

Transfer Learning

Makoto Yamada
myamada@i.kyoto-u.ac.jp

Kyoto University

July/9/2018



Review: Supervised Learning

Problem formulation of supervised learning.

- Input vector: $\mathbf{x} = [x_1, x_2, \dots, x_d]^\top \in \mathbb{R}^d$
- Output: $y \in \mathbb{R}$
- $(\mathbf{x}_i, y_i) \stackrel{\text{i.i.d.}}{\sim} p(\mathbf{x}, y)$
- Labeled data: $\{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n)\}$
- Model: $f(\mathbf{x}; \mathbf{w}) = \mathbf{w}^\top \mathbf{x}$. (Linear model)

Risk: $R(\mathbf{w}) = \iint \text{loss}(y, f(\mathbf{x}; \mathbf{w})) p(\mathbf{x}, y) d\mathbf{x} dy$

Empirical Risk: $R_{\text{emp}}(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \text{loss}(y_i, f(\mathbf{x}_i; \mathbf{w}))$

Empirical Risk Minimization (ERM): $\hat{\mathbf{w}} = \operatorname{argmin}_{\mathbf{w}} R_{\text{emp}}(\mathbf{w})$

Supervised Learning:

- Training $\{(\mathbf{x}_i, y_i)\}_{i=1}^n \stackrel{\text{i.i.d.}}{\sim} p_{\text{tr}}(\mathbf{x}, y)$
- Test $(\mathbf{x}, y) \stackrel{\text{i.i.d.}}{\sim} p_{\text{te}}(\mathbf{x}, y)$ (Not observed during training)
- $p_{\text{tr}} = p_{\text{te}}$ (Training and test distributions are same)

Semi-supervised Learning:

- Training $\{(\mathbf{x}_i, y_i)\}_{i=1}^n \stackrel{\text{i.i.d.}}{\sim} p_{\text{tr}}(\mathbf{x}, y), \{\mathbf{x}_i\}_{i=n+1}^{n+m} \stackrel{\text{i.i.d.}}{\sim} p_{\text{tr}}(\mathbf{x})$.
- Test $(\mathbf{x}, y) \stackrel{\text{i.i.d.}}{\sim} p_{\text{te}}(\mathbf{x}, y)$ (Not observed during training)
- $p_{\text{tr}} = p_{\text{te}}$ (Training and test distributions are same)

If $p_{\text{tr}} \neq p_{\text{te}}$, supervised method and semi-supervised method do not perform well. A possible answer would be **Transfer Learning!**

Types of Transfer Learning

Key idea: Reduce generalization error in **test** data. (not in training data)

Unsupervised transfer learning

- $\{(\mathbf{x}_i^{\text{tr}}, y_i^{\text{tr}})\}_{i=1}^{n_{\text{tr}}} \stackrel{\text{i.i.d.}}{\sim} p_{\text{tr}}(\mathbf{x}, y),$
- $\{\mathbf{x}_j^{\text{te}}\}_{j=1}^{n_{\text{te}}} \stackrel{\text{i.i.d.}}{\sim} p_{\text{te}}(\mathbf{x}), n_{\text{tr}} \ll n_{\text{te}}$

Supervised transfer learning

- $\{(\mathbf{x}_i^{\text{tr}}, y_i^{\text{tr}})\}_{i=1}^{n_{\text{tr}}} \stackrel{\text{i.i.d.}}{\sim} p_{\text{tr}}(\mathbf{x}, y)$
- $\{(\mathbf{x}_j^{\text{te}}, y_j^{\text{te}})\}_{j=1}^{n_{\text{te}}} \stackrel{\text{i.i.d.}}{\sim} p_{\text{te}}(\mathbf{x}, y), n_{\text{te}} \ll n_{\text{tr}}$

Semi-supervised transfer learning

- $\{(\mathbf{x}_i^{\text{tr}}, y_i^{\text{tr}})\}_{i=1}^{n_{\text{tr}}} \stackrel{\text{i.i.d.}}{\sim} p_{\text{tr}}(\mathbf{x}, y)$
- $\{(\mathbf{x}_j^{\text{te}}, y_j^{\text{te}})\}_{j=1}^{n_{\text{te}}} \stackrel{\text{i.i.d.}}{\sim} p_{\text{te}}(\mathbf{x}), n_{\text{te}} \ll n_{\text{tr}}$
- $\{\mathbf{x}_j^{\text{te}}\}_{j=n_{\text{te}}+1}^{n_{\text{te}}+n'_{\text{te}}} \stackrel{\text{i.i.d.}}{\sim} p_{\text{te}}(\mathbf{x}), n_{\text{tr}} \ll n_{\text{te}}$

Unsupervised Transfer Learning

Key idea: We assume

- It does not need to have **test** label
- Need some assumption

Standard approaches

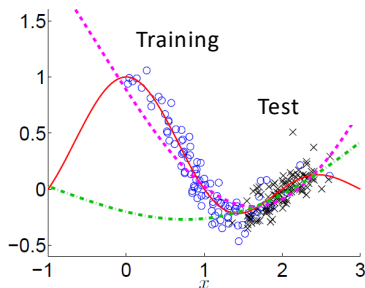
- Importance weighted method (e.g., Covariate shift adaptation)
- Subspace based method.

Unsupervised Transfer Learning: Covariate shift adaptation

Problem setup:

- $\{(\mathbf{x}_i^{\text{tr}}, y_i^{\text{tr}})\}_{i=1}^{n_{\text{tr}}} \stackrel{\text{i.i.d.}}{\sim} p_{\text{tr}}(\mathbf{x}, y),$
- $\{\mathbf{x}_j^{\text{te}}\}_{j=1}^{n_{\text{te}}} \stackrel{\text{i.i.d.}}{\sim} p_{\text{te}}(\mathbf{x}), \quad n_{\text{tr}} \ll n_{\text{te}}$

Key idea: Learning a function so that error in test data is minimized under the assumption $p_{\text{tr}}(y|\mathbf{x}) = p_{\text{te}}(y|\mathbf{x})$



Unsupervised Transfer Learning

The risk can be written as

$$\begin{aligned} J(\mathbf{w}) &= \iint L(y, f(\mathbf{x})) p_{\text{te}}(\mathbf{x}, y) d\mathbf{x} dy \\ &= \iint L(y, f(\mathbf{x})) \frac{p_{\text{te}}(\mathbf{x}, y)}{p_{\text{tr}}(\mathbf{x}, y)} p_{\text{tr}}(\mathbf{x}, y) d\mathbf{x} dy \\ &= \iint L(y, f(\mathbf{x})) \frac{p_{\text{te}}(y|\mathbf{x}) p_{\text{te}}(\mathbf{x})}{p_{\text{tr}}(y|\mathbf{x}) p_{\text{tr}}(\mathbf{x})} p_{\text{tr}}(y, \mathbf{x}) d\mathbf{x} dy \\ &= \iint L(y, f(\mathbf{x})) \frac{p_{\text{te}}(\mathbf{x})}{p_{\text{tr}}(\mathbf{x})} p_{\text{tr}}(y, \mathbf{x}) d\mathbf{x} dy \\ &\approx \frac{1}{n_{\text{tr}}} \sum_{i=1}^{n_{\text{tr}}} L(y_i^{\text{tr}}, f(\mathbf{x}_i^{\text{tr}})) \frac{p_{\text{te}}(\mathbf{x}_i^{\text{tr}})}{p_{\text{tr}}(\mathbf{x}_i^{\text{tr}})} \end{aligned}$$

Actually, it is a **weighted maximum likelihood** problem. Note $\frac{p_{\text{te}}(\mathbf{x}_i^{\text{tr}})}{p_{\text{tr}}(\mathbf{x}_i^{\text{tr}})}$ is a ratio of probability densities (density-ratio)

Unsupervised Transfer Learning: Covariate shift adaptation

Exponentially-flattened Importance weighted empirical risk minimization (IW-ERM):

$$\min_{f \in \mathcal{F}} \frac{1}{n_{\text{tr}}} \sum_{i=1}^{n_{\text{tr}}} L(y_i^{\text{tr}}, f(\mathbf{x}_i^{\text{tr}})) \left(\frac{p_{\text{te}}(\mathbf{x}_i^{\text{tr}})}{p_{\text{tr}}(\mathbf{x}_i^{\text{tr}})} \right)^{\tau}$$

where $0 \leq \tau \leq 1$ is a tuning parameter for stabilizing the covariate shift adaptation.

- $\tau = 0 \rightarrow \text{ERM}$
- $0 < \tau < 1 \rightarrow \text{Intermediate}$
- $\tau = 1 \text{ IW-ERM}$

Setting τ to $0 < \tau < 1$ is practically useful.

Unsupervised Transfer Learning: Covariate shift adaptation

Relative Importance weighted empirical risk minimization (RIW-ERM):

$$\min_{f \in \mathcal{F}} \frac{1}{n_{\text{tr}}} \sum_{i=1}^{n_{\text{tr}}} L(y_i^{\text{tr}}, f(\mathbf{x}_i^{\text{tr}})) \frac{p_{\text{te}}(\mathbf{x}_i^{\text{tr}})}{(1 - \alpha)p_{\text{tr}}(\mathbf{x}_i^{\text{tr}}) + \alpha p_{\text{tr}}(\mathbf{x}_i^{\text{tr}})}$$

where $0 \leq \alpha \leq 1$ is a tuning parameter for stabilizing the covariate shift adaptation.

- $\alpha = 0 \rightarrow \text{ERM}$
- $0 < \alpha < 1 \rightarrow \text{Intermediate}$
- $\alpha = 1 \text{ IW-ERM}$

$$r_{\alpha}(\mathbf{x}) = \frac{p_{\text{te}}(\mathbf{x})}{(1 - \alpha)p_{\text{tr}}(\mathbf{x}) + \alpha p_{\text{tr}}(\mathbf{x})} < \frac{1}{1 - \alpha}$$

The density ratio is bounded above by $1/(1 - \alpha)$. Next lecture, I will introduce how to estimate the density ratio.

Unsupervised Transfer Learning: Importance Weighted Least Squares

The importance weighted least squares problem can be written as

$$\min_{\mathbf{w}} J(\mathbf{w}) = \frac{1}{n_{\text{tr}}} \sum_{i=1}^{n_{\text{tr}}} r(\mathbf{x}_i^{\text{tr}}) \|y_i^{\text{tr}} - \mathbf{w}^\top \mathbf{x}_i^{\text{tr}}\|_2^2,$$

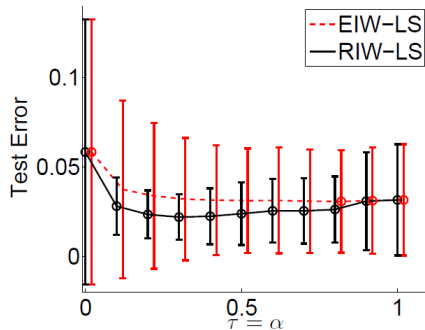
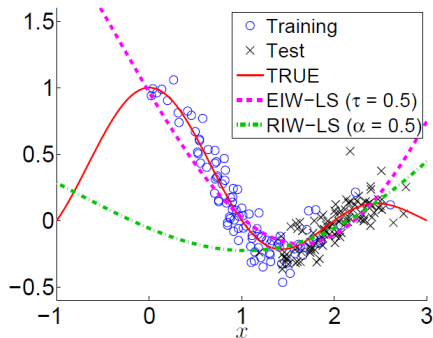
where $r(\mathbf{x})$ is a weight function (e.g., density-ratio).

Take the derivative w.r.t. \mathbf{w} and equating it to zero.

$$\begin{aligned} \frac{\partial J(\mathbf{w})}{\partial \mathbf{w}} &= -\frac{2}{n_{\text{tr}}} \sum_{i=1}^{n_{\text{tr}}} r(\mathbf{x}_i^{\text{tr}}) (y_i^{\text{tr}} - \mathbf{w}^\top \mathbf{x}_i^{\text{tr}}) \mathbf{x}_i^{\text{tr}} = \mathbf{0} \\ \hat{\mathbf{w}} &= \left(\sum_{i=1}^{n_{\text{tr}}} r(\mathbf{x}_i^{\text{tr}}) \mathbf{x}_i^{\text{tr}} \mathbf{x}_i^{\text{tr}^\top} \right)^{-1} \sum_{i=1}^{n_{\text{tr}}} r(\mathbf{x}_i^{\text{tr}}) y_i^{\text{tr}} \mathbf{x}_i^{\text{tr}} \end{aligned}$$

Covariate Shift Adaptation: Synthetic Example

Comparison of EIW-LS and RIW-LS:



Supervised Transfer Learning

Problem formulation:

- $\{(\mathbf{x}_i^{\text{tr}}, y_i^{\text{tr}})\}_{i=1}^{n_{\text{tr}}} \stackrel{\text{i.i.d.}}{\sim} p_{\text{tr}}(\mathbf{x}, y)$
- $\{(\mathbf{x}_j^{\text{te}}, y_j^{\text{te}})\}_{j=1}^{n_{\text{te}}} \stackrel{\text{i.i.d.}}{\sim} p_{\text{te}}(\mathbf{x}, y), n_{\text{te}} \ll n_{\text{tr}}$

We assume to have a large number of training samples and a small number of paired target labeled samples.

- Frustratingly easy domain adaptation
- Multi-task Learning
- Fine-tuning (Deep Learning)

Supervised Transfer Learning: Importance Weight

Naive approach: Pooling training and test samples

$$\begin{aligned} J(\mathbf{w}) &= \iint \text{loss}(y, f(\mathbf{x}; \mathbf{w})) p_{\text{te}}(\mathbf{x}, y) d\mathbf{x} dy \\ &= \alpha \iint \text{loss}(y, f(\mathbf{x}; \mathbf{w})) p_{\text{tr}}(\mathbf{x}, y) d\mathbf{x} dy \\ &\quad + (1 - \alpha) \iint \text{loss}(y, f(\mathbf{x}; \mathbf{w})) p_{\text{te}}(\mathbf{x}, y) d\mathbf{x} dy \\ &\simeq \frac{\alpha}{n_{\text{tr}}} \sum_{i=1}^{n_{\text{tr}}} \text{loss}(y_i^{\text{tr}}, f(\mathbf{x}_i^{\text{tr}}; \mathbf{w})) + \frac{(1 - \alpha)}{n_{\text{te}}} \sum_{j=1}^{n_{\text{te}}} \text{loss}(y_j^{\text{te}}, f(\mathbf{x}_j^{\text{te}}; \mathbf{w})), \end{aligned}$$

where $0 \leq \alpha \leq 1$ is a tuning parameter to control trade off between source and target errors. Let us define the pooled input matrix and the pooled output vector:

Supervised Transfer Learning: Multi-task Learning

Problem formulation:

- Task1: $\{(\mathbf{x}_i^{(1)}, y_i^{(1)})\}_{i=1}^{n_1} \stackrel{\text{i.i.d.}}{\sim} p(\mathbf{x}, y)$
- Task2: $\{(\mathbf{x}_j^{(2)}, y_j^{(2)})\}_{j=1}^{n_2} \stackrel{\text{i.i.d.}}{\sim} p'(\mathbf{x}, y)$
- ...
- TaskM: $\{(\mathbf{x}_j^{(M)}, y_j^{(M)})\}_{j=1}^{n_M} \stackrel{\text{i.i.d.}}{\sim} p'(\mathbf{x}, y)$
- Linear Models:

$$f_1(\mathbf{x}^{(1)}) = \mathbf{w}_1^\top \mathbf{x}^{(1)}, f_2(\mathbf{x}^{(2)}) = \mathbf{w}_2^\top \mathbf{x}^{(2)}, \dots, f_M(\mathbf{x}^{(M)}) = \mathbf{w}_M^\top \mathbf{x}^{(M)}$$

$$\min_{\mathbf{w}_1, \dots, \mathbf{w}_M} \sum_{m=1}^M \frac{1}{n_m} \sum_{i=1}^{n_m} \text{loss}(y_i^{(m)}, f_m(\mathbf{x}^{(m)})) + \lambda R(\mathbf{w}_1, \dots, \mathbf{w}_M).$$

where $R(\mathbf{w}_1, \dots, \mathbf{w}_M)$ is a regularizer.

- $\lambda = 0$: Independently optimize \mathbf{w} s
- $\lambda > 0$: We share some information among models.

Supervised Transfer Learning: Multi-task Learning

Multi-task learning optimization (Graph-Laplacian).

$$\min_{\mathbf{w}_1, \dots, \mathbf{w}_M} \sum_{m=1}^M \frac{1}{n_m} \sum_{i=1}^{n_m} \text{loss}(y_i^{(m)}, f_m(\mathbf{x}_i^{(m)})) + \lambda \sum_{m=1}^M \sum_{m'=1}^M r_{m,m'} \|\mathbf{w}_m - \mathbf{w}_{m'}\|_2^2.$$

where $r_{m,m'} \geq 0$ is a model parameter (similarity between models). If $r_{m,m'} > 0$, we make \mathbf{w}_m and $\mathbf{w}_{m'}$ close.

Supervised Transfer Learning: Multi-task Learning

Other approach: Explicitly including shared parameter. We decompose

$$\mathbf{w}_m = \mathbf{w}_0 + \mathbf{v}_m$$

That is

- $f_1(\mathbf{x}^{(1)}) = (\mathbf{w}_0 + \mathbf{v}_1)^\top \mathbf{x}^{(1)},$
- $f_2(\mathbf{x}^{(2)}) = (\mathbf{w}_0 + \mathbf{v}_2)^\top \mathbf{x}^{(2)},$
- ...
- $f_M(\mathbf{x}^{(M)}) = (\mathbf{w}_0 + \mathbf{v}_M)^\top \mathbf{x}^{(M)}$

where \mathbf{w}_0 is a common factor for all models.

For squared-loss, we can write the problem as

$$\min_{\mathbf{w}_1, \dots, \mathbf{w}_M} \frac{1}{2} \sum_{m=1}^M \frac{1}{n_m} \sum_{i=1}^{n_m} \left(y_i^{(m)} - (\mathbf{w}_0 + \mathbf{v}_m)^\top \mathbf{x}_i^{(m)} \right)^2 + \lambda (\|\mathbf{w}_0\|_2^2 + \sum_{m=1}^M \|\mathbf{v}_m\|_2^2)$$

Supervised Transfer Learning: Frustratingly easy domain adaptation

A **frustratingly easy** feature augmentation approach:

$$\begin{aligned}\mathbf{z}^{\text{tr}} &= [\mathbf{x}^{\text{tr}\top} \quad \mathbf{x}^{\text{tr}\top} \quad \mathbf{0}_d^\top]^\top, \\ \mathbf{z}^{\text{te}} &= [\mathbf{x}^{\text{te}\top} \quad \mathbf{0}_d^\top \quad \mathbf{x}^{\text{te}\top}]^\top,\end{aligned}$$

The inner product of \mathbf{z} in the same domain is give as

$$\begin{aligned}\mathbf{z}^{\text{tr}\top} \mathbf{z}^{\text{tr}} &= 2\mathbf{x}^{\text{tr}\top} \mathbf{x}^{\text{tr}}, \\ \mathbf{z}^{\text{te}\top} \mathbf{z}^{\text{te}} &= 2\mathbf{x}^{\text{te}\top} \mathbf{x}^{\text{te}},\end{aligned}$$

while we have

$$\mathbf{z}^{\text{tr}\top} \mathbf{z}^{\text{te}} = \mathbf{x}^{\text{tr}\top} \mathbf{x}^{\text{te}}, .$$

Then, we train a supervised learning method with the transformed vectors \mathbf{z} . **Super easy!!!!**

Supervised Transfer Learning: Multi-task Learning

Actually, supervised transfer learning can be regarded as a **two-task** learning problem. First task is for training and second task is for test. Let us denote the transformed vectors as

$$\begin{aligned}\mathbf{z}^{\text{tr}} &= [\mathbf{x}^{\text{tr}\top} \quad \mathbf{x}^{\text{tr}\top} \quad \mathbf{0}_d^\top]^\top, \\ \mathbf{z}^{\text{te}} &= [\mathbf{x}^{\text{te}\top} \quad \mathbf{0}_d^\top \quad \mathbf{x}^{\text{te}\top}]^\top.\end{aligned}$$

And, we consider a linear regression problem: The model parameter of the linear model can be written as

$$\mathbf{w} = [\mathbf{w}_0^\top \quad \mathbf{v}_1^\top \quad \mathbf{v}_2^\top]^\top \in \mathbb{R}^{3d}$$

Supervised Transfer Learning: Multi-task Learning

$$\begin{aligned} J(\mathbf{w}) &= \frac{1}{2n_{\text{tr}}} \sum_{i=1}^{n_{\text{tr}}} \|y_i^{\text{tr}} - \mathbf{z}_i^{\text{tr}\top} \mathbf{w}\|_2^2 + \frac{1}{2n_{\text{te}}} \sum_{i=1}^{n_{\text{te}}} \|y_i^{\text{te}} - \mathbf{z}_i^{\text{te}\top} \mathbf{w}\|_2^2 + \lambda \|\mathbf{w}\|_2^2 \\ &= \frac{1}{2} \sum_{m=1}^M \frac{1}{n_m} \sum_{i=1}^{n_m} \left(y_i^{(m)} - (\mathbf{w}_0 + \mathbf{v}_m)^\top \mathbf{x}_i^{(m)} \right)^2 + \lambda (\|\mathbf{w}_0\|_2^2 + \sum_{m=1}^M \|\mathbf{v}_m\|_2^2), \end{aligned}$$

where we use

$$\begin{aligned} \mathbf{w}^\top \mathbf{z}^{\text{tr}} &= (\mathbf{w}_0 + \mathbf{v}_1)^\top \mathbf{x}^{\text{tr}}, & \mathbf{w}^\top \mathbf{z}^{\text{te}} &= (\mathbf{w}_0 + \mathbf{v}_2)^\top \mathbf{x}^{\text{te}} \\ \mathbf{x}^{\text{tr}} &= \mathbf{x}^{(1)}, & \mathbf{x}^{\text{te}} &= \mathbf{x}^{(2)}, \\ \|\mathbf{w}\|_2^2 &= \|\mathbf{w}_0\|_2^2 + \sum_{m=1}^2 \|\mathbf{v}_m\|_2^2. \end{aligned}$$

Frustratingly easy domain adaptation is a multi-task learning.

- **Semi-supervised learning**. Use unlabeled samples and assume the data distribution of unlabeled data is same as training.
- Weighted Maximum Likelihood, Graph-based method.
- **Transfer Learning**. Use samples from test data. **Training and test distributions are different**.
- Covariate shift adaptation, frustratingly easy domain adaptation.