

# Covariate-shift Robust Adaptive Transfer Learning for High-dimensional Regression

Zelin He

Dept. of Statistics,  
The Pennsylvania State University

- 1 Introduction
  - Background and Motivation
  - Key Contributions
- 2 Covariate-shift Robust Transfer Learning
  - TransFusion: Transfer Learning with a Fused Regularization
  - D-TransFusion: TransFusion under Distributed Setting
- 3 Adaptive Covariate-shift Robust Transfer Learning
  - AdaTrans: Feature-wise Adaptive Transfer Learning
  - Oracle AdaTrans and Theoretical Guarantee of AdaTrans
- 4 Numerical Studies
  - Simulation Examples for TransFusion
  - Simulation Examples for AdaTrans
- 5 Conclusion

# Outline



- 1 Introduction
  - Background and Motivation
  - Key Contributions

# Transfer Learning

■ **Target sample:**  $(\mathbf{X}^{(0)}, \mathbf{y}^{(0)}) \sim P^{(0)}(\mathbf{x}, y) = P^{(0)}(y|\mathbf{x})P^{(0)}(\mathbf{x})$ .

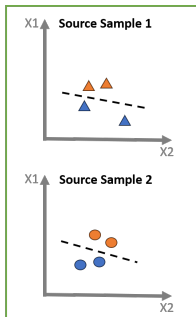
■ **Source samples:**

$$(\mathbf{X}^{(k)}, \mathbf{y}^{(k)}) \sim P^{(k)}(\mathbf{x}, y) = P^{(k)}(y|\mathbf{x})P^{(k)}(\mathbf{x}), \quad k = 1, \dots, K.$$

■ **Goal of transfer learning:**

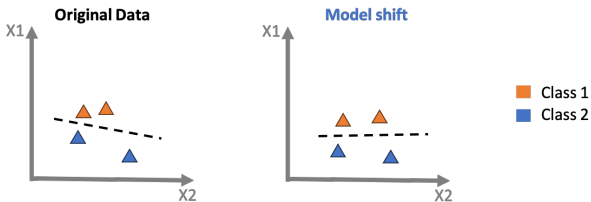
Learn the target model  $P^{(0)}(y|\mathbf{x})$ , by incorporating source information.

E.g. Classification Problem



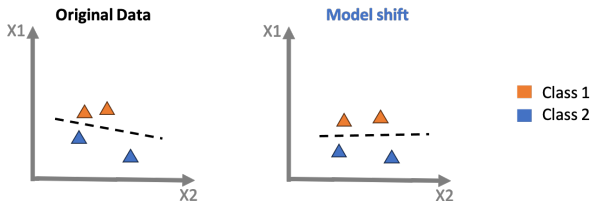
# Model shift in Transfer Learning

- **Model shift:**  $P^{(i)}(y|x) \neq P^{(j)}(y|x)$ ,  $i, j = 0, 1, \dots, K$ .



# Model shift in Transfer Learning

- **Model shift:**  $P^{(i)}(y|x) \neq P^{(j)}(y|x)$ ,  $i, j = 0, 1, \dots, K$ .



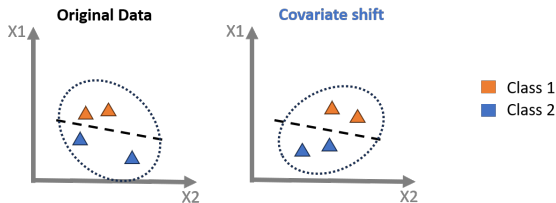
- **Typical strategy to deal with HD model shift:**

## “Pooling training + debiasing”

- Linear model: Li et al. (2022a,b)
- Generalized linear model: Tian et al. (2022)
- Quantile regression model: Jin et al. (2022), Cao and Song (2024)
- ...

# Covariate Shift in Transfer Learning

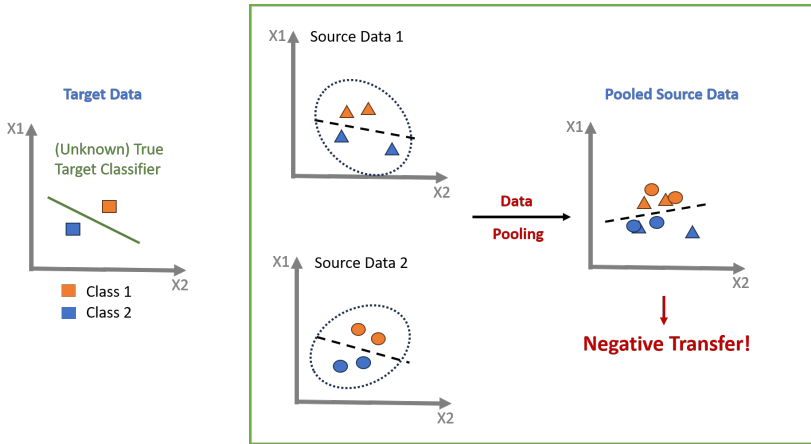
- **covariate shift:**  $P^{(i)}(x) \neq P^{(j)}(x), i, j = 0, 1, \dots, K.$



Under high-dimensionality, covariate shift is often inevitable.

# Covariate Shift in Transfer Learning

## ■ Impact of covariate shift on pooling training:





# Covariate Shift in Transfer Learning



## ■ Existing works to deal with covariate shift:

- **Domain adaptation:** Chen et al. (2016), Redko et al. (2020), etc.
  - typically assumes no model shift;
  - struggles with high-dimensional covariates.
- **Constrained  $\ell_1$ -minimization:** Li et al. (2023), etc.
  - involves multiple nonsmooth constraints;
  - restricted parameter space/strong theoretical assumptions;
  - computationally intractable.

## ■ Existing works to deal with covariate shift:

- **Domain adaptation:** Chen et al. (2016), Redko et al. (2020), etc.
  - typically assumes no model shift;
  - struggles with high-dimensional covariates.
- **Constrained  $\ell_1$ -minimization:** Li et al. (2023), etc.
  - involves multiple nonsmooth constraints;
  - restricted parameter space/strong theoretical assumptions;
  - computationally intractable.

## ⇒ **Our first question:**

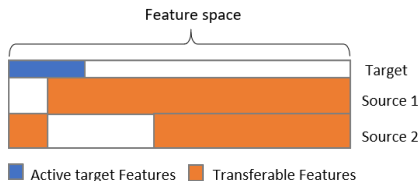
*How to develop a **computationally tractable** method that effectively handles model shift in **high-dimensional** transfer learning, while being **robust to covariate shift**?*

# Feature-specific Transferable Structure



## ■ Feature-specific transferable structure:

- In high-dimensional transfer learning, the transferable structure often varies across features within the same source sample.
- E.g. Whole brain functional connectivity pattern analysis (Li et al., 2018): Each source may have a distinct set of non-transferable features due to variations in brain conditions.

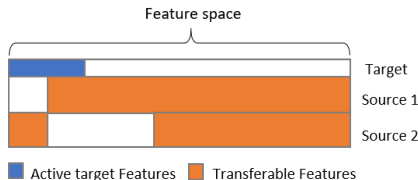


# Feature-specific Transferable Structure



## ■ Feature-specific transferable structure:

- In high-dimensional transfer learning, the transferable structure often varies across features within the same source sample.
- E.g. Whole brain functional connectivity pattern analysis (Li et al., 2018): Each source may have a distinct set of non-transferable features due to variations in brain conditions.



## ⇒ Our second question:

*Is it possible to **auto-detect nontransferable features**/learn transferable structure for each source while transferring source information?*

# Outline



- 1 Introduction
  - Background and Motivation
  - Key Contributions

# Key Contributions of the work

In the **high-dimensional regression** context:

- To tackle the **covariate shift** issue, we
  - 1 develop a new transfer learning framework via fused regularization, named **TransFusion**;
  - 2 extend TransFusion to a distributed setting, called **D-TransFusion**.
- To further address the **feature-specific transferable structure**, we
  - 3 propose a feature-adaptive transfer learning framework, named **AdaTrans**, to auto-detect non-transferable feature.
- For each method, we establish the corresponding non-asymptotic bound.

- 2 Covariate-shift Robust Transfer Learning
  - TransFusion: Transfer Learning with a Fused Regularization
  - D-TransFusion: TransFusion under Distributed Setting

# High-dimensional Linear Regression Model



Sample-level **target** model (with sample size  $n_T$ ):

$$\mathbf{y}^{(0)} = \mathbf{X}^{(0)}\boldsymbol{\beta}^{(0)} + \boldsymbol{\epsilon}^{(0)},$$

Sample-level **source** model (with sample size  $n_S$ ):

$$\mathbf{y}^{(k)} = \mathbf{X}^{(k)}\boldsymbol{\beta}^{(k)} + \boldsymbol{\epsilon}^{(k)} \equiv \mathbf{X}^{(k)}(\boldsymbol{\beta}^{(0)} + \boldsymbol{\delta}^{(k)}) + \boldsymbol{\epsilon}^{(k)}, \quad k = 1, \dots, K.$$

- $E(\boldsymbol{\epsilon}^{(k)}) = \mathbf{0}$ ,  $\text{Cov}(\boldsymbol{\epsilon}^{(k)}) = \sigma^2 \mathbf{I}$ ,  $\boldsymbol{\epsilon}^{(k)} \perp \mathbf{X}^{(k)}$ ,  $k = 0, 1, \dots, K$ .
- $\mathbf{X}_i^{(k)}$  i.i.d sub-Gaussian across  $i$ , with covariance matrix  $\boldsymbol{\Sigma}^{(k)}$ .
- $p \geq n_S \geq n_T$ ,  $\boldsymbol{\beta}^{(0)} \in \mathbb{R}^p$  is high-dimensional yet sparse.
- **Covariate shift**:  $\text{Cov}(\mathbf{X}_i^{(k)}) = \boldsymbol{\Sigma}^{(k)}$  varies across  $k$ .
- **Model shift**:  $\boldsymbol{\delta}^{(k)} \in \mathbb{R}^p$  varies across  $k$ ; take  $\boldsymbol{\delta}^{(0)} = \mathbf{0}$ .



# Objective function of TransFusion

**Formulation:** Estimate  $\beta := ((\beta^{(0)})^\top, (\beta^{(1)})^\top, \dots, (\beta^{(K)})^\top)^\top \in \mathbb{R}^{(K+1)p}$  by solving

$$\operatorname{argmin}_{\beta \in \mathbb{R}^{(K+1)p}} \left\{ \frac{1}{2N} \sum_{k=0}^K \|\mathbf{y}^{(k)} - \mathbf{X}^{(k)} \beta^{(k)}\|_2^2 + \lambda_0 \left( \|\beta^{(0)}\|_1 + \sum_{k=1}^K \nu \|\beta^{(k)} - \beta^{(0)}\|_1 \right) \right\},$$

**Reformulation:** Estimate  $((\beta^{(0)})^\top, (\delta^{(1)})^\top, \dots, (\delta^{(K)})^\top)^\top \in \mathbb{R}^{(K+1)p}$  by solving

$$\operatorname{argmin}_{\beta \in \mathbb{R}^{(K+1)p}} \left\{ \frac{1}{2N} \left( \sum_{k=0}^K \|\mathbf{y}^{(k)} - \mathbf{X}^{(k)} (\beta^{(0)} + \delta^{(k)})\|_2^2 \right) + \lambda_0 \left( \|\beta^{(0)}\|_1 + \sum_{k=1}^K \nu \|\delta^{(k)}\|_1 \right) \right\},$$

- $N = Kn_S + n_T$  is the total sample size,  $\lambda_0$  and  $\nu$  are the tuning parameters.
- The first term measures the average fitness of the models.
- $\lambda_0$  and  $\lambda_0\nu$  respectively take charge of achieving sparsity of  $\beta^{(0)}$  and shrinking the contrast  $\delta^{(k)} = \beta^{(k)} - \beta^{(0)}$  for information transfer.

# Two-step TransFusion Estimator

## Step 1. Joint training:

- 1 Obtain  $\hat{\beta} := ((\hat{\beta}^{(0)})^\top, (\hat{\beta}^{(1)})^\top, \dots, (\hat{\beta}^{(K)})^\top)^\top \in \mathbb{R}^{(K+1)p}$  using a newly advocated proximal gradient descent-based algorithm with message-passing iteration.
- 2 Construct the first-step TransFusion estimator

$$\hat{\beta}_{\text{TF1}}^{(0)} = \sum_{k=1}^K \frac{n_S}{N} \hat{\beta}^{(k)} + \frac{n_T}{N} \hat{\beta}^{(0)}.$$

# Two-step TransFusion Estimator

## Step 1. Joint training:

- 1 Obtain  $\hat{\beta} := ((\hat{\beta}^{(0)})^\top, (\hat{\beta}^{(1)})^\top, \dots, (\hat{\beta}^{(K)})^\top)^\top \in \mathbb{R}^{(K+1)p}$  using a newly advocated proximal gradient descent-based algorithm with message-passing iteration.
- 2 Construct the first-step TransFusion estimator

$$\hat{\beta}_{\text{TF1}}^{(0)} = \sum_{k=1}^K \frac{n_S}{N} \hat{\beta}^{(k)} + \frac{n_T}{N} \hat{\beta}^{(0)}.$$

## Step 2. Local debiasing:

- 3 If necessary, correct the bias of  $\hat{\beta}_{\text{TF1}}^{(0)}$  using the target sample through

$$\hat{\delta} \in \operatorname{argmin}_{\delta \in \mathbb{R}^p} \left\{ \frac{1}{2n_T} \left\| \mathbf{y}^{(0)} - \mathbf{X}^{(0)} \hat{\beta}_{\text{TF1}}^{(0)} - \mathbf{X}^{(0)} \delta \right\|_2^2 + \tilde{\lambda} \|\delta\|_1 \right\},$$

- 4 Define the second-step TransFusion estimator

$$\hat{\beta}_{\text{TF2}}^{(0)} = \hat{\beta}_{\text{TF1}}^{(0)} + \hat{\delta}.$$

## Theorem (Error rate of TransFusion estimators)

Consider the parameter space

$$\Theta(s, h) = \left\{ B = (\beta^{(0)}, \beta^{(1)}, \dots, \beta^{(K)}) : \|\beta^{(0)}\|_0 \leq s, \|\beta^{(k)} - \beta^{(0)}\|_1 \leq h \right\}.$$

Under mild conditions,

- if  $n_S \gg (K^2 h^2 \vee s) \log p$ , with a proper choice of  $\lambda_0$ , with probability at least  $1 - c_1 \exp(-c_2 n_T) - c_3 \exp(-c_4 \log p)$ , we have

$$\|\hat{\beta}_{TF1}^{(0)} - \beta^{(0)}\|_2^2 \lesssim \frac{s \log p}{N} + h \sqrt{\frac{\log p}{n_S}} + \varepsilon_D^2.$$

- if  $n_T \gtrsim s \log p$ ,  $n_S \gg K^2 (h^2 \vee s) \log p$  and  $h \sqrt{\log p / n_T} = o(1)$ , with probability at least  $1 - c_2 \exp(-c_3 \log p)$ , we have

$$\|\hat{\beta}_{TF2}^{(0)} - \beta^{(0)}\|_2^2 \lesssim \frac{s \log p}{N} + h \sqrt{\frac{\log p}{n_T}},$$

# Theoretical Guarantee for TransFusion



## Theorem (Error rate of TransFusion estimators)

$$\|\hat{\beta}_{TF1}^{(0)} - \beta^{(0)}\|_2^2 \lesssim \frac{s \log p}{N} + h \sqrt{\frac{\log p}{n_S}} + \epsilon_D^2;$$
$$\|\hat{\beta}_{TF2}^{(0)} - \beta^{(0)}\|_2^2 \lesssim \frac{s \log p}{N} + h \sqrt{\frac{\log p}{n_T}}.$$

## Theorem (Error rate of TransFusion estimators)

$$\|\hat{\beta}_{TF1}^{(0)} - \beta^{(0)}\|_2^2 \lesssim \frac{s \log p}{N} + h \sqrt{\frac{\log p}{n_S}} + \varepsilon_D^2;$$
$$\|\hat{\beta}_{TF2}^{(0)} - \beta^{(0)}\|_2^2 \lesssim \frac{s \log p}{N} + h \sqrt{\frac{\log p}{n_T}}.$$

- $s \log p / N$ : rate of estimating  $\beta^{(0)}$  based on the full sample with size  $N$ .

# Theoretical Guarantee for TransFusion



## Theorem (Error rate of TransFusion estimators)

$$\|\hat{\beta}_{TF1}^{(0)} - \beta^{(0)}\|_2^2 \lesssim \frac{s \log p}{N} + h \sqrt{\frac{\log p}{n_S}} + \epsilon_D^2;$$
$$\|\hat{\beta}_{TF2}^{(0)} - \beta^{(0)}\|_2^2 \lesssim \frac{s \log p}{N} + h \sqrt{\frac{\log p}{n_T}}.$$

- $s \log p / N$ : rate of estimating  $\beta^{(0)}$  based on the full sample with size  $N$ .
- $h \sqrt{\log p / n_S}$ : rate of estimating  $\delta^{(k)}$ 's based on source samples.

## Theorem (Error rate of TransFusion estimators)

$$\|\hat{\beta}_{TF1}^{(0)} - \beta^{(0)}\|_2^2 \lesssim \frac{s \log p}{N} + h \sqrt{\frac{\log p}{n_S}} + \varepsilon_D^2;$$
$$\|\hat{\beta}_{TF2}^{(0)} - \beta^{(0)}\|_2^2 \lesssim \frac{s \log p}{N} + h \sqrt{\frac{\log p}{n_T}}.$$

- $s \log p / N$ : rate of estimating  $\beta^{(0)}$  based on the full sample with size  $N$ .
- $h \sqrt{\log p / n_S}$ : rate of estimating  $\delta^{(k)}$ 's based on source samples.
- $\varepsilon_D = \|\sum_{k=1}^K \frac{n_S}{N} (\beta^{(k)} - \beta^{(0)})\|_1$ : **task diversity** which measures the bias introduced by averaging.



## Theorem (Error rate of TransFusion estimators)

$$\|\hat{\beta}_{TF1}^{(0)} - \beta^{(0)}\|_2^2 \lesssim \frac{s \log p}{N} + h \sqrt{\frac{\log p}{n_S}} + \varepsilon_D^2;$$
$$\|\hat{\beta}_{TF2}^{(0)} - \beta^{(0)}\|_2^2 \lesssim \frac{s \log p}{N} + h \sqrt{\frac{\log p}{n_T}}.$$

- $s \log p / N$ : rate of estimating  $\beta^{(0)}$  based on the full sample with size  $N$ .
- $h \sqrt{\log p / n_S}$ : rate of estimating  $\delta^{(k)}$ 's based on source samples.
- $\varepsilon_D = \|\sum_{k=1}^K \frac{n_S}{N} (\beta^{(k)} - \beta^{(0)})\|_1$ : **task diversity** which measures the bias introduced by averaging.
- $h \sqrt{\log p / n_T}$ : rate of estimating  $\delta^{(k)}$ 's using the target sample.

## Theorem (Error rate of TransFusion estimators)

$$\|\hat{\beta}_{TF1}^{(0)} - \beta^{(0)}\|_2^2 \lesssim \frac{s \log p}{N} + h \sqrt{\frac{\log p}{n_S}} + \varepsilon_D^2;$$
$$\|\hat{\beta}_{TF2}^{(0)} - \beta^{(0)}\|_2^2 \lesssim \frac{s \log p}{N} + h \sqrt{\frac{\log p}{n_T}}.$$

- $s \log p / N$ : rate of estimating  $\beta^{(0)}$  based on the full sample with size  $N$ .
- $h \sqrt{\log p / n_S}$ : rate of estimating  $\delta^{(k)}$ 's based on source samples.
- $\varepsilon_D = \|\sum_{k=1}^K \frac{n_S}{N} (\beta^{(k)} - \beta^{(0)})\|_1$ : **task diversity** which measures the bias introduced by averaging.
- $h \sqrt{\log p / n_T}$ : rate of estimating  $\delta^{(k)}$ 's using the target sample.
- **Baseline**: target-only lasso with rate  $O(s \log p / n_T)$ .

## Theorem (Error rate of TransFusion estimators)

$$\|\hat{\beta}_{TF1}^{(0)} - \beta^{(0)}\|_2^2 \lesssim \frac{s \log p}{N} + h \sqrt{\frac{\log p}{n_S}} + \varepsilon_D^2;$$
$$\|\hat{\beta}_{TF2}^{(0)} - \beta^{(0)}\|_2^2 \lesssim \frac{s \log p}{N} + h \sqrt{\frac{\log p}{n_T}}.$$

- $s \log p / N$ : rate of estimating  $\beta^{(0)}$  based on the full sample with size  $N$ .
- $h \sqrt{\log p / n_S}$ : rate of estimating  $\delta^{(k)}$ 's based on source samples.
- $\varepsilon_D = \|\sum_{k=1}^K \frac{n_S}{N} (\beta^{(k)} - \beta^{(0)})\|_1$ : **task diversity** which measures the bias introduced by averaging.
- $h \sqrt{\log p / n_T}$ : rate of estimating  $\delta^{(k)}$ 's using the target sample.
- **Baseline**: target-only lasso with rate  $O(s \log p / n_T)$ .
- $\hat{\beta}_{TF1}^{(0)}$  is preferred when  $\varepsilon_D$  or  $n_T$  is small, while  $\hat{\beta}_{TF2}^{(0)}$  is preferred when  $\varepsilon_D$  is non-negligible and  $n_T$  is adequately large.

# Why TF is Robust to Covariate Shifts?



**TransFusion** utilizes task-specific parameters with a series of fused regularizers, resulting in

$$\hat{\beta}_{TF1}^{(0)} - \beta^{(0)} \rightarrow \frac{1}{K} \sum_{k=1}^K (\beta^{(k)} - \beta^{(0)}).$$

# Why TF is Robust to Covariate Shifts?



**TransFusion** utilizes task-specific parameters with a series of fused regularizers, resulting in

$$\hat{\beta}_{TF1}^{(0)} - \beta^{(0)} \rightarrow \frac{1}{K} \sum_{k=1}^K (\beta^{(k)} - \beta^{(0)}).$$

**Pooling training** first pools all data and trains a common model, resulting in

$$\hat{\beta}_{\text{Pooling}}^{(0)} - \beta^{(0)} \rightarrow \left( \sum_{k=1}^K \Sigma^{(k)} \right)^{-1} \sum_{k=1}^K \Sigma^{(k)} (\beta^{(k)} - \beta^{(0)}).$$

# Why TF is Robust to Covariate Shifts?



**TransFusion** utilizes task-specific parameters with a series of fused regularizers, resulting in

$$\hat{\beta}_{TF1}^{(0)} - \beta^{(0)} \rightarrow \frac{1}{K} \sum_{k=1}^K (\beta^{(k)} - \beta^{(0)}).$$

**Pooling training** first pools all data and trains a common model, resulting in

$$\hat{\beta}_{\text{Pooling}}^{(0)} - \beta^{(0)} \rightarrow \left( \sum_{k=1}^K \Sigma^{(k)} \right)^{-1} \sum_{k=1}^K \Sigma^{(k)} (\beta^{(k)} - \beta^{(0)}).$$

The bias can be amplified by a factor (Li et al., 2022)

$$C_{\Sigma} := 1 + \max_{j \leq p} \max_k \left\| e_j^{\top} \left( \Sigma^{(k)} - \Sigma^{(0)} \right) \left( \sum_{1 \leq k \leq K} \frac{1}{K} \Sigma^{(k)} \right)^{-1} \right\|_1,$$

which could diverge with rate  $\sqrt{p}$ .

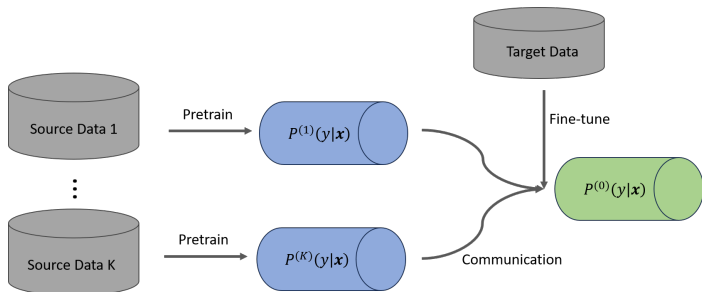
- 2 Covariate-shift Robust Transfer Learning
  - TransFusion: Transfer Learning with a Fused Regularization
  - D-TransFusion: TransFusion under Distributed Setting

# Distributed Transfer Learning in One-shot



**Distributed setting:** Source datasets are distributed across different machines.

- Privacy concern: raw data cannot be shared across machines;
- Communication bottleneck: inter-machine data communication is a significant source of latency;
- Pretraining & Fine-tuning strategy: quickly adapt to downstream tasks.





# Two-step D-TransFusion Estimator

**Step 1:** The  $k$ th source machine computes an initial estimator  $\tilde{\beta}^{(k)}$  using  $(\mathbf{X}^{(k)}, \mathbf{y}^{(k)})$  and sends it to the target machine. Then the target machine solves

$$\hat{\beta}_D \in \operatorname{argmin}_{\beta \in \mathbb{R}^{(K+1)p}} \left\{ \frac{n_S}{2N} \sum_{k=1}^K \|\tilde{\beta}^{(k)} - \beta^{(k)}\|_2^2 + \frac{1}{2N} \|\mathbf{y}^{(0)} - \mathbf{X}^{(0)}\beta^{(0)}\|_2^2 + \lambda_0 \left( \|\beta^{(0)}\|_1 + \sum_{k=1}^K \nu \|\tilde{\beta}^{(k)} - \beta^{(0)}\|_1 \right) \right\},$$

and obtains  $\hat{\beta}_{D\text{-TF1}}^{(0)} = \frac{n_S}{N} \sum_{k=1}^K \hat{\beta}_D^{(k)} + \frac{n_T}{N} \hat{\beta}_D^{(0)}$ .

# Two-step D-TransFusion Estimator

**Step 1:** The  $k$ th source machine computes an initial estimator  $\tilde{\beta}^{(k)}$  using  $(\mathbf{X}^{(k)}, \mathbf{y}^{(k)})$  and sends it to the target machine. Then the target machine solves

$$\hat{\beta}_D \in \operatorname{argmin}_{\beta \in \mathbb{R}^{(K+1)p}} \left\{ \frac{n_S}{2N} \sum_{k=1}^K \|\tilde{\beta}^{(k)} - \beta^{(k)}\|_2^2 + \frac{1}{2N} \|\mathbf{y}^{(0)} - \mathbf{X}^{(0)}\beta^{(0)}\|_2^2 + \lambda_0 \left( \|\beta^{(0)}\|_1 + \sum_{k=1}^K \nu \|\tilde{\beta}^{(k)} - \beta^{(0)}\|_1 \right) \right\},$$

and obtains  $\hat{\beta}_{\text{D-TF1}}^{(0)} = \frac{n_S}{N} \sum_{k=1}^K \hat{\beta}_D^{(k)} + \frac{n_T}{N} \hat{\beta}_D^{(0)}$ .

**Step 2:** The target node corrects  $\hat{\beta}_{\text{D-TF1}}^{(0)}$  on its local sample  $(\mathbf{X}^{(0)}, \mathbf{y}^{(0)})$  by

$$\hat{\delta}_D \in \operatorname{argmin}_{\delta \in \mathbb{R}^p} \left\{ \frac{1}{2n_T} \|\mathbf{y}^{(0)} - \mathbf{X}^{(0)}\hat{\beta}_{\text{D-TF1}}^{(0)} - \mathbf{X}^{(0)}\delta\|_2^2 + \tilde{\lambda} \|\delta\|_1 \right\},$$

and outputs the second-step estimator  $\hat{\beta}_{\text{D-TF2}}^{(0)} = \hat{\beta}_{\text{D-TF1}}^{(0)} + \hat{\delta}_D$ .

# Two-step D-TransFusion Estimator

**Step 1:** The  $k$ th source machine computes an initial estimator  $\tilde{\beta}^{(k)}$  using  $(\mathbf{X}^{(k)}, \mathbf{y}^{(k)})$  and sends it to the target machine. Then the target machine solves

$$\hat{\beta}_D \in \operatorname{argmin}_{\beta \in \mathbb{R}^{(K+1)p}} \left\{ \frac{n_S}{2N} \sum_{k=1}^K \|\tilde{\beta}^{(k)} - \beta^{(k)}\|_2^2 + \frac{1}{2N} \|\mathbf{y}^{(0)} - \mathbf{X}^{(0)}\beta^{(0)}\|_2^2 + \lambda_0 \left( \|\beta^{(0)}\|_1 + \sum_{k=1}^K \nu \|\tilde{\beta}^{(k)} - \beta^{(0)}\|_1 \right) \right\},$$

and obtains  $\hat{\beta}_{\text{D-TF1}}^{(0)} = \frac{n_S}{N} \sum_{k=1}^K \hat{\beta}_D^{(k)} + \frac{n_T}{N} \hat{\beta}_D^{(0)}$ .

**Step 2:** The target node corrects  $\hat{\beta}_{\text{D-TF1}}^{(0)}$  on its local sample  $(\mathbf{X}^{(0)}, \mathbf{y}^{(0)})$  by

$$\hat{\delta}_D \in \operatorname{argmin}_{\delta \in \mathbb{R}^p} \left\{ \frac{1}{2n_T} \left\| \mathbf{y}^{(0)} - \mathbf{X}^{(0)} \hat{\beta}_{\text{D-TF1}}^{(0)} - \mathbf{X}^{(0)} \delta \right\|_2^2 + \tilde{\lambda} \|\delta\|_1 \right\},$$

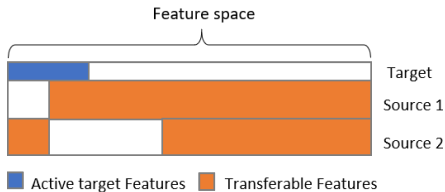
and outputs the second-step estimator  $\hat{\beta}_{\text{D-TF2}}^{(0)} = \hat{\beta}_{\text{D-TF1}}^{(0)} + \hat{\delta}_D$ .

- Only one-shot communication with the summary statistic is required.
- D-TransFusion has the same rate as TransFusion under mild conditions.

- 3 Adaptive Covariate-shift Robust Transfer Learning
  - AdaTrans: Feature-wise Adaptive Transfer Learning
  - Oracle AdaTrans and Theoretical Guarantee of AdaTrans

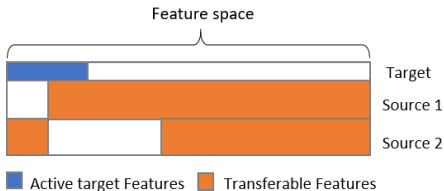
# Intuition of AdaTrans

On base of TransFusion, consider the **feature-specific transferable structure**:



# Intuition of AdaTrans

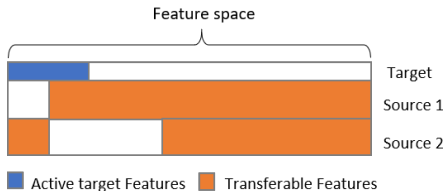
On base of TransFusion, consider the **feature-specific transferable structure**:



The transferability of the  $j$ th feature in the  $k$ -th source task can be assessed by the magnitude of model shift  $\delta_j^{(k)}$ . Ideally, we should...

# Intuition of AdaTrans

On base of TransFusion, consider the **feature-specific transferable structure**:



The transferability of the  $j$ th feature in the  $k$ -th source task can be assessed by the magnitude of model shift  $\delta_j^{(k)}$ . Ideally, we should...

- apply stronger penalties to transferable features with negligible  $\delta_j^{(k)}$ ;  
 → shrink  $\delta_j^{(k)}$  to 0, i.e. pool  $\beta_j^{(k)}$  and  $\beta_j^{(0)}$ , if the  $j$ -th feature from the  $k$ -th source is informative/transferable
- prevents excessive penalties to non-transferable features with large  $\delta_j^{(k)}$ .  
 → prevent introducing bias from non-transferable signals

# Key Idea of AdaTrans

Estimate  $\beta^{(0)}$  by solving

$$\operatorname{argmin}_{\beta \in \mathbb{R}^{(K+1)p}} \left\{ \frac{1}{2N} \sum_{k=0}^K \|\mathbf{y}^{(k)} - \mathbf{X}^{(k)}(\beta^{(0)} + \delta^{(k)})\|_2^2 + \sum_{j=1}^p \hat{w}_{0j} |\beta_j^{(0)}| + \sum_{k=1}^K \sum_{j=1}^p \hat{w}_{kj} |\delta_j^{(k)}| \right\},$$



# Key Idea of AdaTrans

Estimate  $\beta^{(0)}$  by solving

$$\operatorname{argmin}_{\beta \in \mathbb{R}^{(K+1)p}} \left\{ \frac{1}{2N} \sum_{k=0}^K \|\mathbf{y}^{(k)} - \mathbf{X}^{(k)}(\beta^{(0)} + \delta^{(k)})\|_2^2 + \sum_{j=1}^p \hat{w}_{0j} |\beta_j^{(0)}| + \sum_{k=1}^K \sum_{j=1}^p \hat{w}_{kj} |\delta_j^{(k)}| \right\},$$

Choice of weight?

# Key Idea of AdaTrans

Estimate  $\beta^{(0)}$  by solving

$$\operatorname{argmin}_{\beta \in \mathbb{R}^{(K+1)p}} \left\{ \frac{1}{2N} \sum_{k=0}^K \|\mathbf{y}^{(k)} - \mathbf{X}^{(k)}(\beta^{(0)} + \delta^{(k)})\|_2^2 + \sum_{j=1}^p \hat{w}_{0j} |\beta_j^{(0)}| + \sum_{k=1}^K \sum_{j=1}^p \hat{w}_{kj} |\delta_j^{(k)}| \right\},$$

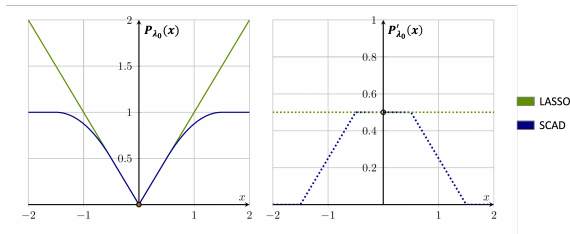
**Choice of weight?**  $\Rightarrow$  **Folded-concave penalty function**  $\mathcal{P}_{\lambda_0}(\cdot)$ :

# Key Idea of AdaTrans

Estimate  $\beta^{(0)}$  by solving

$$\operatorname{argmin}_{\beta \in \mathbb{R}^{(K+1)p}} \left\{ \frac{1}{2N} \sum_{k=0}^K \|\mathbf{y}^{(k)} - \mathbf{X}^{(k)}(\beta^{(0)} + \delta^{(k)})\|_2^2 + \sum_{j=1}^p \hat{w}_{0j} |\beta_j^{(0)}| + \sum_{k=1}^K \sum_{j=1}^p \hat{w}_{kj} |\delta_j^{(k)}| \right\},$$

**Choice of weight?**  $\Rightarrow$  **Folded-concave penalty function**  $\mathcal{P}_{\lambda_0}(\cdot)$ :

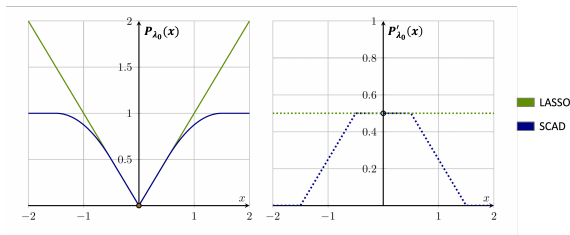


# Key Idea of AdaTrans

Estimate  $\beta^{(0)}$  by solving

$$\operatorname{argmin}_{\beta \in \mathbb{R}^{(K+1)p}} \left\{ \frac{1}{2N} \sum_{k=0}^K \|\mathbf{y}^{(k)} - \mathbf{X}^{(k)}(\beta^{(0)} + \delta^{(k)})\|_2^2 + \sum_{j=1}^p \hat{w}_{0j} |\beta_j^{(0)}| + \sum_{k=1}^K \sum_{j=1}^p \hat{w}_{kj} |\delta_j^{(k)}| \right\},$$

**Choice of weight?**  $\Rightarrow$  **Folded-concave penalty function**  $\mathcal{P}_{\lambda_0}(\cdot)$ :

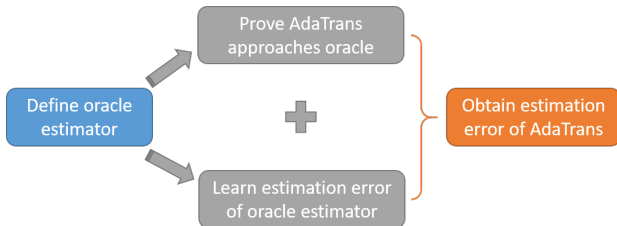


Borrowing the idea of local linear approximation, take  $\hat{w}_{0j} \propto \mathcal{P}'_{\lambda_0}(\hat{\beta}_{\text{init},j}^{(0)})$  and  $\hat{w}_{kj} \propto \mathcal{P}'_{\lambda_0}(\hat{\delta}_{\text{init},j}^{(k)})$ , where  $\hat{\beta}_{\text{init},j}^{(0)}$  and  $\hat{\delta}_{\text{init},j}^{(k)}$  are initial estimators of  $\beta_j^{(0)}$  and  $\delta_j$ .

- 3 Adaptive Covariate-shift Robust Transfer Learning
  - AdaTrans: Feature-wise Adaptive Transfer Learning
  - Oracle AdaTrans and Theoretical Guarantee of AdaTrans

# Developing Theory for AdaTrans

Recall that under certain conditions, folded-concave penalization can obtain an **oracle estimator**, where the sparsity and transferable structures are known.



# Oracle AdaTrans Estimator



How to define oracle estimator for AdaTrans?

# Oracle AdaTrans Estimator

How to define oracle estimator for AdaTrans?

1 Define **sparsity structure**:

- Active target feature set:  $S_0 = \{j : \beta_j^{(0)} \neq 0\}$ ,
- Inactive target feature set:  $S_0^c = \{j : \beta_j^{(0)} = 0\}$ ;



# Oracle AdaTrans Estimator

How to define oracle estimator for AdaTrans?

1 Define **sparsity structure**:

- Active target feature set:  $S_0 = \{j : \beta_j^{(0)} \neq 0\}$ ,
- Inactive target feature set:  $S_0^c = \{j : \beta_j^{(0)} = 0\}$ ;

2 Define **transferability structure**:

- **Non-transferable** set:  $S_k = \{j : \delta_j^{(k)} \neq 0\}$ ,  $k = 1, \dots, K$ ,
- Transferable set:  $S_k^c = \{j : \delta_j^{(k)} = 0\}$ ,  $k = 1, \dots, K$ .

# Oracle AdaTrans Estimator

How to define oracle estimator for AdaTrans?

**1** Define **sparsity structure**:

- Active target feature set:  $S_0 = \{j : \beta_j^{(0)} \neq 0\}$ ,
- Inactive target feature set:  $S_0^c = \{j : \beta_j^{(0)} = 0\}$ ;

**2** Define **transferability structure**:

- **Non-transferable** set:  $S_k = \{j : \delta_j^{(k)} \neq 0\}$ ,  $k = 1, \dots, K$ ,
- Transferable set:  $S_k^c = \{j : \delta_j^{(k)} = 0\}$ ,  $k = 1, \dots, K$ .

**3** Define **oracle AdaTrans estimator**  $\hat{\beta}_{\text{ora}}^{(0)}, \hat{\delta}_{\text{ora}}^{(1)}, \dots, \hat{\delta}_{\text{ora}}^{(K)}$  via

$$\begin{aligned} \min_{\beta^{(0)}, \{\delta^{(k)}\}_{k=1}^K} & \frac{1}{N} \sum_{k=0}^K \|\mathbf{y}^{(k)} - \mathbf{X}^{(k)}(\beta^{(0)} + \delta^{(k)})\|_2^2 \\ \text{s.t.} & \beta_{S_0^c}^{(0)} = 0, \delta_{S_k^c}^{(k)} = 0, \forall k = 1, \dots, K. \end{aligned} \quad (1)$$

# Oracle AdaTrans Estimator

## Theorem (Solution of oracle AdaTrans estimator)

If  $|S_0| < n_T$  and  $\max_{1 \leq k \leq K} |S_k| < n_S$ , the solution to problem (1) satisfies

$$\hat{\beta}_{\text{ora}, S_0}^{(0)} = [\tilde{\mathbf{X}}_{S_0}^\top \tilde{\mathbf{X}}_{S_0}]^{-1} \tilde{\mathbf{X}}_{S_0}^\top \mathbf{y} \quad \text{and} \quad \hat{\beta}_{\text{ora}, S_0^c}^{(0)} = \mathbf{0}. \quad (2)$$

- $\tilde{\mathbf{X}}_{S_0} = ((\mathbf{X}_{S_0}^{(0)})^\top, (\tilde{\mathbf{X}}_{S_0}^{(1)})^\top, \dots, (\tilde{\mathbf{X}}_{S_0}^{(K)})^\top)^\top$ .
- $\tilde{\mathbf{X}}_{S_0}^{(k)} = (\mathbf{I} - \mathbf{H}_{S_k}^{(k)}) \mathbf{X}_{S_0}^{(k)}$ , where  $\mathbf{H}_{S_k}^{(k)} := \mathbf{X}_{S_k}^{(k)} [(\mathbf{X}_{S_k}^{(k)})^\top \mathbf{X}_{S_k}^{(k)}]^{-1} (\mathbf{X}_{S_k}^{(k)})^\top$ .

# Oracle AdaTrans Estimator

## Theorem (Solution of oracle AdaTrans estimator)

If  $|S_0| < n_T$  and  $\max_{1 \leq k \leq K} |S_k| < n_S$ , the solution to problem (1) satisfies

$$\hat{\beta}_{\text{ora}, S_0}^{(0)} = [\tilde{\mathbf{X}}_{S_0}^\top \tilde{\mathbf{X}}_{S_0}]^{-1} \tilde{\mathbf{X}}_{S_0}^\top \mathbf{y} \quad \text{and} \quad \hat{\beta}_{\text{ora}, S_0^c}^{(0)} = \mathbf{0}. \quad (2)$$

- $\tilde{\mathbf{X}}_{S_0} = ((\mathbf{X}_{S_0}^{(0)})^\top, (\tilde{\mathbf{X}}_{S_0}^{(1)})^\top, \dots, (\tilde{\mathbf{X}}_{S_0}^{(K)})^\top)^\top$ .
- $\tilde{\mathbf{X}}_{S_0}^{(k)} = (\mathbf{I} - \mathbf{H}_{S_k}^{(k)}) \mathbf{X}_{S_0}^{(k)}$ , where  $\mathbf{H}_{S_k}^{(k)} := \mathbf{X}_{S_k}^{(k)} [(\mathbf{X}_{S_k}^{(k)})^\top \mathbf{X}_{S_k}^{(k)}]^{-1} (\mathbf{X}_{S_k}^{(k)})^\top$ .
- $\tilde{\mathbf{X}}_{S_0}^{(k)}$  is indeed the projection of the **active target feature** onto the **null space** of the **non-transferable** feature in the  $k$ -th source sample.

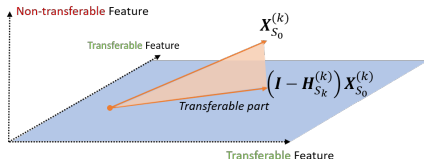
# Oracle AdaTrans Estimator

## Theorem (Solution of oracle AdaTrans estimator)

If  $|S_0| < n_T$  and  $\max_{1 \leq k \leq K} |S_k| < n_S$ , the solution to problem (1) satisfies

$$\hat{\beta}_{\text{ora}, S_0}^{(0)} = [\tilde{\mathbf{X}}_{S_0}^\top \tilde{\mathbf{X}}_{S_0}]^{-1} \tilde{\mathbf{X}}_{S_0}^\top \mathbf{y} \quad \text{and} \quad \hat{\beta}_{\text{ora}, S_0^c}^{(0)} = \mathbf{0}. \quad (2)$$

- $\tilde{\mathbf{X}}_{S_0} = ((\mathbf{X}_{S_0}^{(0)})^\top, (\tilde{\mathbf{X}}_{S_0}^{(1)})^\top, \dots, (\tilde{\mathbf{X}}_{S_0}^{(K)})^\top)^\top$ .
- $\tilde{\mathbf{X}}_{S_0}^{(k)} = (\mathbf{I} - \mathbf{H}_{S_k}^{(k)}) \mathbf{X}_{S_0}^{(k)}$ , where  $\mathbf{H}_{S_k}^{(k)} := \mathbf{X}_{S_k}^{(k)} [(\mathbf{X}_{S_k}^{(k)})^\top \mathbf{X}_{S_k}^{(k)}]^{-1} (\mathbf{X}_{S_k}^{(k)})^\top$ .
- $\tilde{\mathbf{X}}_{S_0}^{(k)}$  is indeed the projection of the **active target feature** onto the **null space** of the **non-transferable** feature in the  $k$ -th source sample.



# Estimation Error of Oracle AdaTrans

## Theorem (Estimation error of oracle AdaTrans)

If  $|S_0| < n_T$ ,  $\max_{1 \leq k \leq K} |S_k| < n_S$  and  $N \geq \log p$ , the error of  $\hat{\beta}_{\text{ora}}^{(0)}$  satisfies

$$\|\hat{\beta}_{\text{ora}}^{(0)} - \beta^{(0)}\|_2 \lesssim \kappa_F \left\| \left( \frac{\mathbf{X}_{S_0}^\top \mathbf{X}_{S_0}}{N} \right)^{-1} \right\|_\infty \sqrt{\frac{s \log s}{N}}, \quad (3)$$

with probability larger than  $1 - \exp(-c_1 \log p)$ , where  $\mathbf{X}_{S_0}$  is column-submatrix indexed by  $S_0$  of the full-sample design matrix  $\mathbf{X}$ , and

$$\kappa_F := \frac{\left\| [\tilde{\mathbf{X}}_{S_0}^\top \tilde{\mathbf{X}}_{S_0}]^{-1} \tilde{\mathbf{X}}_{S_0}^\top \boldsymbol{\epsilon} \right\|_\infty}{\left\| [\mathbf{X}_{S_0}^\top \mathbf{X}_{S_0}]^{-1} \mathbf{X}_{S_0}^\top \boldsymbol{\epsilon} \right\|_\infty}.$$

# Estimation Error of Oracle AdaTrans

## Theorem (Estimation error of oracle AdaTrans)

If  $|S_0| < n_T$ ,  $\max_{1 \leq k \leq K} |S_k| < n_S$  and  $N \geq \log p$ , the error of  $\hat{\beta}_{\text{ora}}^{(0)}$  satisfies

$$\|\hat{\beta}_{\text{ora}}^{(0)} - \beta^{(0)}\|_2 \lesssim \kappa_F \left\| \left( \frac{\mathbf{X}_{S_0}^\top \mathbf{X}_{S_0}}{N} \right)^{-1} \right\|_\infty \sqrt{\frac{s \log s}{N}}, \quad (3)$$

with probability larger than  $1 - \exp(-c_1 \log p)$ , where  $\mathbf{X}_{S_0}$  is column-submatrix indexed by  $S_0$  of the full-sample design matrix  $\mathbf{X}$ , and

$$\kappa_F := \frac{\left\| [\tilde{\mathbf{X}}_{S_0}^\top \tilde{\mathbf{X}}_{S_0}]^{-1} \tilde{\mathbf{X}}_{S_0}^\top \epsilon \right\|_\infty}{\left\| [\mathbf{X}_{S_0}^\top \mathbf{X}_{S_0}]^{-1} \mathbf{X}_{S_0}^\top \epsilon \right\|_\infty}.$$

- $\kappa_F$  measures the transferability of source datasets. For  $k = 1, \dots, K$ ,
  - if  $\mathbf{X}_{S_k}^{(k)} \perp \mathbf{X}_{S_0}^{(k)}$ , all active features are transferable, then  $\kappa_F = 1$ ;
  - if  $S_0 \subset S_k$ , all active features are non-transferable, then  $\kappa_F \asymp \sqrt{N/n_T}$ , and the final rate becomes  $\sqrt{s \log s / n_T}$ .

## Theorem (Oracle property of AdaTrans)

Consider the parametric space

$$\Theta_1 = \left\{ \left\| \delta_{S_k}^{(k)} \right\|_{\min} \geq h_k^\wedge, \left\| \delta_{S_k^c}^{(k)} \right\|_{\max} = 0, k = 1, \dots, K; \left\| \beta_{S_0}^{(0)} \right\|_{\min} \geq h_0^\wedge, \left\| \beta_{S_0^c}^{(0)} \right\|_{\max} = 0 \right\}.$$

Suppose for some  $a > a_2 \geq 0$ , the initial estimators satisfy

$$\left\| \hat{\beta}_{init}^{(0)} - \beta^{(0)} \right\|_{\infty} \leq \frac{a_2}{2} \lambda_0, \quad \left\| \hat{\delta}_{init}^{(k)} - \delta^{(k)} \right\|_{\infty} \leq \frac{a_2}{2} \lambda_1;$$

the minimal target signal  $h_0^\wedge \geq a \lambda_0 \gtrsim \sqrt{\frac{\log p}{N}}$ , and the non-transferable signal  $h_k^\wedge \geq a \lambda_1 \gtrsim \sqrt{\frac{n_S \log p}{N}}$ , and  $n_S \gtrsim \log p$ . Then by choosing  $w_{0j} = \mathcal{P}'_{\lambda_0}(\hat{\beta}_{init,j}^{(0)})/\lambda_0$  and  $w_{kj} = \mathcal{P}'_{\lambda_1}(\hat{\delta}_{init,j}^{(k)})/\lambda_1$ , with probability larger than  $1 - \exp(-c_1 \log p)$ , we obtain the oracle AdaTrans.



- 4 Numerical Studies
  - Simulation Examples for TransFusion
  - Simulation Examples for AdaTrans

# Simulation settings for TransFusion

Recall the regression models

$$\mathbf{y}^{(0)} = \mathbf{X}^{(0)}\boldsymbol{\beta}^{(0)} + \boldsymbol{\epsilon}^{(0)}, \mathbf{y}^{(k)} = \mathbf{X}^{(k)}(\boldsymbol{\beta}^{(0)} + \boldsymbol{\delta}^{(k)}) + \boldsymbol{\epsilon}^{(k)}, k = 1, \dots, K.$$

## General setup:

- Target task:  $n_T = 150$ ,  $s = 10$ ,  $\boldsymbol{\beta}^{(0)} = (\mathbf{1}_s^\top, \mathbf{0}_{p-s}^\top)^\top$ ,  $\epsilon_i^{(0)} \sim N(0, 1)$ .
- Source task:  $n_S = 200$ ,  $K \in \{1, 3, 5, 7, 9\}$ ,  $\epsilon_i^{(k)} \sim N(0, 1)$ .

## Model shift:

- $\delta_j^{(k)} \sim N(0.1, 0.2^2)$  for  $1 \leq j \leq 50$  and  $\delta_j^{(k)} = 0$  otherwise.

## Covariate shift:

- Homogeneous design (without covariate shift): Each  $\mathbf{X}_i^{(k)} \sim N(0, \mathbf{I})$ .
- Heterogeneous design (with covariate shift): Each  $\mathbf{X}_i^{(k)} \sim N(0, \boldsymbol{\Sigma}^{(k)})$ , with  $\boldsymbol{\Sigma}^{(k)} = (\mathbf{A}^{(k)})^\top (\mathbf{A}^{(k)}) + \mathbf{I}$ , where  $\mathbf{A}^{(k)}$  is a random matrix with each entry equals 0.3 with probability 0.3 and equals 0 with probability 0.7.

# Methods to be compared

- Lasso (baseline): LASSO regression on the target task.
- TransLasso (first-step) (Li et al., 2022): pooled estimator.
- TransLasso (two-step) (Li et al., 2022): debiased estimator.
- TransHDGLM (Li et al., 2023).
- TransFusion (first-step): the first step TransFusion estimator  $\hat{\beta}_{TF1}^{(0)}$ .
- TransFusion (two-step): the debiased TransFusion estimator  $\hat{\beta}_{TF2}^{(0)}$ .

# Simulation Results: TransFusion

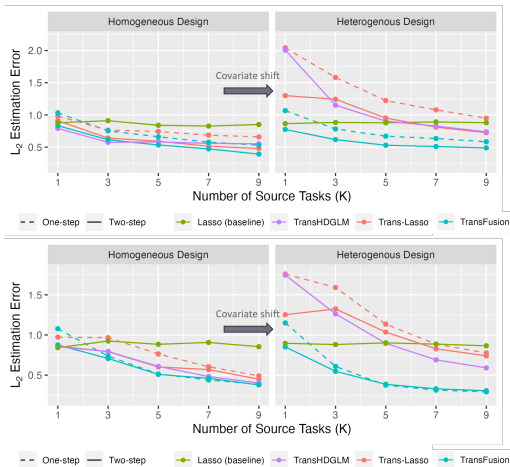


Figure: Estimation errors with/without covariate shift. Upper panel: task diversity  $\epsilon_D \neq 0$ ; lower panel:  $\epsilon_D = 0$ .

# Simulation Results: D-TransFusion

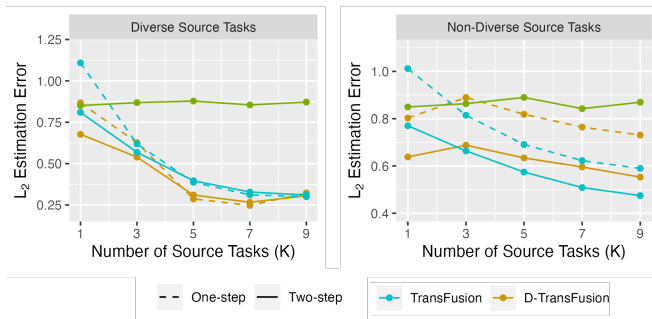


Figure: Estimation errors with  $\epsilon_D = 0$  (left panel) and  $\epsilon_D \neq 0$  (right panel).

- 4 Numerical Studies
  - Simulation Examples for TransFusion
  - Simulation Examples for AdaTrans

# Simulation Settings for AdaTrans

Recall the regression models

$$\mathbf{y}^{(0)} = \mathbf{X}^{(0)}\boldsymbol{\beta}^{(0)} + \boldsymbol{\epsilon}^{(0)}, \mathbf{y}^{(k)} = \mathbf{X}^{(k)}(\boldsymbol{\beta}^{(0)} + \boldsymbol{\delta}^{(k)}) + \boldsymbol{\epsilon}^{(k)}, k = 1, \dots, K.$$

## General setup:

- Target task:  $n_T = 50$ ,  $s = 8$ ,  $\boldsymbol{\beta}^{(0)} = (\mathbf{1}_s^\top, \mathbf{0}_{p-s}^\top)^\top$ ,  $\epsilon_i^{(0)} \sim N(0, 1)$ .
- Source task:  $n_S = 200$ ,  $K = 2$ ,  $\epsilon_i^{(k)} \sim N(0, 1)$ .

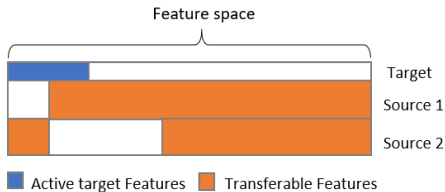
**Covariate shift:** Same as TransFusion.

# Simulation Settings for AdaTrans

## Model shift:

We generate two source samples with non-overlapping transferable features:

- First source: the non-transferable  $\delta^{(k)}$  is nonzero for the first  $s/2$  elements;
- Second source: the non-transferable  $\delta^{(k)}$  is nonzero from  $(s/2 + 1)$ -th to 25th elements.





# Methods to be compared

- Lasso (baseline): LASSO regression on the target task.
- TransGLM (Tian and Feng, 2022): TransLasso with source detection.
- AdaTrans: AdaTrans estimator.
- Oracle AdaTrans: Oracle AdaTrans estimator.

# Simulation Results: AdaTrans

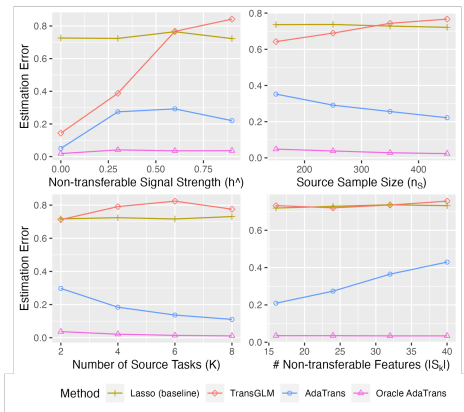


Figure: Estimation errors of different transfer learning methods.

- AdaTrans can also auto-detect and filter out non-transferable features.

We proposed a new transfer learning framework that is robust to covariate shift and adaptive to feature-specific transferable structure.

- **TransFusion:** Conducting a fused-regularization based “joint training + debiasing” to achieve covariate-shift robustness.
- **D-TransFusion:** Incorporating intermediate estimators from different machines into TransFusion with one-shot communication.
- **AdaTrans:** Utilizing folded-concave penalization to auto-detect transferable structure while estimating parameters.
- Non-asymptotic bounds of estimation errors for all proposed estimators are established.

*Thank You!*