

# On the Rates of Information Transfer in Multi-Task Learning using Random Matrix Theory

24/02/2020 at 5:58pm

## 1 Problem Setup

In the multi-task learning (MTL) problem, we are given the input of  $k$  tasks  $(X_1, Y_1), (X_2, Y_2), \dots, (X_k, Y_k)$ . We use a shared module  $B$  for all tasks and a separate module  $\{A_i\}_{i=1}^k$  for each task. This corresponds to minimizing the following optimization objective.

$$f(B; A_1, \dots, A_k) = \sum_{i=1}^k \|X_i B A_i - Y_i\|_F^2. \quad (1.1)$$

When the models are the same for all tasks, we use a simplified objective as follows.

$$f(w) = \sum_{i=1}^k \|X_i w - Y_i\|_F^2. \quad (1.2)$$

Equation (1.1) is a slightly more complicated parametrization of equation (1.2), in that each  $A_i$  can select a model from the subspace of  $B$  to fit  $(X_i, Y_i)$ .

### 1.1 Hypothesis

Our hypothesis is that the heterogeneity among the multiple tasks can be categorized into two classes, *covariate shift* and *model shift*.

#### 1.1.1 Covariate Shift

A natural setting is when the covariance matrices  $X_i^\top X_i$  are different, i.e. having different spectrum. This is also known as covariate shift in the literature. Our hypothesis is that the covariate shift can slow down the convergence of learning the true  $\theta$  as a function of the number of data points.

In this setting, we assume that the covariates of task  $i$  are drawn from  $\Sigma_i \in \mathbb{R}^{d \times d}$ , but the models are the same across all the tasks, i.e.

$$y_i = X_i \beta + \varepsilon_i, \text{ with } \frac{1}{n_i} X_i^\top X_i \sim \Sigma_i \quad (1.3)$$

#### 1.1.2 Model Shift

The most general version is when the covariates are different between different tasks and the single-task models are also different. Here the hypothesis is that the optimal  $B$  is captured by a low-rank approximation of the single-task models.

We remove the assumption that the models are the same. The  $i$ -th task data can be viewed as generated by a separate model  $\beta_i \in \mathbb{R}^d$ .

$$y_i = X_i \beta_i + \varepsilon_i. \quad (1.4)$$

## 1.2 Objectives

We will designate the  $k$ -th task as the target. Our goal is to come up with an estimator  $\hat{\beta}$  to provide accurate predictions for the target task. Concretely, we focus on two objectives.

- Estimation error for the target model  $\beta_t$ : we consider their distance

$$\mathbb{E}_{\varepsilon} \left[ \left\| \hat{\beta} - \beta_t \right\|^2 \right].$$

- Test error for the target task:

$$\mathbb{E}_{x \sim \Sigma_k} \left[ (x^\top \hat{\beta} - x^\top \beta_t)^2 \right].$$

## 2 Results for Two Tasks

We would like to get insight on how covariate and model shifts result in slower rates of transfer. For the case of two tasks, we can get precise rates in the low-dimensional setting using random matrix theory. Since there are only two tasks, we call task 1 the source task and task 2 the target task, i.e.  $\beta_1 = \beta_s$  and  $\beta_2 = \beta_t$ .

### 2.1 Covariate Shift

First, we show that if  $\beta_s$  and  $\beta_t$  are equal, then combining the source and target task always helps. The estimation error using the source and target together is

$$e(\hat{\beta}_{s,t}) = \sigma^2 \cdot \text{Tr}[(X_1^\top X_1 + X_2^\top X_2)^{-1}]. \quad (2.1)$$

The estimation error using the target alone is

$$e(\hat{\beta}_t) = \sigma^2 \cdot \text{Tr}[X_2^\top X_2^{-1}]. \quad (2.2)$$

**Proposition 2.1.** *When there is no model shift, adding the source task data always reduces the estimation error, i.e.*

$$e(\hat{\beta}_{s,t}) \leq e(\hat{\beta}_t).$$

*Proof.* This is simply because  $\text{Tr}[(X_1^\top X_1 + X_2^\top X_2)^{-1}] \leq \text{Tr}[X_1^\top X_1^{-1}]$ .  $\square$

Next, we calculate the amount of improvement by comparing equation (2.1) to equation (2.2). We study the spectrum of the random matrix model:

$$Q = \Sigma_1^{1/2} X_1^\top X_1 \Sigma_1^{1/2} + \Sigma_2^{1/2} X_2^\top X_2 \Sigma_2^{1/2},$$

where  $\Sigma_{1,2}$  are  $p \times p$  deterministic covariance matrices, and  $X_1 = (x_{ij})_{1 \leq i \leq n_1, 1 \leq j \leq p}$  and  $X_2 = (x_{ij})_{n_1+1 \leq i \leq n_1+n_2, 1 \leq j \leq p}$  are  $n_1 \times p$  and  $n_2 \times p$  random matrices, respectively, where the entries  $x_{ij}$ ,  $1 \leq i \leq n_1 + n_2 \equiv n$ ,  $1 \leq j \leq p$ , are real independent random variables satisfying

$$\mathbb{E} x_{ij} = 0, \quad \mathbb{E} |x_{ij}|^2 = n^{-1}. \quad (2.3)$$

For now, we assume that the random variables  $x_{ij}$  are i.i.d. Gaussian, but we know that universality holds for generally distributed entries. We shall consider the high-dimensional setting such that

$$\gamma_n := \frac{p}{n} \rightarrow \gamma, \quad c_n := \frac{n_1}{n} \rightarrow c, \quad \text{as } n \rightarrow \infty,$$

for some constants  $\gamma \in (0, \infty)$  and  $c \in (0, 1)$ .

**Lemma 2.2.** *Put in the statement* We have with high probability  $1 - o(1)$ ,

$$\frac{1}{p} \text{Tr}(Q^{-1}) = \frac{1}{p} \text{Tr} \left[ \frac{1}{a_3 \Sigma_1 + a_4 \Sigma_2} \right] + O(n^{-1+\varepsilon})$$

for any constant  $\varepsilon > 0$ , where  $a_{3,4}$  are found using equations in (2.7).

We assume that  $\Sigma_1^{-1/2} \Sigma_2$  has eigendecomposition

$$\Sigma_1^{-1/2} \Sigma_2^{1/2} = O D O^T, \quad D = \text{diag}(d_1, \dots, d_p). \quad (2.4)$$

Then by the rotational invariance of Gaussian matrices, we have

$$\tilde{Q} \stackrel{d}{=} \Sigma_1^{1/2} O \tilde{Q} O^T \Sigma_1^{1/2}, \quad \tilde{Q} := X_1^T X_1 + D X_2^T X_2 D.$$

Thus we study the spectrum of  $\tilde{Q}$  instead. We define  $\mathcal{G}(z) := (\tilde{Q} - z)^{-1}$  for  $z \in \mathbb{C}_+$ . With some random matrix tools, we have that

$$\mathcal{G}(z) \approx \text{Diag} \left( \frac{1}{-z(1 + m_3(z) + d_i^2 m_4(z))} \right)_{1 \leq i \leq p} = \frac{1}{-z(1 + m_3(z) + D^2 m_4(z))}$$

in certain sense. Here  $m_{3,4}(z)$  satisfy the following self-consistent equations

$$\frac{n_1}{n} \frac{1}{m_3} = -z + \frac{1}{n} \sum_{i=1}^p \frac{1}{1 + m_3 + d_i^2 m_4}, \quad \frac{n_2}{n} \frac{1}{m_4} = -z + \frac{1}{n} \sum_{i=1}^p \frac{d_i^2}{1 + m_3 + d_i^2 m_4} \quad (2.5)$$

When  $z \rightarrow 0$ , we shall have

$$m_3(z) = -\frac{a_3}{z} + O(1), \quad m_4(z) = -\frac{a_4}{z} + O(1), \quad a_3, a_4 > 0.$$

Then for  $z \rightarrow 0$ , the equations in (2.5) are reduced to

$$\frac{n_1}{n} \frac{1}{a_3} = 1 + \frac{1}{n} \sum_{i=1}^p \frac{1}{a_3 + d_i^2 a_4}, \quad \frac{n_2}{n} \frac{1}{a_4} = 1 + \frac{1}{n} \sum_{i=1}^p \frac{d_i^2}{a_3 + d_i^2 a_4}. \quad (2.6)$$

First, it is easy to see that these equations are equivalent to

$$a_3 + a_4 = 1 - \gamma_n, \quad a_3 + \frac{1}{n} \sum_{i=1}^p \frac{a_3}{a_3 + d_i^2 [(1 - \gamma_n) - a_3]} = c_n. \quad (2.7)$$

Furthermore, we have

$$\begin{aligned} \text{Tr}(Q^{-1}) &= \lim_{z \rightarrow 0} \text{Tr} \left[ \Sigma_1^{-1/2} O \mathcal{G}(z) O^T \Sigma_1^{-1/2} \right] = \text{Tr} \left[ \Sigma_1^{-1/2} O \left( \frac{1}{a_3 + D^2 a_4} \right) O^T \Sigma_1^{-1/2} \right] \\ &= \text{Tr} \left[ \Sigma_1^{-1/2} \frac{1}{a_3 + \Sigma_1^{-1} \Sigma_2 a_4} \Sigma_1^{-1/2} \right] = \text{Tr} \left[ \frac{1}{a_3 \Sigma_1 + a_4 \Sigma_2} \right]. \end{aligned}$$

## 2.2 Model Shift

In this case,  $\beta_s$  and  $\beta_t$  are different. If we put the two tasks together, we get

$$\hat{\beta}_{s,t} = (X_1^\top X_1 + X_2^\top X_2)^{-1} (X_1^\top Y_1 + X_2^\top Y_2) \quad (2.8)$$

$$= (X_1^\top X_1 + X_2^\top X_2)^{-1} ((X_1^\top X_1 \beta_s + X_2^\top X_2 \beta_t) + (X_1^\top \varepsilon_1 + X_2^\top \varepsilon_2)) \quad (2.9)$$

Hence

$$\mathbb{E}_\varepsilon \left[ \left\| \hat{\beta}_{s,t} - \beta_t \right\|^2 \right] = \left\| (X_1^\top X_1 + X_2^\top X_2)^{-1} X_1^\top X_1 (\beta_s - \beta_t) \right\|^2 + \sigma^2 \text{Tr} \left[ (X_1^\top X_1 + X_2^\top X_2)^{-1} \right] \quad (2.10)$$

We are interested in the eigenvalues of

$$(X_1^\top X_1)^{-1} (X_2^\top X_2),$$

which is called a generalized Fisher matrix. As in the previous setting, we write out their covariance explicitly and consider

$$(\Sigma_1^{1/2} X_1^\top X_1 \Sigma_1^{1/2})^{-1} \Sigma_2^{1/2} X_2^\top X_2 \Sigma_2^{1/2},$$

where  $\Sigma_{1,2}$  are  $p \times p$  deterministic covariance matrices, and  $X_1 = (x_{ij})_{1 \leq i \leq n_1, 1 \leq j \leq p}$  and  $X_2 = (x_{ij})_{n_1+1 \leq i \leq n_1+n_2, 1 \leq j \leq p}$  are  $n_1 \times p$  and  $n_2 \times p$  random matrices, respectively, where the entries  $x_{ij}$ ,  $1 \leq i \leq n_1 + n_2 \equiv n$ ,  $1 \leq j \leq p$ , are real independent random variables satisfying (2.3).

For now, we assume that the random variables  $x_{ij}$  are i.i.d. Gaussian, and that  $\Sigma_1^{-1/2} \Sigma_2$  has singular value decomposition

$$\Sigma_1^{-1/2} \Sigma_2^{1/2} = O_1 D O_2, \quad D = \text{diag}(d_1, \dots, d_p).$$

Then it is equivalent to study

$$Q := (X_1^\top X_1)^{-1} D X_2^\top X_2 D,$$

which has the same nonzero eigenvalues of

$$\mathcal{Q} := X_2 D (X_1^\top X_1)^{-1} D X_2^\top,$$

i.e.  $\mathcal{Q}$  has the same nonzero eigenvalues of  $Q$ , but has  $(n_2 - p)$  more zero eigenvalues.

We can study it using the linearization matrix (for my own purpose right now)

$$H = \begin{pmatrix} -zI & X_2 D & 0 \\ D X_2^\top & 0 & X_1^\top \\ 0 & X_1 & I \end{pmatrix}, \quad G(z) := H(z)^{-1}.$$

Then the  $(1,1)$ -th block is equal to  $\mathcal{G}(z) := (Q - z)^{-1}$ . We denote

$$m_1(z) := \frac{1}{n} \text{Tr} \mathcal{G}(z), \quad m(z) = \frac{1}{p} \left[ n m_1(z) + \frac{n_2 - p}{z} \right]. \quad (2.11)$$

We can show that it satisfies the following self-consistent equations together with another  $m_3(z)$ :

$$\frac{n_2}{n} \frac{1}{m_1} = -z + \frac{1}{n} \sum_{i=1}^p \frac{d_i^2}{m_3 + d_i^2 m_1}, \quad \frac{n_1}{n} \frac{1}{m_3} = 1 + \frac{1}{n} \sum_{i=1}^p \frac{1}{m_3 + d_i^2 m_1}. \quad (2.12)$$

For these two equations, we can obtain one single equation for  $m_1(z)$ :

$$m_3 = 1 - \frac{p}{n} + z m_1, \quad \frac{n_2}{n} \frac{1}{m_1} = -z + \frac{1}{n} \sum_{i=1}^p \frac{d_i^2}{1 - \frac{p}{n} + (z + d_i^2) m_1}. \quad (2.13)$$

One can solve the above equation for  $m_1$  with positive imaginary parts, and then calculate  $m(z)$  using (2.11).

With  $m(z)$ , we can define

$$\rho_c(z) := \frac{1}{\pi} \lim_{\eta \downarrow 0} m(z).$$

It will be a compact supported probability density which gives the eigenvalue distribution of  $Q$ . Moreover, we have

$$\frac{1}{p} \text{Tr} \frac{1}{(\Sigma_1^{1/2} X_1^\top X_1 \Sigma_1^{1/2})^{-1} \Sigma_2^{1/2} X_2^\top X_2 \Sigma_2^{1/2} + 1} \left( \frac{1}{(\Sigma_1^{1/2} X_1^\top X_1 \Sigma_1^{1/2})^{-1} \Sigma_2^{1/2} X_2^\top X_2 \Sigma_2^{1/2} + 1} \right)^T \approx \int \frac{\rho_c(x)}{(x+1)^2} dx.$$

Moreover, the right edge  $\lambda_+$  of  $\rho_c$  gives the location of the largest eigenvalue, while the left edge  $\lambda_-$  of  $\rho_c$  gives the location of the smallest eigenvalue (**we have ways to determine the edges by solving some equations, state them later**). They will provide the upper and lower bounds on the operator norm:

$$\frac{1}{1 + \lambda_+} \leq \frac{1}{(\Sigma_1^{1/2} X_1^T X_1 \Sigma_1^{1/2})^{-1} \Sigma_2^{1/2} X_2^T X_2 \Sigma_2^{1/2} + 1} \leq \frac{1}{1 + \lambda_-}.$$