

A combination of Transfer learning and Double machine learning

19331053 纪传宇

2022 年 8 月 21 日

1 Pre-knowledge

Interactive Regression Model (IRM) is a model to estimate causal effects

1.1 Interactive Regression Model (IRM)

We consider estimation of average treatment effects when treatment effects are fully heterogeneous and the treatment variable is binary, $D \in \{0, 1\}$. We consider vectors (Y, D, X) such that

$$\begin{aligned} Y &= g_0(D, X) + U, \quad \mathbb{E}(U \mid X, D) = 0 \\ D &= m_0(X) + V, \quad \mathbb{E}(V \mid X) = 0 \end{aligned} \tag{1}$$

Since D is not additively separable, this model is more general than the partially linear model for the case of binary D . A common target parameter

of interest in this model is the average treatment effect (ATE),

$$\theta_0 = \mathbb{E} [g_0(1, X) - g_0(0, X)] . \quad (2)$$

Another common target parameter is the average treatment effect for the treated (ATTE),

$$\theta_0 = \mathbb{E} [g_0(1, X) - g_0(0, X) \mid D = 1] . \quad (3)$$

1.2 Key Ingredients of the Double Machine Learning Inference Approach

The confounding factors X affect the treatment variable D via the propensity score, $m_0(Z) := E[D \mid X]$, and the outcome variable via the function $g_0(D, Z)$. Both of these functions are unknown and potentially complicated, and we consider estimating these functions via the use of ML methods.

We proceed to set up moment conditions with scores that obey a type of orthogonality with respect to nuisance functions. Specifically, we make use of scores $\psi(W, \theta, \eta)$ that satisfy the identification condition

$$E\psi(W, \theta_0, \eta_0) = 0 \quad (4)$$

The first key input of the inference procedure is using a score function $\psi(W; \theta; \eta)$ that satisfies (4), with $\bar{\theta}_0$ being the unique solution, and that obeys the Neyman orthogonality condition

$$\partial_\eta E\psi(W, \theta_0, \eta)|_{\eta=\eta_0} = 0, \quad (5)$$

where $W=(Y, D, X)$, θ_0 is the parameter of interest, η denotes nuisance functions with population value η_0 and $\partial_\eta f|_{\eta=\eta_0}$ denotes the derivative of f with respect to η (the Gateaux derivative operator).

Neyman-orthogonal scores are readily available for both the ATE and ATTE one can use the doubly robust/efficient scores of Robins and Rotnitzky (1995) and Hahn (1998) which are automatically Neyman orthogonal.

For estimating the ATE, we employ

$$\psi(W, \theta, \eta) := g(1, X) - g(0, X) + \frac{D(Y - g(1, X))}{m(X)} - \frac{(1 - D)(Y - g(0, X))}{1 - m(X)} - \theta \quad (6)$$

with

$$\begin{aligned} \eta(X) &:= (g(0, X), g(1, X), m(X))' \\ \eta_0(X) &:= (g_0(0, X), g_0(1, X), m_0(X))' \end{aligned} \quad (7)$$

where $\eta(X)$ is the nuisance parameter with true value denoted by $\eta_0(X)$ consisting of P -square integrable functions.

The second key input is to use a form of sample splitting at the stage of producing the estimator of the main parameter θ_0 , which allows us to avoid biases arising from overfitting.

1.3 Double Machine Learning for Estimation of a Causal Parameter

We describe the estimator of θ_0 using random sample $(W_i)_{i=1}^N$. The algorithm makes use of a form of sample splitting, which we call cross-fitting. It builds on the ideas in, e.g., Angrist and Krueger (1995). The use of sample-splitting is a crucial ingredient to the approach that helps avoid overfitting which can easily result from the application of complex, flexible methods such as boosted linear and tree models, random forests, and various ensemble and hybrid ML methods. A. Algorithm: Estimation by K -fold Cross-Fitting

Step 1: Let K be a fixed integer. Form a K -fold random partition of $\{1, \dots, N\}$ by dividing it into equal parts $(I_k)_{k=1}^K$ each of size $n := N / K$, assuming that N is a multiple of K. For each set I_k , let I_k^c denote all observation indices that are not in I_k .

Step 2: Construct K estimators

$$\check{\theta}_0(I_k, I_k^c), \quad k = 1, \dots, K, \quad (8)$$

that employ the machine learning estimators

$$\begin{aligned} \hat{\eta}_0(I_k^c) = & (\hat{g}_0(0, Z; I_k^c), \hat{g}_0(1, Z; I_k^c) \\ & \hat{m}_0(Z; I_k^c), \frac{1}{N-n} \sum_{i \in I_k^c} D_i)' , \end{aligned} \quad (9)$$

of the nuisance parameters

$$\eta_0(X) = (g_0(0, X), g_0(1, X), m_0(X), E[D])' , \quad (10)$$

and where each estimator $\check{\theta}_0(I_k, I_k^c)$ is defined as the root θ of

$$\frac{1}{n} \sum_{i \in I_k} \psi(W, \theta, \hat{\eta}_0(I_k^c)) = 0, \quad (11)$$

for the score ψ defined in (6) for the ATE.

Step 3: Average the K estimators to obtain the final estimator:

$$\tilde{\theta}_0 = \frac{1}{K} \sum_{k=1}^K \check{\theta}_0(I_k, I_k^c) \quad (12)$$

An approximate standard error for this estimator is $\hat{\sigma}/\sqrt{N}$, where

$$\hat{\sigma}^2 = \frac{1}{N} \sum_{i=1}^N \hat{\psi}_i^2, \hat{\psi}_i := \psi(W_i, \tilde{\theta}_0, \hat{\eta}_0(I_{k(i)}^c)) , \quad (13)$$

and

$$k(i) := \{k \in \{1, \dots, K\} : i \in I_k\} . \quad (14)$$

An approximate $(1-\alpha) \times 100$ percent confidence interval is

$$CI_n := \left[\tilde{\theta}_0 \pm \Phi^{-1}(1 - \alpha/2) \hat{\sigma} / \sqrt{N} \right] . \quad (15)$$

1.4 Transfer learning method and algorithm

The main strategy is to first transfer the information from those sources by pooling all the data to obtain a rough estimator, then correct the bias in the second step using the target data. More specifically, we fit a GLM with ℓ_1 -penalty by pooled samples first, then fit the contrast in the second step

using only the target by another ℓ_1 -regularization. The detailed algorithm (\mathcal{A} -Trans-GLM) is presented in Algorithm 1. The transferring step could be understood as to solve the following equation w.r.t. $\mathbf{w} \in \mathbb{R}^p$:

$$\sum_{k \in \{0\} \cup \mathcal{A}} \left[\left(\mathbf{X}^{(k)} \right)^T \mathbf{y}^{(k)} - \sum_{i=1}^{n_k} \psi' \left((\mathbf{w})^T \mathbf{x}_i^{(k)} \right) \mathbf{x}_i^{(k)} \right] = \mathbf{0}_p \quad (16)$$

which converges to the solution of its population version under certain conditions

$$\sum_{k \in \{0\} \cup \mathcal{A}} \alpha_k \mathbb{E} \left\{ \left[\psi' \left((\mathbf{w}^{\mathcal{A}})^T \mathbf{x}^{(k)} \right) - \psi' \left((\mathbf{w}^{(k)})^T \mathbf{x}^{(k)} \right) \right] \mathbf{x}^{(k)} \right\} = \mathbf{0}_p \quad (17)$$

where $\alpha_k = \frac{n_k}{n_{\mathcal{A}} + n_0}$. Notice that in the linear case, $\mathbf{w}^{\mathcal{A}}$ can be explicitly expressed as a linear transformation of the true parameter $\mathbf{w}^{(k)}$, i.e. $\mathbf{w}^{\mathcal{A}} = \Sigma^{-1} \sum_{k \in \{0\} \cup \mathcal{A}} \alpha_k \Sigma^{(k)} \mathbf{w}^{(k)}$, where $\Sigma^{(k)} = \mathbb{E} \left[\mathbf{x}^{(k)} (\mathbf{x}^{(k)})^T \right]$ and $\Sigma = \sum_{k \in \{0\} \cup \mathcal{A}} \alpha_k \mathbb{E} \left[\mathbf{x}^{(k)} (\mathbf{x}^{(k)})^T \right]$ (Li et al., 2021).

We can use Algorithm 1 if we are certain about which sources to transfer.

Algorithm 1: \mathcal{A} -Trans-GLM

Input: target data $(\mathbf{X}^{(0)}, \mathbf{y}^{(0)})$, source data $\{(\mathbf{X}^{(k)}, \mathbf{y}^{(k)})\}_{k=1}^K$, penalty parameters

$\lambda_{\mathbf{w}}$ and λ_{δ} , transferring set \mathcal{A}

Output: the estimated coefficient vector $\hat{\beta}$

1 **Transferring step:** Compute

$$\hat{\mathbf{w}}^{\mathcal{A}} \leftarrow \arg \min_{\mathbf{w}} \left\{ \frac{1}{n_{\mathcal{A}} + n_0} \sum_{k \in \{0\} \cup \mathcal{A}} \left[-(\mathbf{y}^{(k)})^T \mathbf{X}^{(k)} \mathbf{w} + \sum_{i=1}^{n_k} \psi(\mathbf{w}^T \mathbf{x}_i^{(k)}) \right] + \lambda_{\mathbf{w}} \|\mathbf{w}\|_1 \right\}$$

2 **Debiasing step:** Compute

$$\hat{\delta}^{\mathcal{A}} \leftarrow \arg \min_{\delta} \left\{ -\frac{1}{n_0} (\mathbf{y}^{(0)})^T \mathbf{X}^{(0)} (\hat{\mathbf{w}}^{\mathcal{A}} + \delta) + \frac{1}{n_0} \sum_{i=1}^{n_0} \psi((\hat{\mathbf{w}}^{\mathcal{A}} + \delta)^T \mathbf{x}_i^{(0)}) + \lambda_{\delta} \|\delta\|_1 \right\}$$

3 Let $\hat{\beta} \leftarrow \hat{\mathbf{w}}^{\mathcal{A}} + \hat{\delta}^{\mathcal{A}}$

4 Output $\hat{\beta}$

Fig. 1. Transfer Learning \mathcal{A} -Trans-GLM

We propose a simple, algorithm-free, and data-driven method to determine an informative transferring set \mathcal{A} . We call this approach a transferable

source detection algorithm and refer to it as Trans-GLM

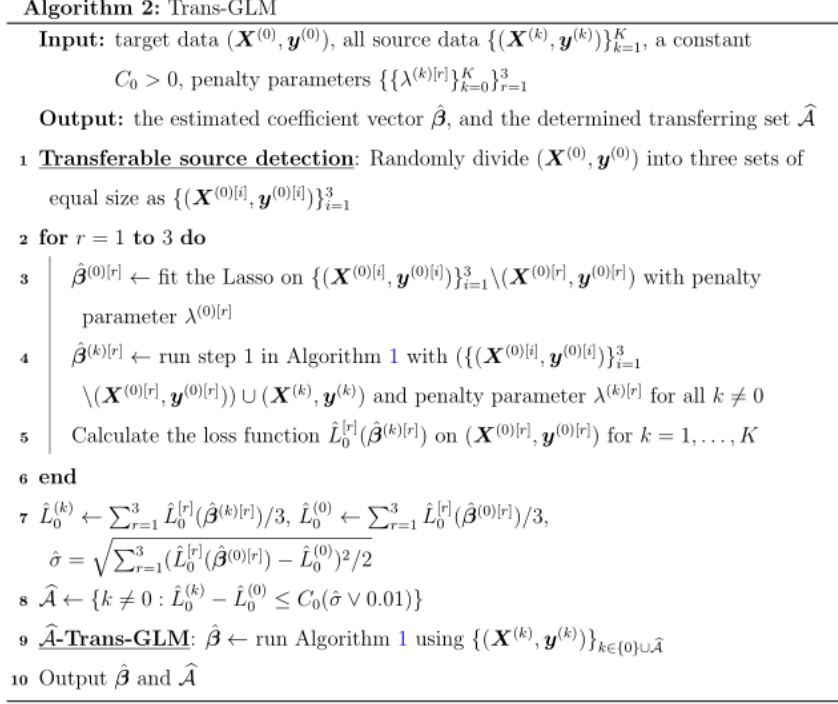


Fig. 2. Transfer Learning Trans-GLM

2 Combination of Transfer learning and Double Machine Learning

In IRM model, We consider the case where treatment effects are fully heterogeneous and the treatment variable, D , is binary, $D \in \{0, 1\}$. We let Y denote the outcome variable of interest and X denote a set of control variables. We then model random vector (Y, D, X) as

$$\begin{aligned} (1) Y &= g_0(D, X) + \zeta, & E[\zeta \mid X, D] &= 0, \\ (2) D &= m_0(X) + \nu, & E[\nu \mid X] &= 0. \end{aligned} \tag{18}$$

The second equation keeps track of confounding, namely the dependence of D on covariates/controls. The characteristics X affect variable D via the function $m_0(X)$ and the outcome variable via the function $g_0(X)$.

Here are three hypothesis

(1) X is a high-dimensional vector that $X = (X_1, \dots, X_p)$ consists of other confounding covariates

(2) m_0 is a logistic model with respect to a dictionary of basis functions with respect to X and the end value is binary.

(3) g_0 is linear with respect to a dictionary of basis functions with respect to X and D

2.1 Method 1

We can use the second key of Double Machine which is cross fitting to estimate the parameter of m_0 .

In "Online Appendix for "Double/Debiased/Neyman Machine Learning of Treatment Effects"" Then report results based on five simple methods for estimating the nuisance functions used in forming the orthogonal estimating equations. They consider three tree-based methods, labeled "Random Forest", "Reg. Tree", and "Boosting", one ℓ_1 -penalization based method, labeled "Lasso", and a neural network method, labeled "Neural Net".

Because we assume that m_0 is logistic with respect to a dictionary of basis functions with respect to X . We choose the "Lasso" method to estimate the parameter of m_0 .

Through using "rlassologit" in "Lasso" method combining corss-fitting, we get the estimation of parameters .

We then use the parameters and original X to fix D and get \hat{D}

Because g_0 is linear with respect to a dictionary of basis functions with respect to X and D . So g_0 can be represented as $g_0(X, D) = \beta_0 + \beta_1 x_1 + \dots + \beta_n x_n + \theta d$

We use Transfer Learning under Generalized Linear Models and R pack-

age glmtrans to estimate the parameters $\hat{\beta}_i$ $i=1,2,\dots,n$ and $\hat{\theta}$

2.2 Method 2

D is binary, so we can consider D is a reponse variable then in order to estimate the parameters in m_0 we use Transfer Learning under Generalized Linear Models in Logistic case which the reponse variable in this model is 0-1 variable.

Then we use the parameters and original X to fix D to get \hat{D}

Because g_0 is linear with respect to a dictionary of basis functions with respect to X and D. So g_0 can be represented as $g_0(X, D) = \beta_0 + \beta_1 x_1 + \dots + \beta_n x_n + \theta d$

We use Transfer Learning under Generalized Linear Models and R package glmtrans to estimate the parameters $\hat{\beta}_i$ $i=1,2,\dots,n$ and $\hat{\theta}$.

3 Simulations

MSEtheta denote through estimating m_0 , the MSE value of $\hat{\beta}$ and real β

MSEerror1 denote through estimating m_0 , the different quantity between \hat{D} and D

MSEtotal denote through estimating m_0 and Transfer learning to estimate g_0 , the MSE value of $\hat{\alpha}$ and setting α plus $\hat{\theta}$ and setting θ

MSEerror2 denote through estimating m_0 and Transfer learning to estimate g_0 , the MSE value of $\hat{\theta}$ and setting θ

3.1 Dataset 5

The data generating process is defined as

$$d_i = 1 \left\{ \frac{\exp(x'_i \beta)}{1 + \exp(x'_i \beta)} > v_i \right\} \quad (19)$$

$$y_i = \theta d_i + x'_i \alpha + \zeta_i$$

with $v_i \sim \mathcal{U}(0.5, 0.2)$, $\zeta_i \sim \mathcal{N}(0, 1)$ and covariates $\mathbf{x}_i \stackrel{i.i.d.}{\sim} \mathcal{N}(\mathbf{0}_p, \mathbf{1}_p)$, $\beta = (3 \cdot \mathbf{1}_s, \mathbf{0}_{p-s})^T$ where s is set to be 5, $\theta = 2$ and $\alpha = (2 \cdot \mathbf{1}_t, \mathbf{0}_{p-t})^T$ where t is set to be 5, p is set to be 30.

We set the target sample size and source sample size for each $k = 1, 5$ are variable. The dimension $p=31$ for both target and source data. The sample size is set to be 100, 150, ..., 450, 500.

For each sample size, we run 5 simulations and average the results.

3.1.1 Method 1

The four indicators are as follow:

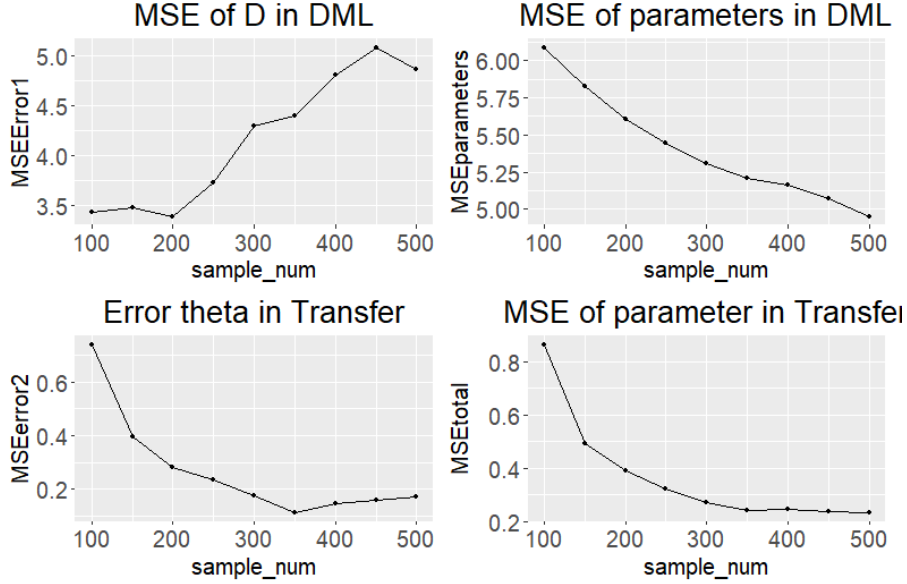


Fig. 3. MSEtheta MSEError1 MSEtotal and MSEError2 in dataset 5

From Fig.3, it can be seen that when fold is set to be 5,6,7, the MSE of D and MSE of parameters in DML are smaller than other fold. But after Transfer learning, the error in θ are very large, by contrast, the other parameters are smaller.