Introduction



Lack of theoretical understanding

- The representations may NOT be the same! But most of theoretical studies impose such an assumption.
- As the number of tasks grow, there can be outlier tasks or adversarial attacks.

Multi-task Learning

Problem setting

- T tasks, data $\{\boldsymbol{x}_i^{(t)}, y_i^{(t)}\}_{i=1}^n$ from the t-th task, with $\boldsymbol{x}_i^{(t)} \in \mathbb{R}^p, y_i^{(t)} \in \mathbb{R};$
- There exists an unknown subset $S \subseteq [T]$, such that for all $t \in S$,

$$y_i^{(t)} = (\boldsymbol{x}_i^{(t)})^T \boldsymbol{\beta}^{(t)*} + \boldsymbol{\epsilon}_i^{(t)}, \quad i = 1:n,$$

where $\boldsymbol{\beta}^{(t)*} = \boldsymbol{A}^{(t)*}\boldsymbol{\theta}^{(t)*}, \ \boldsymbol{A}^{(t)*} \in \mathcal{O}^{p \times r} = \{\boldsymbol{A} \in \mathbb{R}^{p \times r} : \boldsymbol{A}^T \boldsymbol{A} = \boldsymbol{I}_r\},\$ $\boldsymbol{\theta}^{(t)*} \in \mathbb{R}^r, r \leq p$, and $\{\epsilon_i^{(t)}\}_{i=1}^n$ are i.i.d. zero-mean sub-Gaussian $\{\boldsymbol{x}_i^{(t)}\}_{i=1}^n$;

- $\min_{\overline{A} \in \mathcal{O}^{p \times r}} \max_{t \in S} \|A^{(t)*}(A^{(t)*})^T \overline{A}(\overline{A})^T\|_2 \le h;$
- Data $\{\{\boldsymbol{x}_i^{(t)}, y_i^{(t)}\}_{i=1}^n\}_{t \notin S}$ from outlier tasks ~ an arbitrary distribution \mathbb{Q}_{S^c} ;

 $oldsymbol{eta}^{(5)*}$

• Goal: jointly learning $\{\beta^{(t)*}\}_{t\in S}$

 $m{eta}^{(3)*} = m{A}^{(3)*} m{ heta}^{(3)*}$

 $oldsymbol{eta}^{(4)*} = oldsymbol{A}^{(4)*}oldsymbol{ heta}^{(4)*}$

 $\beta^{(2)*} = A^{(1)*} \theta^{(2)*}$

 $\beta^{(1)*} = A^{(1)*} \theta^{(1)*}$

LEARNING FROM SIMILAR LINEAR REPRESENTATIONS: ADAPTIVITY, MINIMAXITY, AND ROBUSTNESS

Ye Tian[†], Yuqi Gu[†], and Yang Feng[‡]

[†]Department of Statistics, Columbia University [‡]Department of Biostatistics, School of Global Public Health, New York University

> **Two-step algorithm:** (*r* known) Set $\lambda \approx \sqrt{r(p + \log T)}$ and $\gamma \approx \sqrt{p + \log T}$ • $\{\widehat{A}^{(t)}\}_{t=1}^{T}, \{\widehat{\theta}^{(t)}\}_{t=1}^{T}, \widehat{\overline{A}} \in \underset{\{A^{(t)}\}_{t=1}^{T}, \{\theta^{(t)}\}_{t=1}^{T}, \overline{A}}{\operatorname{arg\,min}}$ $\{\sum_{t=1}^{T} \frac{1}{n} \sum_{i=1}^{n} [y_{i}^{(t)} - (\mathbf{x}_{i}^{(t)})^{T} \mathbf{A}^{(t)} \mathbf{\theta}^{(t)}]^{2} + \frac{\lambda}{\sqrt{n}} \|\mathbf{A}^{(t)} (\mathbf{A}^{(t)})^{T} - \overline{\mathbf{A}} (\overline{\mathbf{A}})^{T} \|_{2}\};$ • $\widehat{\boldsymbol{\beta}}^{(t)} = \underset{\boldsymbol{\beta} \in \mathbb{R}^p}{\operatorname{arg\,min}} \left\{ \frac{1}{n} \sum_{i=1}^n [y_i^{(t)} - (\boldsymbol{x}_i^{(t)})^T \boldsymbol{\beta}]^2 + \frac{\gamma}{\sqrt{n}} \|\boldsymbol{\beta} - \widehat{\boldsymbol{A}}^{(t)} \widehat{\boldsymbol{\theta}}^{(t)}\|_2 \right\} \text{ for } t \in [T]$

Assumptions

Notations: Coefficient matrix $\boldsymbol{B}_{S}^{*} \in \mathbb{R}^{p \times |S|}$, each column of which is a coefficient vector in $\{\boldsymbol{\beta}^{(t)*}\}_{t\in S}$. $\boldsymbol{\Sigma}^{(t)} \coloneqq \mathbb{E}[\tilde{\boldsymbol{x}}^{(t)}(\boldsymbol{x}^{(t)})^T]$.

A.1 $\{\boldsymbol{x}_i^{(t)}\}_{i=1}^n$ are i.i.d. sub-Gaussian, and $0 < c \leq \lambda_{\min}(\boldsymbol{\Sigma}^{(t)}) \leq \lambda_{\max}(\boldsymbol{\Sigma}^{(t)}) \leq C$. **A.2** (Task diversity) $\max_{t \in S} \|\boldsymbol{\theta}^{(t)*}\|_2 \leq C < \infty$, and $\sigma_r(\boldsymbol{B}_S^*/\sqrt{T}) \geq \frac{c}{\sqrt{r}}$. **A.3** (Not many outlier tasks) $\frac{|S^c|}{T}r^{3/2} \leq a$ small constant c.

Upper bounds: When $n \gtrsim p + \log T$, $\forall S \subseteq [T]$ and an arbitrary \mathbb{Q}_{S^c} , w.h.p.,

$$\max_{t \in S} \|\widehat{\boldsymbol{\beta}}^{(t)} - \boldsymbol{\beta}^{(t)*}\|_2 \lesssim \left(r\sqrt{\frac{p}{nT}} + \sqrt{rh} + \sqrt{r}\sqrt{\frac{r + \log T}{n}} + \sqrt{\frac{p}{n}} \cdot \frac{|S^c|}{T} r^{3/2} \right)$$

If outlier tasks in S^c also follow linear model (1), when

tasks III D also follow IIIear IIIodel (1), w.II.p.,

$$\max_{t \in [T]} \|\widehat{\boldsymbol{\beta}}^{(t)} - \boldsymbol{\beta}^{(t)*}\|_2 \lesssim \sqrt{\frac{p + \log T}{n}}$$

(1)

Lower bounds:
$$\forall \{\widehat{\boldsymbol{\beta}}^{(t)}\}_{t=1}^{T}, \exists S \subseteq [T], \{\boldsymbol{\beta}^{(t)}\}_{t\in S}, \mathbb{Q}_{S^{c}}, \text{w.p.} \geq 1/10,$$

$$\max_{t\in S} \|\widehat{\boldsymbol{\beta}}^{(t)} - \boldsymbol{\beta}^{(t)*}\|_{2} \gtrsim \sqrt{\frac{pr}{nT}} + h \wedge \sqrt{\frac{p+\log T}{n}} + \sqrt{\frac{r+\log T}{n}} + \frac{\epsilon n}{\sqrt{nT}}$$

If outlier tasks in S^c also follow linear model (1), $\forall \{\widehat{\boldsymbol{\beta}}^{(t)}\}_{t=1}^T$, $\exists \{\boldsymbol{\beta}^{(t)}\}_{t=1}^T$, w.p. $\geq 1/10,$

$$\max_{t \in [T]} \|\widehat{\boldsymbol{\beta}}^{(t)} - \boldsymbol{\beta}^{(t)*}\|_2 \gtrsim \sqrt{\frac{p + \log T}{n}}.$$

Transferring to New Tasks

Problem setting:

- Data $\{(\boldsymbol{x}_{i}^{(0)}, y_{i}^{(0)})\}_{i=1}^{n_{0}}$ from a new task, generated from model (1);
- $\max_{t \in S} \| \mathbf{A}^{(t)*} (\mathbf{A}^{(t)*})^T \mathbf{A}^{(0)*} (\mathbf{A}^{(0)*})^T \|_2 \le h;$
- Goal: learning $\beta^{(0)*}$

Two-step algorithm: (*r* known) Take \overline{A} from MTL algorithm, $\gamma \simeq \sqrt{p + \log T}$

- $\widehat{\boldsymbol{\theta}}^{(0)} \in \operatorname*{arg\,min}_{\boldsymbol{\theta} \in \mathbb{R}^r} \left\{ \frac{1}{n_0} \sum_{i=1}^{n_0} [y_i^{(0)} (\boldsymbol{x}_i^{(0)})^T \widehat{\overline{\boldsymbol{A}}} \boldsymbol{\theta}]^2 \right\};$
- $\widehat{\boldsymbol{\beta}}^{(0)} = \operatorname*{arg\,min}_{\boldsymbol{\beta} \in \mathbb{R}^p} \left\{ \frac{1}{n_0} \sum_{i=1}^{n_0} [y_i^{(0)} (\boldsymbol{x}_i^{(0)})^T \boldsymbol{\beta}]^2 + \frac{\gamma}{\sqrt{n_0}} \|\boldsymbol{\beta} \widehat{\overline{\boldsymbol{A}}} \widehat{\boldsymbol{\theta}}^{(0)}\|_2 \right\} \text{ for } t \in [T]$

Assumptions

A.4 $\{\boldsymbol{x}_i^{(0)}\}_{i=1}^{n_0}$ are i.i.d. sub-Gaussian, and $0 < c \leq \lambda_{\min}(\boldsymbol{\Sigma}^{(0)}) \leq \lambda_{\max}(\boldsymbol{\Sigma}^{(0)}) \leq C$. **A.5** $\|\theta^{(0)*}\|_2 \leq C$.

Upper bound: When $n, n_0 \gtrsim p + \log T, \forall S \subseteq [T]$ and arbitrary \mathbb{Q}_{S^c} , w.h.p.,



$$\|\widehat{\boldsymbol{\beta}}^{(0)} - \boldsymbol{\beta}^{(0)*}\|_{2} \lesssim \left(r\sqrt{\frac{p}{nT}} + \sqrt{rh} + \sqrt{r}\sqrt{\frac{r+\log T}{n}} + \sqrt{\frac{p}{n}} \cdot \frac{|S^{c}|}{T}r\right)$$
Lower bound: $\forall \widehat{\boldsymbol{\beta}}^{(0)}, \exists S \subseteq [T], \{\boldsymbol{\beta}^{(t)}\}_{t \in \{0\} \cup S}, \mathbb{Q}_{S^{c}}, \text{w.p.} \geq 1/10,$

$$2$$
;

Adaptation to Unknown Intrinsic Dimension

Intuition:

- It suffices to estimate r well when h and the proportion of outlier task $|S^c|/T$ are small;
- $\sigma_r(\boldsymbol{B}_S^*/\sqrt{T}) \gtrsim 1/\sqrt{r}, \ \sigma_{r+1}(\boldsymbol{B}_S^*/\sqrt{T}) \lesssim h \lesssim \sqrt{\frac{p+\log T}{nr}} \lesssim 1/\sqrt{r}$
- Do a thresholding to estimate r



MTL Simulations

No outlier tasks: n = 100, T = 6, p = 20, r = 3



🕨 RL-MTL-oracle 🔶 RL-MTL-adaptive 🔶 RL-MTL-naive 🔶 Single-task linear regressior

Estimation error $\max_{t \in S} \|\widehat{\boldsymbol{\beta}}^{(t)} - \boldsymbol{\beta}^{(t)*}\|_2$





(a) Estimation error $\max_{t \in S} \|\widehat{\beta}^{(t)} - \beta^{(t)*}\|_2$

(b) Estimation error $\max_{t \in S^c} \|\widehat{\beta}^{(t)} - \beta^{(t)*}\|_2$









