mathbf{T}_t^l)^{-1} + \frac{\kappa_{t,n}^l}{\kappa_{t,n}^l + 1} (\mathbf{m}_{t,n}^l - \mathbf{x}) (\mathbf{m}_{t,n}^l - \mathbf{x})'. \end{aligned} \quad (20)$$

# Optimal Bayesian Transfer Learning (OBTL) Classifier

- Let  $c_t^l$  be the prior probability that the target sample  $\mathbf{x}$  belongs to the class  $l \in \{1, \dots, L\}$ . Since  $0 < c_t^l < 1$  and  $\sum_{l=1}^L c_t^l = 1$ , a Dirichlet prior is assumed for the  $c_t^l$ :

$$(c_t^1, \dots, c_t^L) \sim \text{Dir}(L, \xi_t), \quad (21)$$

where  $\xi_t = (\xi_t^1, \dots, \xi_t^L)$  are the concentration parameters, where  $\xi_t^l > 0$  for all  $l \in \{1, \dots, L\}$ .

- The posterior of  $c_t^l$ 's is also another Dirichlet distribution:

$$\begin{aligned} \pi^* &= (c_t^1, \dots, c_t^L | \mathbf{n}) \sim \text{Dir}(L, \xi_t + \mathbf{n}) \\ &= \text{Dir}(L, \xi_t^1 + n_t^1, \dots, \xi_t^L + n_t^L), \end{aligned} \quad (22)$$

with the posterior mean of  $c_t^l$  as

$$E_{\pi^*}(c_t^l) = \frac{\xi_t^l + n_t^l}{N_t + \xi_t^0}, \quad (23)$$

where  $N_t = \sum_{l=1}^L n_t^l$  and  $\xi_t^0 = \sum_{l=1}^L \xi_t^l$ .

- The optimal Bayesian transfer learning (OBTL) classifier for any new unlabeled sample  $\mathbf{x}$  in the target domain is defined as:

## OBTL

$$\Psi_{\text{OBTL}}(\mathbf{x}) = \arg \max_{l \in \{1, \dots, L\}} E_{\pi^*}(c_t^l) O_{\text{OBTL}}(\mathbf{x}|l). \quad (24)$$

# OBC in Target Domain

- The effective class-conditional densities  $p(\mathbf{x}|l) = O_{\text{OBC}}(\mathbf{x}|l)$  for OBC are derived as:

$$O_{\text{OBC}}(\mathbf{x}|l) = \pi^{-\frac{d}{2}} \left( \frac{\kappa_{t,n}^l}{\kappa_{t,n}^l + 1} \right)^{\frac{d}{2}} \Gamma_d \left( \frac{\nu^l + n_t^l + 1}{2} \right) \Gamma_d^{-1} \left( \frac{\nu^l + n_t^l}{2} \right) |\mathbf{M}_{\mathbf{x}}^l|^{\frac{\nu^l + n_t^l + 1}{2}} |\mathbf{M}_{t,n}^l|^{-\frac{\nu^l + n_t^l}{2}}, \quad (25)$$

where

$$\left( \mathbf{M}_{\mathbf{x}}^l \right)^{-1} = \left( \mathbf{M}_{t,n}^l \right)^{-1} + \frac{\kappa_{t,n}^l}{\kappa_{t,n}^l + 1} (\mathbf{m}_{t,n}^l - \mathbf{x})(\mathbf{m}_{t,n}^l - \mathbf{x})', \quad (26)$$

$$\kappa_{t,n}^l = \kappa_t^l + n_t^l, \quad \nu_{t,n}^l = \nu^l + n_t^l, \quad \mathbf{m}_{t,n}^l = \frac{\kappa_t^l \mathbf{m}_t^l + n_t^l \bar{\mathbf{x}}_t^l}{\kappa_t^l + n_t^l}, \quad (27)$$

$$\left( \mathbf{M}_{t,n}^l \right)^{-1} = \left( \mathbf{M}_t^l \right)^{-1} + \mathbf{S}_t^l + \frac{\kappa_t^l n_t^l}{\kappa_t^l + n_t^l} (\mathbf{m}_t^l - \bar{\mathbf{x}}_t^l)(\mathbf{m}_t^l - \bar{\mathbf{x}}_t^l)',$$

with the corresponding sample mean and covariance:  $\bar{\mathbf{x}}_t^l = \frac{1}{n_t^l} \sum_{i=1}^{n_t^l} \mathbf{x}_{t,i}^l$ ,  $\mathbf{S}_t^l = \sum_{i=1}^{n_t^l} (\mathbf{x}_{t,i}^l - \bar{\mathbf{x}}_t^l)(\mathbf{x}_{t,i}^l - \bar{\mathbf{x}}_t^l)'$ .

- The OBC is defined as:  $\Psi_{\text{OBC}}(\mathbf{x}) = \operatorname{argmax}_{l \in \{1, \dots, L\}} E_{\pi^*} (c_l^l) O_{\text{OBC}}(\mathbf{x}|l)$ .

## Theorem

If  $\mathbf{M}_{t,n}^l = \mathbf{0}$  for all  $l \in \{1, \dots, L\}$ , then

$$\Psi_{\text{OBTL}}(\mathbf{x}) = \Psi_{\text{OBC}}(\mathbf{x}), \quad (28)$$

meaning that if there is no interaction between the source and target domains in all the classes a priori, then the OBTL classifier turns to the OBC classifier in the target domain.

# Laplace Approximation of Gauss Hypergeometric

- The Gauss hypergeometric function has the following integral representation:

$${}_2F_1(a, b; c; \mathbf{X}) = B_d^{-1}(a, c - a) \times \int_{0_d < \mathbf{Y} < \mathbf{I}_d} |\mathbf{Y}|^{a - \frac{d+1}{2}} |\mathbf{I}_d - \mathbf{Y}|^{c-a - \frac{d+1}{2}} |\mathbf{I}_d - \mathbf{X}\mathbf{Y}|^{-b} d\mathbf{Y}, \quad (29)$$

which is valid under the following conditions:  $\mathbf{X} \in \mathbf{C}^{d \times d}$  is symmetric and satisfies  $\text{Re}(\mathbf{X}) < \mathbf{I}_d$ ,  $\text{Re}(a) > \frac{d-1}{2}$ , and  $\text{Re}(c - a) > \frac{d-1}{2}$ .  $B_d(\alpha, \beta)$  is the multivariate beta function

$$B_d(\alpha, \beta) = \frac{\Gamma_d(\alpha)\Gamma_d(\beta)}{\Gamma_d(\alpha + \beta)}. \quad (30)$$

- The Laplace approximation is one common solution to approximate the integral

$$I = \int_{\mathbf{y} \in D} h(\mathbf{y}) \exp(-\lambda g(\mathbf{y})) d\mathbf{y}, \quad (31)$$

where  $D \subseteq \mathbf{R}^d$  is an open set and  $\lambda$  is a real parameter. If  $g(\lambda)$  has a unique minimum over  $D$  at point  $\hat{\mathbf{y}} \in D$ , then the Laplace approximation to  $I$  is given by

$$\tilde{I} = (2\pi)^{\frac{d}{2}} \lambda^{-\frac{d}{2}} |g''(\hat{\mathbf{y}})|^{-\frac{1}{2}} h(\hat{\mathbf{y}}) \exp(-\lambda g(\hat{\mathbf{y}})), \quad (32)$$

where  $g''(y) = \frac{\partial^2 g(y)}{\partial y \partial y^T}$  is the Hessian of  $g(y)$ .

- The calibrated Laplace approximation of Gauss hypergeometric functions of matrix argument:

$$\begin{aligned}
 {}_2\hat{F}_1(a, b; c; \mathbf{X}) &= \frac{{}_2\tilde{F}_1(a, b; c; \mathbf{X})}{{}_2\tilde{F}_1(a, b; c; \mathbf{0})} = c^{cd - \frac{d(d+1)}{4}} R_{2,1}^{-\frac{1}{2}} \\
 &\times \prod_{i=1}^d \left\{ \left( \frac{\hat{y}_i}{a} \right)^a \left( \frac{1 - \hat{y}_i}{c - a} \right)^{c-a} (1 - x_i \hat{y}_i)^{-b} \right\},
 \end{aligned} \tag{33}$$

where

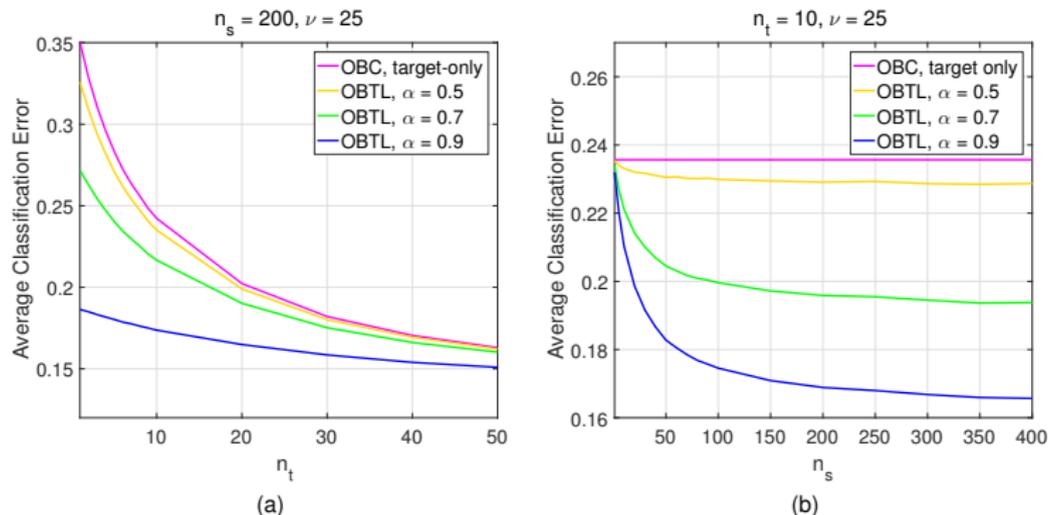
$$R_{2,1} = \prod_{i=1}^d \prod_{j=1}^d \left\{ \frac{\hat{y}_i \hat{y}_j}{a} + \frac{(1 - \hat{y}_i)(1 - \hat{y}_j)}{c - a} - \frac{bx_i x_j \hat{y}_i \hat{y}_j (1 - \hat{y}_i)(1 - \hat{y}_j)}{(1 - x_i \hat{y}_i)(1 - x_j \hat{y}_j) a (c - a)} \right\}, \tag{34}$$

where  $\hat{y}_i$  is defined as

$$\hat{y}_i = \frac{2a}{\sqrt{\tau^2 - 4ax_i(c - b) - \tau}}, \tag{35}$$

where  $\tau = x_i(b - a) - c$  and  $\mathbf{X} = \text{diag}\{x_1, \dots, x_d\}$ .

# Experiment results: synthetic data



**Figure:** (a) Average classification error versus the number of target training data per class,  $n_t$ . The dimension is  $d = 10$ , number of source training data per class is  $n_s = 200$ , and there are  $L = 2$  classes in each domain, (b) Average classification error versus the number of source training data per class,  $n_s$ . The dimension is  $d = 10$ , number of target training data per class is  $n_t = 10$ , and there are  $L = 2$  classes in each domain.

# Experiment results: synthetic data

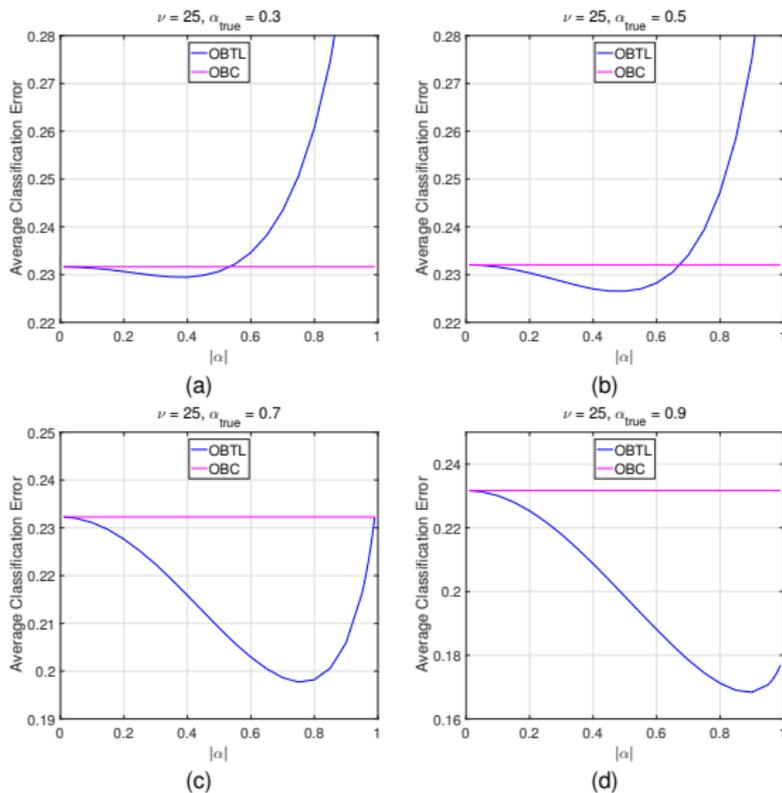


Figure: Average classification error vs  $|\alpha|$

# Experiment results: benchmark image datasets

- Office and Caltech dataset
- Labeled images from four domains: Amazon website, DSLR camera, Webcam, and Caltech dataset
- Labels are office stuff like laptop, backpack, calculator, ...



# Experiment results: benchmark image datasets

**Table:** Accuracy for different source and target domains in *Office+Clatech256* dataset. Domain names are denoted as a: *amazon*, w: *webcam*, d: *dslr*, c: *Caltech256*. Red shows the best accuracy and blue shows the second best accuracy in each column. The results of the first six methods has been adopted from [1]. Similar to [1], we also used the simulation setup of [6] for the OBTL's results.

	a $\rightarrow$ w	a $\rightarrow$ d	a $\rightarrow$ c	w $\rightarrow$ a	w $\rightarrow$ d	w $\rightarrow$ c	d $\rightarrow$ a	d $\rightarrow$ w	d $\rightarrow$ c	c $\rightarrow$ a	c $\rightarrow$ w	c $\rightarrow$ d	Mean
1-NN-t	34.5	33.6	19.7	29.5	35.9	18.9	27.1	33.4	18.6	29.2	33.5	34.1	29.0
SVM-t	63.7	57.2	32.2	46.0	56.5	29.7	45.3	62.1	32.0	45.1	60.2	56.3	48.9
HFA [7]	57.4	55.1	31.0	<b>56.5</b>	56.5	29.0	42.9	60.5	30.9	43.8	58.1	55.6	48.1
MMDT [6]	64.6	56.7	36.4	47.7	67.0	32.2	46.9	74.1	34.1	49.4	63.8	56.5	52.5
CDLS [8]	<b>68.7</b>	<b>60.4</b>	35.3	51.8	60.7	33.5	50.7	68.5	34.9	50.9	<b>66.3</b>	<b>59.8</b>	53.5
ILS (1-NN) [1]	59.7	49.8	<b>43.6</b>	54.3	<b>70.8</b>	<b>38.6</b>	<b>55.0</b>	<b>80.1</b>	<b>41.0</b>	<b>55.1</b>	62.9	56.2	<b>55.6</b>
<b>OBTL</b>	<b>72.1</b>	<b>60.5</b>	<b>42.4</b>	<b>54.7</b>	<b>76.5</b>	<b>37.7</b>	<b>53.9</b>	<b>84.8</b>	<b>40.2</b>	<b>54.8</b>	<b>70.6</b>	<b>61.2</b>	<b>59.1</b>

# OBTL for Count Data

- We use the Negative Binomial model for the feature-label distribution in each domain:

$$\mathbf{x}_{z,i,j}^l \sim \text{NB}(\mu_{z,i}^l, r_{z,i}^l), \quad (36)$$

with the probability mass function (PMF)

$$P(\mathbf{x}_{z,i,j}^l = k | \mu_{z,i}^l, r_{z,i}^l) = \frac{\Gamma(k + r_{z,i}^l)}{\Gamma(r_{z,i}^l)\Gamma(k + 1)} \left( \frac{\mu_{z,i}^l}{\mu_{z,i}^l + r_{z,i}^l} \right)^k \left( \frac{r_{z,i}^l}{\mu_{z,i}^l + r_{z,i}^l} \right)^{r_{z,i}^l}, \quad (37)$$

where  $z \in \{s, t\}$  denotes the source,  $s$ , or target,  $t$ , domains;  $\mu_{z,i}^l$  and  $r_{z,i}^l$  are respectively the mean and shape of the gene  $i$  in domain  $z$  and class  $l$ . The shape parameter is the inverse of the dispersion parameter in Negative Binomial model, which controls the amount of variance. The mean and variance of  $\mathbf{x}_{z,i,j}^l$  are

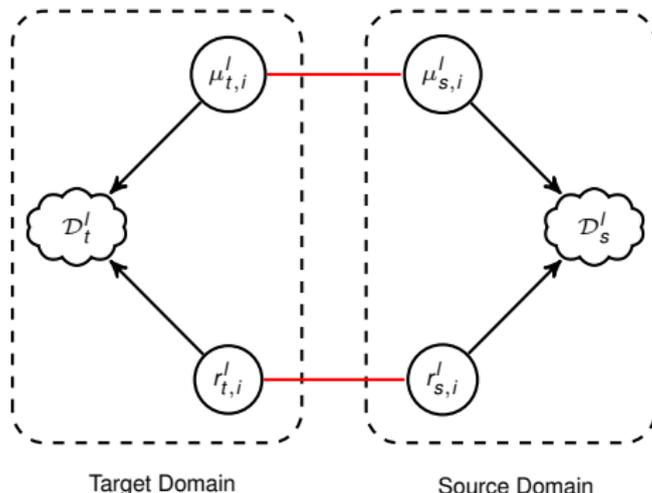
$$\begin{aligned} E(\mathbf{x}_{z,i,j}^l) &= \mu_{z,i}^l, \\ \text{Var}(\mathbf{x}_{z,i,j}^l) &= \mu_{z,i}^l + \frac{(\mu_{z,i}^l)^2}{r_{z,i}^l}. \end{aligned} \quad (38)$$

# Priors and Posteriors

- Let  $\mu = \left\{ \mu_{\{s,t\},\{1:d\}}^{\{1:L\}} \right\}$  and  $r = \left\{ r_{\{s,t\},\{1:d\}}^{\{1:L\}} \right\}$  denote respectively all the mean and shape parameters of the  $d$  genes in  $L$  classes and two domains  $s$  and  $t$ . The prior is factorized as

$$p(\mu, r) = \prod_{l=1}^L \prod_{i=1}^d p\left(\mu_{s,i}^l, \mu_{t,i}^l\right) p\left(r_{s,i}^l, r_{t,i}^l\right). \quad (39)$$

- No closed-form posteriors in this model.
- Hamiltonian Monte Carlo (HMC) method is used for posterior sampling, which outperforms other MCMC methods in that it eliminates all the tuning steps.



# Joint Prior

## Lemma

If  $\Lambda \sim W_2(\mathbf{M}, \nu)$ ,  $\Lambda = \begin{pmatrix} \lambda_{11} & \lambda_{12} \\ \lambda_{12} & \lambda_{22} \end{pmatrix}$ , and  $\mathbf{M} = \begin{pmatrix} m_{11} & m_{12} \\ m_{12} & m_{22} \end{pmatrix}$ , then  $\lambda_{ij} \sim m_{ij} \chi_\nu^2$  for  $i = 1, 2$ , where  $\chi_\nu^2$  denotes the Chi-squared distribution with  $\nu$  degrees of freedom. As a result, their mean and variance are  $E(\lambda_{ij}) = \nu m_{ij}$  and  $\text{Var}(\lambda_{ij}) = 2\nu m_{ij}^2$  for  $i = 1, 2$ . The covariance and correlation between  $\lambda_{11}$  and  $\lambda_{22}$  are respectively

$$\text{Cov}(\lambda_{11}, \lambda_{22}) = 2\nu m_{12}^2, \quad \rho_\lambda = \frac{m_{12}^2}{m_{11} m_{22}}. \quad (40)$$

## Theorem ([3])

Let  $\Lambda = \begin{pmatrix} \lambda_{11} & \lambda_{12} \\ \lambda_{12} & \lambda_{22} \end{pmatrix}$  be a  $2 \times 2$  Wishart random matrix with  $\nu \geq 2$  degrees of freedom and positive-definite scale matrix  $\mathbf{M} = \begin{pmatrix} m_{11} & m_{12} \\ m_{12} & m_{22} \end{pmatrix}$ . The joint distribution of the two diagonal entries  $\lambda_{11}$  and  $\lambda_{22}$  have the density function given by

$$\begin{aligned} p(\lambda_{11}, \lambda_{22}) &= K \exp\left(-\frac{1}{2} (m_{11}^{-1} + c_2 f^2) \lambda_{11}\right) \exp\left(-\frac{1}{2} c_2^{-1} \lambda_{22}\right) \\ &\times (\lambda_{11})^{\frac{\nu}{2}-1} (\lambda_{22})^{\frac{\nu}{2}-1} {}_0F_1\left(\frac{\nu}{2}; \frac{1}{4}g\right), \end{aligned} \quad (41)$$

where  $c_2 = m_{22} - m_{12}^2 m_{11}^{-1}$ ,  $f = c_2^{-1} m_{12} m_{11}^{-1}$ ,  $g = f^2 \lambda_{11} \lambda_{22}$ ,  $K^{-1} = 2^\nu \Gamma^2\left(\frac{\nu}{2}\right) |\mathbf{M}|^{\frac{\nu}{2}}$ , and  ${}_0F_1$  is the generalized hypergeometric function.

# Joint Prior

Here  ${}_0F_1(b; x) = \sum_{k=0}^{\infty} \frac{x^k}{(b)_k k!}$  is called confluent hypergeometric limit function, which is closely related to the Bessel functions:

$$J_{\alpha}(x) = \frac{\left(\frac{x}{2}\right)^{\alpha}}{\Gamma(\alpha+1)} {}_0F_1\left(\alpha+1; -\frac{1}{4}x^2\right). \quad (42)$$

Now, we can define the joint priors of both mean and shape parameters in terms of correlations between two domains:

$$\begin{aligned} \rho(\mu_{s,i}^l, \mu_{t,i}^l) &= K_{\mu,i}^l \exp\left(-\frac{\mu_{s,i}^l}{2m_{s,i}^l(1-\rho_{\mu,i}^l)}\right) \exp\left(-\frac{\mu_{t,i}^l}{2m_{t,i}^l(1-\rho_{\mu,i}^l)}\right) (\mu_{s,i}^l)^{\frac{\nu_{\mu}}{2}-1} (\mu_{t,i}^l)^{\frac{\nu_{\mu}}{2}-1} \\ &\times {}_0F_1\left(\frac{\nu_{\mu}}{2}; \frac{\rho_{\mu,i}^l}{4m_{s,i}^l m_{t,i}^l (1-\rho_{\mu,i}^l)^2} \mu_{s,i}^l \mu_{t,i}^l\right), \end{aligned} \quad (43)$$

$$\begin{aligned} \rho(r_{s,i}^l, r_{t,i}^l) &= K_{r,i}^l \exp\left(-\frac{r_{s,i}^l}{2s_{s,i}^l(1-\rho_{r,i}^l)}\right) \exp\left(-\frac{r_{t,i}^l}{2s_{t,i}^l(1-\rho_{r,i}^l)}\right) (r_{s,i}^l)^{\frac{\nu_r}{2}-1} (r_{t,i}^l)^{\frac{\nu_r}{2}-1} \\ &\times {}_0F_1\left(\frac{\nu_r}{2}; \frac{\rho_{r,i}^l}{4s_{s,i}^l s_{t,i}^l (1-\rho_{r,i}^l)^2} r_{s,i}^l r_{t,i}^l\right), \end{aligned} \quad (44)$$

# Effective Class-Conditional Densities

- Effective class-conditional density for any new test data in target domain is defined as:

$$\rho(\mathbf{x}|l) = \int_{\mu_t^l, r_t^l} \rho(\mathbf{x}|\mu_t^l, r_t^l) \pi^*(\mu_t^l, r_t^l) d\mu_t^l dr_t^l \quad (45)$$

for  $l \in \{1, \dots, L\}$ , where  $\pi^*(\mu_t^l, r_t^l) = \rho(\mu_t^l, r_t^l | \mathcal{D}_t^l, \mathcal{D}_s^l)$  is the posterior of  $(\mu_t^l, r_t^l)$  upon observation of  $\mathcal{D}_t^l$  and  $\mathcal{D}_s^l$ .

- There is no closed form solution for the effective densities. Posterior samples from HMC sampling are used to approximate these effective densities. Suppose we have  $N$  posterior samples from all of  $d$  genes in  $L$  classes. Then the approximation is given by:

$$\rho(\mathbf{x}|l) = \frac{1}{N} \sum_{j=1}^N \prod_{i=1}^d \rho(\mathbf{x}_i | \bar{\mu}_{t,i,j}^l, \bar{r}_{t,i,j}^l) \quad (46)$$

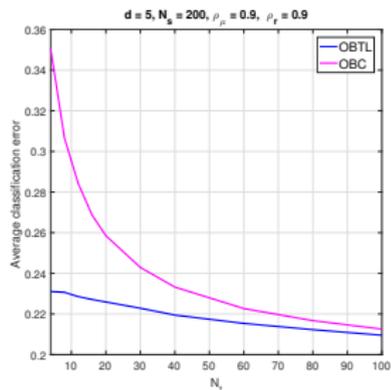
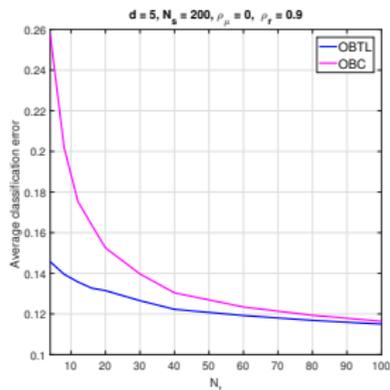
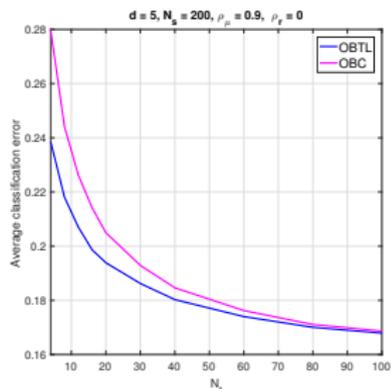
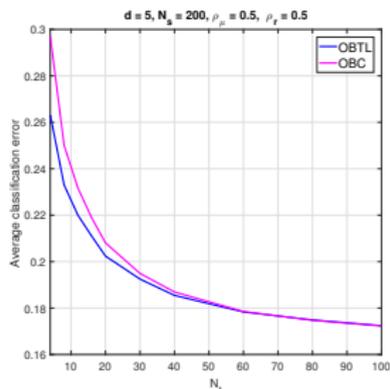
where  $\bar{\mu}_{t,i,j}^l$  and  $\bar{r}_{t,i,j}^l$  are the  $j$ -th posterior sample of gene  $i$  in class  $l$  of target domain for the mean and shape parameters, respectively.

- The OBTL is given by:

## OBTL

$$\Psi_{\text{OBTL}}(\mathbf{x}) = \arg \max_{l \in \{1, \dots, L\}} E_{\pi^*} (c_t^l) \rho(\mathbf{x}|l). \quad (47)$$

# Experiment results: synthetic data



# Experiment results: RNA-seq data

- We classify two kinds of lung cancer: **LUAD** and **LUSC**
- Data are extracted from The Cancer Genome Atlas (TCGA)
- Two RNA-seq measurements: **RNA-seq** and **RNA-seq-v2**. These have different distributions for each genes, so assume two domains:
- Target domain: RNA-seq. LUAD: 125 tumor samples. LUSC: 223 tumor samples
- Source domain: RNA-seq-v2. LUAD: 515 tumor samples. LUSC: 501 tumor samples
- Experiment setup: we randomly generate 50 splits of training (from source and target) and test (only from target) data. We assume  $n'_s = 100$  and  $n'_t = 5$  and number of test data per target class is 100 in each split.
- The average classification error is given for different values of correlations of mean and shape parameters between source and target domains.
- The average error of the OBC is also given for the sake of comparisons.
- Two different sets of features of size  $d = 10$  are picked.

# Experiment results: RNA-seq data

Case 1:

OBC error = 0.1453

OBTL error:

	$\rho_r = 0.5$	$\rho_r = 0.7$	$\rho_r = 0.9$	$\rho_r = 0.99$
$\rho_\mu = 0.5$	0.1187	0.1184	0.1153	0.1136
$\rho_\mu = 0.7$	0.1193	0.1175	0.1149	0.1139
$\rho_\mu = 0.9$	0.1162	0.1141	0.1130	0.1122
$\rho_\mu = 0.99$	0.1167	0.1127	0.1107	0.1111

Case 2:

OBC error = 0.1936

OBTL error:

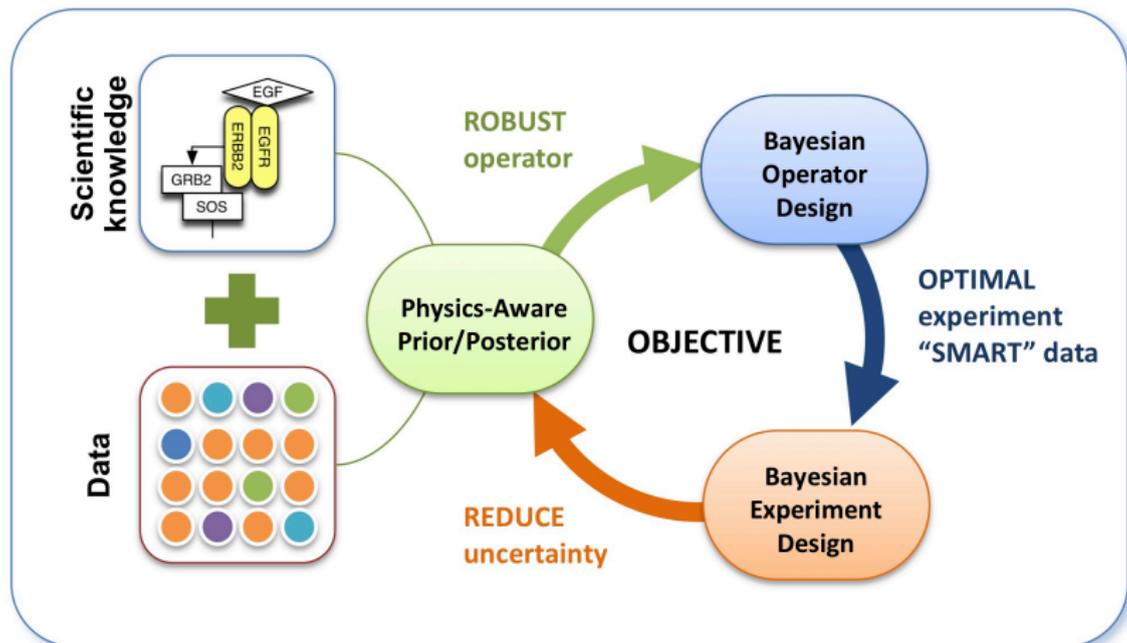
	$\rho_r = 0.5$	$\rho_r = 0.7$	$\rho_r = 0.9$	$\rho_r = 0.99$
$\rho_\mu = 0.5$	0.1654	0.1640	0.1588	0.1543
$\rho_\mu = 0.7$	0.1678	0.1628	0.1571	0.1540
$\rho_\mu = 0.9$	0.1646	0.1619	0.1569	0.1531
$\rho_\mu = 0.99$	0.1631	0.1607	0.1561	0.1513

## Conclusions [9]

- We formulate a Bayesian transfer learning framework to transfer source domain knowledge and data for learning in target domain.
- Our Bayesian framework directly models the **feature-label distributions** in source and target domains.
- The “transferability” across domains can be characterized by a **joint prior distribution** on model parameters of feature-label distributions across domains.
- We derive the Optimal Bayesian Transfer Learning (OBTL) classifier for both continuous and count data with efficient computational solutions.

# Future Research

- Such a Bayesian transfer learning framework enables the closed-loop learning to design experiments for “smart” data and scientific knowledge acquisition.



# References



S. Herath, M. Harandi, and F. Porikli, "Learning an invariant Hilbert space for domain adaptation," in *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017, pp. 3956–3965.



R. J. Muirhead, *Aspects of multivariate statistical theory*. John Wiley & Sons, 2009.



K. Halvorsen, V. Ayala, and E. Fierro, "On the marginal distribution of the diagonal blocks in a blocked Wishart random matrix," *International Journal of Analysis*, vol. 2016, pp. 1–5, 2016.



D. K. Nagar and J. C. Mosquera-Benitez, "Properties of matrix variate hypergeometric function distribution," *Applied Mathematical Sciences*, vol. 11, no. 14, pp. 677–692, 2017.



A. K. Gupta, D. K. Nagar, and L. E. Sánchez, "Properties of matrix variate confluent hypergeometric function distribution," *Journal of Probability and Statistics*, vol. 2016, 2016.



J. Hoffman, E. Rodner, T. Darrell, J. Donahue, and K. Saenko, "Efficient learning of domain-invariant image representations," in *International Conference on Learning Representations (ICLR)*, 2013.



L. Duan, D. Xu, and I. Tsang, "Learning with augmented features for heterogeneous domain adaptation," *ICML*, 2012.



Y.-H. Hubert Tsai, Y.-R. Yeh, and Y.-C. Frank Wang, "Learning cross-domain landmarks for heterogeneous domain adaptation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 5081–5090.



A. Karbalayghareh, X. Qian, and E. R. Dougherty, "Optimal Bayesian transfer learning," *IEEE Transactions on Signal Processing*, vol. 66, pp. 3724–3739, 2018.