# **Near-Optimal Model Discrimination**

arxiv.org/abs/2012.02901

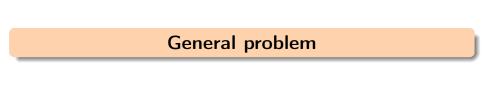
Dmitrii M. Ostrovskii Mohamed Ndaoud Adel Javanmard Meisam Razaviyayn

University of Southern California

USC Probability and Statistics Seminar February 5, 2021

# Outline

- General problem formulation
- Linear models
- Extensions



### Model discrimination task

- Let  $z \in \mathcal{Z}$  be a random observation distributed according to  $\mathbb{P}_0$  or  $\mathbb{P}_1$ .
- Let  $\theta_0, \theta_1 \in \mathbb{R}^d$  be the **best-fit models** of z according to  $\mathbb{P}_0, \mathbb{P}_1$ , i.e.

$$\theta_k = \underset{\theta \in \mathbb{R}^d}{\operatorname{argmin}} \left\{ L_k(\theta) := \mathbb{E}_{z \sim \mathbb{P}_k} \, \ell_z(\theta) \right\},$$

with strictly convex loss  $\ell_z: \mathbb{R}^d \to \mathbb{R}$  and population risks  $L_0(\cdot), L_1(\cdot)$ .

• Statistician has access to  $\theta^* \in \{\theta_0, \theta_1\}$  (but not to  $\bar{\theta} \in \{\theta_0, \theta_1\} \setminus \theta^*$ ) knows  $\ell_z$ , and observes two i.i.d. samples:

$$Z^0 = (z_1^0,...,z_n^0) \sim \mathbb{P}_0^{\otimes n}, \quad Z^1 = (z_1^1,...,z_n^1) \sim \mathbb{P}_1^{\otimes n}.$$

• Task: distinguish between the two hypotheses

$$\mathcal{H}_0: \{\theta^* = \theta_0\}, \quad \mathcal{H}_1: \{\theta^* = \theta_1\}.$$

### Model discrimination task

- Classical setup: both  $\theta_0, \theta_1$  known; one sample  $Z \sim \mathbb{P}_{\theta}^{\otimes n}$  observed.

  Which  $\theta \in \{\theta_0, \theta_1\}$  corresponds to the sample?

  Two simple hypotheses about  $\theta$ .
- Our setup: we observe both samples but only one model  $\theta^* \in \{\theta_0, \theta_1\}$ .

  Which  $Z \in \{Z^0, Z^1\}$  corresponds to  $\theta^*$ ?

  Two composite hypotheses about  $(\theta_0, \theta_1)$ .
- Statistician has access to  $\theta^* \in \{\theta_0, \theta_1\}$  (but not to  $\bar{\theta} \in \{\theta_0, \theta_1\} \setminus \theta^*$ ) knows  $\ell_z$ , and observes two i.i.d. samples:

$$Z^0 = (z_1^0,...,z_n^0) \sim \mathbb{P}_0^{\otimes n}, \quad Z^1 = (z_1^1,...,z_n^1) \sim \mathbb{P}_1^{\otimes n}.$$

• Task: distinguish between the two hypotheses about  $(\theta_0, \theta_1) \in \mathbb{R}^{2d}$ :

$$\mathcal{H}_0: (\theta_0,\theta_1) \in \{\theta^*\} \times \bar{\Theta}_0 \ \text{ vs. } \mathcal{H}_1: (\theta_0,\theta_1) \in \bar{\Theta}_1 \times \{\theta^*\} \text{ for some } \bar{\Theta}_0,\bar{\Theta}_1.$$

# Separation and sample complexity

$$\mathcal{H}_0: (\theta_0,\theta_1) \in ({\color{blue}\theta^*},\bar{\Theta}_0) \text{ vs. } \mathcal{H}_1: (\theta_0,\theta_1) \in (\bar{\Theta}_1,{\color{blue}\theta^*}) \text{ for some } \bar{\Theta}_0,\bar{\Theta}_1.$$

What are  $\bar{\Theta}_0, \bar{\Theta}_1$ ?

- $\mathbb{R}^d$  not an option: then  $\mathcal{H}_0$  and  $\mathcal{H}_1$  have the common point  $(\theta^*, \theta^*)$ .
- Thus we have to separate  $\bar{\Theta}_0, \bar{\Theta}_1$  from  $\theta^*$ .
- "Prediction-wise" separation:

$$\Delta_0 := L_0(\theta_1) - L_0(\theta_0) > 0, \quad \Delta_1 := L_1(\theta_0) - L_1(\theta_1) > 0.$$

Possible to recast this information in terms of  $\bar{\Theta}_0, \bar{\Theta}_1$ , but hardly useful.

### Main question

Characterize the **sample complexity** of distinguishing between  $\mathcal{H}_0$  and  $\mathcal{H}_1$  with fixed error probabilities of both types (say 2/3) in terms of  $\Delta_0, \Delta_1, ...$ 



## Linear regression setup

Well-specified linear regression: z = (x, y), and  $\mathbb{P}_k$ ,  $k \in \{0, 1\}$ , is given by

$$\mathbb{P}_k$$
:  $\mathbf{x} \sim \mathcal{N}(0, \mathbf{\Sigma}_k)$ ,  $\mathbf{y} = \mathbf{x}^{\top} \theta_k + \mathbf{\xi}$  with  $\mathbf{\xi} \sim \mathcal{N}(0, 1)$ 

- Write  $Z^k = (X^k; Y^k)$ , where  $X^k \in \mathbb{R}^{n \times d}$  and  $Y^k \in \mathbb{R}^n$  for  $k \in \{0, 1\}$ .
- Covariances  $\Sigma_k$  and their estimates:  $\widehat{\Sigma}_k := \frac{1}{n} X^{k \top} X^k$ .
- Population and empirical ranks:  $r_k = \text{rank}(\mathbf{\Sigma}_k)$ , and  $\hat{r}_k = \text{rank}(\hat{\mathbf{\Sigma}}_k)$ .
- Separations and their empirical counterparts:

$$\Delta_{k} = \|\theta_{1} - \theta_{0}\|_{\mathbf{\Sigma}_{k}}^{2} = \|\mathbf{\Sigma}_{k}^{1/2}(\theta_{1} - \theta_{0})\|^{2}$$
$$\widehat{\Delta}_{k} = \|\theta_{1} - \theta_{0}\|_{\widehat{\mathbf{\Sigma}}_{k}}^{2} = \frac{1}{n}\|X^{k}(\theta_{1} - \theta_{0})\|^{2}.$$

#### Basic test

Basic test based on the prediction error of  $\theta^*$  under  $\mathcal{H}_0$  and  $\mathcal{H}_1$ :

$$\mathbb{1}\left\{\|Y^0 - X^0\theta^*\|^2 {-} {\color{red} n} \geqslant \|Y^1 - X^1\theta^*\|^2 {-} {\color{red} n}\right\}.$$

Let  $\xi^k=Y^k-X^k\theta_k\sim\mathcal{N}(0,\boldsymbol{I}_n)$  be the noises. Under  $\mathcal{H}_0:\theta^*=\theta_0$  one has

LHS = 
$$\|\xi^0\|^2 - n$$
,  
RHS =  $\|\xi^1\|^2 - n - 2\langle \xi^1, X_1(\theta_0 - \theta_1) \rangle + \|X_1(\theta_1 - \theta_0)\|^2$ .

• Thus,  $\mathbb{E}[\mathsf{LHS}] = 0$  and  $\mathbb{E}[\mathsf{RHS}|X_1] = \|X_1(\theta_1 - \theta_0)\|^2 = n\widehat{\Delta}_1$ , where

$$\widehat{\Delta}_1 = \frac{1}{n} \|X_1(\theta_0 - \theta_1)\|^2 = \|\theta_0 - \theta_1\|_{\widehat{\Sigma}_1}^2$$

is the empirical counterpart of  $\Delta_1 = \|\theta_1 - \theta_0\|_{\mathbf{\Sigma}_1}^2$ .

ullet This motivates the basic test: type-I error  $\iff$  "fluctuations  $\geqslant n\Delta_1$ ."

### Basic test

$$\mathbb{1}\left\{\|Y^{0} - X^{0}\theta^{*}\|^{2} - n \geqslant \|Y^{1} - X^{1}\theta^{*}\|^{2} - n\right\}.$$

More precisely, LHS  $\sim \chi_n^2 - n$  and RHS $|X_1 \sim \chi_n^2 - n + 2\mathcal{N}(0, n\widehat{\Delta}_1) + n\widehat{\Delta}_1$ .

• Recall tail inequalities:  $\mathbb{P}[\mathcal{N}(0,1) \geqslant t] \leqslant \exp(-t^2)$  and  $\chi^2$ -bound:

$$\mathbb{P}[|\chi_s^2 - s| \geqslant t] \lesssim \exp(-c \min\{t, t^2/s\}),$$

Bound for the (conditional over  $X_0, X_1$ ) type-I error:

$$\begin{split} \mathbb{P}_I &= \mathbb{P}[\mathsf{fluctuations} \geqslant n \widehat{\Delta}_1] \\ &\leqslant \mathbb{P}\left[\chi_n^2 - n \geqslant \frac{n \widehat{\Delta}_1}{3}\right] + \mathbb{P}\left[n - \chi_n^2 \geqslant \frac{n \widehat{\Delta}_1}{3}\right] + \mathbb{P}\left[\mathcal{N}(0, n \widehat{\Delta}_1) \geqslant \frac{n \widehat{\Delta}_1}{6}\right] \\ &\lesssim \exp\left(-\frac{c n^2 \widehat{\Delta}_1^2}{n}\right) + \exp(-c n \widehat{\Delta}_1). \end{split}$$

• Thus, error prob. of both types at most  $\exp(-cn\min\{\Delta,\Delta^2\})$ , where  $\Delta:=\min\{\Delta_0,\Delta_1\}.$ 

If  $\Delta \lesssim 1$ : term  $\exp(-cn\Delta^2)$  dominates  $\Rightarrow O(1/\Delta^2)$  sample complexity.

## Improved test

Idea: decrease  $\chi^2$ -term fluctuations by projecting residuals on signal spaces.

#### Test for linear model

$$\widehat{T} = \mathbb{1}\left\{ \|\mathbf{\Pi}_{X^0}[Y^0 - X^0 \theta^*]\|^2 - \widehat{r}_0 \geqslant \|\mathbf{\Pi}_{X^1}[Y^1 - X^1 \theta^*]\|^2 - \widehat{r}_1 \right\},\,$$

where  $\Pi_X := X(X^{\top}X)^{\dagger}X^{\top}$  is the projector on signal space  $\operatorname{col}(X) \subseteq \mathbb{R}^n$ .

• Recall that  $\widehat{r}_k := \operatorname{rank}(\widehat{\Sigma}_k)$  and  $\widehat{\Sigma} = \frac{1}{n}X^\top X$ , hence indeed  $\dim(\operatorname{col}(X)) = \operatorname{Tr}(\Pi_X) = \operatorname{Tr}[(X^\top X)^\dagger X^\top X] = \operatorname{rank}(X^\top X) = \operatorname{rank}(\widehat{\Sigma}).$ 

## Improved test: analysis

#### Test for linear model

$$\widehat{T} = \mathbb{1}\left\{\|\Pi_{X^0}[Y^0 - X^0\theta^*]\|^2 - \widehat{r}_0 \geqslant \|\Pi_{X^1}[Y^1 - X^1\theta^*]\|^2 - \widehat{r}_1\right\},\,$$

where  $\Pi_X := X(X^\top X)^\dagger X^\top$  is the projector on signal space  $\operatorname{col}(X) \subseteq \mathbb{R}^n$ .

• For this test, under  $\mathcal{H}_0$ , we have

$$\mathsf{LHS}|X_0 \sim \chi_{\widehat{\mathbf{r_0}}}^2 - \widehat{\mathbf{r_0}}, \quad \mathsf{RHS}|X_1 \sim \chi_{\widehat{\mathbf{r_1}}}^2 - \widehat{\mathbf{r_1}} + 2\mathcal{N}(0, n\widehat{\Delta}_1) + n\widehat{\Delta}_1.$$

• Smaller  $\chi^2$  fluctuations since  $\widehat{r}_k \stackrel{a.s.}{=} \min\{r_k, n\} \leqslant n$ . Type-I error prob.:

$$\begin{split} & \mathbb{P}\left[\chi_{\widehat{r_0}}^2 - \widehat{r_0} \geqslant \frac{n\widehat{\Delta}_1}{3}\right] + \mathbb{P}\left[\widehat{r_1} - \chi_{\widehat{r_1}}^2 \geqslant \frac{n\widehat{\Delta}_1}{3}\right] + \mathbb{P}\left[\mathcal{N}(0, n\widehat{\Delta}_1) \geqslant \frac{n\widehat{\Delta}_1}{6}\right] \\ & \lesssim \exp\left(-\frac{cn^2\widehat{\Delta}_1^2}{\widehat{r_0}}\right) + \exp\left(-\frac{cn^2\widehat{\Delta}_1^2}{\widehat{r_1}}\right) + \exp(-cn\widehat{\Delta}_1). \end{split}$$

**Theorem.** Denoting  $r_{max} := max\{r_0, r_1\}$ , we have

$$\max\{P_I,P_{II}\} \lesssim \exp\left(-c\min\left\{n\Delta,\frac{n^2\Delta^2}{\min\{n,r_{\sf max}\}}\right\}\right).$$

## Improved test: sample complexity

### Error probability bound

**Theorem.** Denoting  $r_{max} := max\{r_0, r_1\}$ , we have

$$\max\{P_I, P_{II}\} = \exp\left(-c \min\left\{n\Delta, \frac{n^2\Delta^2}{\min\{n, r_{\mathsf{max}}\}}\right\}\right).$$

### Sample complexity bound

**Lemma** Assume  $\Delta \lesssim 1$ . Then  $-\log(\max\{P_I, P_{II}\}) \gtrsim 1$  is equivalent to

$$n \gtrsim \min\left\{rac{1}{\Delta^2}, rac{\sqrt{r_{\mathsf{max}}}}{\Delta}
ight\}.$$

### Comparison

Basic test: 
$$1 \{ \|Y^0 - X^0 \theta^*\|^2 - n \geqslant \|Y^1 - X^1 \theta^*\|^2 - n \}$$
. 
$$n = O\left(\frac{1}{\Delta^2}\right).$$

Improved test: 
$$\mathbb{1}\left\{\|\mathbf{\Pi}_{X^0}[Y^0 - X^0\theta^*]\|^2 - \widehat{r_0} \geqslant \|\mathbf{\Pi}_{X^1}[Y^1 - X^1\theta^*]\|^2 - \widehat{r_1}\right\}.$$
 
$$n = O\left(\min\left\{\frac{1}{\Delta^2}, \frac{\sqrt{r_{\mathsf{max}}}}{\Delta}\right\}\right).$$

Note:  $\widehat{r}_k \stackrel{a.s.}{=} \min\{r_k, n\}$  and  $\Pi_{X^k}$  projects on  $\operatorname{col}(X^k) \subset \mathbb{R}^n$  of dimension  $\widehat{r}_k$ . Thus, the two tests coincide when  $n \leqslant \min\{r_0, r_1\}$ ,

- Well-sep. regime:  $\Delta \gtrsim \frac{1}{\sqrt{r_{\text{max}}}}$ . Sample complexity  $\lesssim r_{\text{max}}$  and rank-independent. No need for projections if  $r_0 \approx r_1$ .
- III-sep. regime:  $\Delta \ll \frac{1}{\sqrt{r_{\text{max}}}}$ , sample complexity  $\gg r_{\text{max}}$ , projections.

# Testing vs. estimation

$$\text{Improved test: } \mathbb{1}\big\{n\|\theta^* - \widehat{\theta}_0\|_{\widehat{\Sigma}_0}^2 - \widehat{r}_0 \geqslant n\|\theta^* - \widehat{\theta}_1\|_{\widehat{\Sigma}_1}^2 - \widehat{r}_1\big\}.$$

Sample complexity: 
$$O\left(\min\left\{\frac{1}{\Delta^2}, \frac{\sqrt{r_{\max}}}{\Delta}\right\}\right)$$

# Testing vs. estimation

$$\begin{split} \text{Improved test:} \quad & \mathbb{1} \left\{ n \| \theta^* - \widehat{\theta}_0 \|_{\widehat{\Sigma}_0}^2 - \widehat{r_0} \geqslant n \| \theta^* - \widehat{\theta}_1 \|_{\widehat{\Sigma}_1}^2 - \widehat{r}_1 \right\}. \\ & \text{Sample complexity:} \quad & O\left( \min \left\{ \frac{1}{\Delta^2}, \frac{\sqrt{r_{\mathsf{max}}}}{\Delta} \right\} \right) \ll \frac{r_{\mathsf{max}}}{\Delta}. \end{split}$$

• Sample complexity of estimating  $\bar{\theta} = \theta_0 + \theta_1 - \theta^*$  up to  $\Delta$  prediction error (i.e., better than by  $\theta^*$ ) is at least  $\frac{r_{\min}}{\Delta} \left[ \approx \frac{r_{\max}}{\Delta} \text{ when } r_0 \asymp r_1 \right]$ .

#### Non-disclosure

We can discriminate between  $\mathcal{H}_0$  and  $\mathcal{H}_1$  with sample size that does not allow to estimate the complementary model  $\bar{\theta}$  (with better quality than  $\theta^*$ ).

• Rich potential for applications in "privacy-aware ML" (see our paper).

## Interpretation via least-squares

Recall the normal equations for the least-squares estimates  $\widehat{\theta}_0, \widehat{\theta}_1$  of  $\theta_0, \theta_1$ :

$$\widehat{\boldsymbol{\Sigma}}_0\widehat{\boldsymbol{\theta}}_0 = \tfrac{1}{n}\boldsymbol{X}^{0\top}\boldsymbol{Y}^0, \quad \widehat{\boldsymbol{\Sigma}}_1\widehat{\boldsymbol{\theta}}_1 = \tfrac{1}{n}\boldsymbol{X}^{1\top}\boldsymbol{Y}^1.$$

This allows to rewrite the squared norms of the projected residuals:

$$\begin{split} \|\mathbf{\Pi}_{X}[Y - X\theta^{*}]\|^{2} &= (Y - X\theta^{*})^{\top} \mathbf{\Pi}_{X}(Y - X\theta^{*}) \\ &= (X^{\top}Y - X^{\top}X\theta^{*})^{\top} (X^{\top}X)^{\dagger} (X^{\top}Y - X^{\top}X\theta^{*}) \\ &= n^{2} (\widehat{\mathbf{\Sigma}}(\widehat{\theta} - \theta^{*}))^{\top} (X^{\top}X)^{\dagger} \widehat{\mathbf{\Sigma}}(\widehat{\theta} - \theta^{*}) \\ &= n(\widehat{\theta} - \theta^{*})^{\top} \widehat{\mathbf{\Sigma}} \widehat{\mathbf{\Sigma}}^{\dagger} \widehat{\mathbf{\Sigma}}(\widehat{\theta} - \theta^{*}) = n(\widehat{\theta} - \theta^{*})^{\top} \widehat{\mathbf{\Sigma}}(\widehat{\theta} - \theta^{*}) \\ &= n \|\widehat{\theta} - \theta^{*}\|_{\widehat{\mathbf{\Sigma}}}^{2}. \end{split}$$

Thus, our test amounts to  $\mathbb{1}\{\|\theta^*-\widehat{\theta}_0\|_{\widehat{\widehat{\Sigma}}_0}^2-\frac{\widehat{r}_0}{n}\geqslant \|\theta^*-\widehat{\theta}_1\|_{\widehat{\widehat{\Sigma}}_1}^2-\frac{\widehat{r}_1}{n}\}.$ 

- We compare the empirical prediction distances from  $\widehat{\theta}^*$  to  $\widehat{\theta}_0$  and  $\widehat{\theta}_1$  after debiasing them under the matching hypothesis.
- **NB**: we don't require  $\widehat{\theta}_0$ ,  $\widehat{\theta}_1$  to be unique (i.e.  $n \ge r_{\text{max}}$ ).

## Lower bound: key ideas

We need to prove two bounds:

$$\inf_{\widehat{T}} \sup_{\|\theta_1 - \theta_0\|_{\widehat{\Sigma}}^2 \geqslant \Delta} P_I(\widehat{T}) + P_{II}(\widehat{T}) \gtrsim \max \left\{ \exp(-cn\Delta), \ \exp\left(-c\frac{n^2\Delta^2}{\min\{n,r\}}\right) \right\}.$$

• First bound: easier problem with known  $\bar{\theta}$  and simple hypotheses:

$$\mathcal{H}_0^o:(\theta_0,\theta_1)=(\theta^*,\bar{\theta}),\quad \text{vs.}\quad \mathcal{H}_1^o:(\theta_0,\theta_1)=(\bar{\theta},\theta^*).$$

• Likelihood-ratio (LR) test

$$T_{\mathsf{LR}} = \mathbb{1}\{\|Y^0 - X^0\theta^*\|^2 + \|Y^1 - X^1\bar{\theta}\|^2 \ge \|Y^0 - X^0\bar{\theta}\|^2 + \|Y^1 - X^1\theta^*\|^2\}$$

is optimal (w.r.t. sum of errors) by the Neyman-Pearson lemma.

- Second bound captures dependence on the rank. Bayesian approach:
  - Put Gaussian prior on  $\bar{\theta}$ , lower-bound max $\{\mathbb{P}_I,\mathbb{P}_{II}\}$  for the Bayes test.
  - Lower-bounding is technical and requires that  $\widehat{\Sigma}_0, \widehat{\Sigma}_1$  commute.
  - We achieve this by sampling x and  $\tilde{x}$  from  $\sqrt{r}\{\pm e_1,...,\pm e_r\}$ .

# Beyond linear models

#### In our paper:

- General result for parametric models in asymptotic regime  $n \to \infty$  with fixed  $r_0, r_1$  and  $n\Delta \to \lambda$ .
- Technical result for generalized linear models (GLMs) allowing for heavy tails and misspecification.
- Same general picture:

$$\max\{P_I, P_{II}\} = \exp\left(-c \min\left\{n\Delta, \frac{n^2\Delta^2}{\max\{\rho_0, \rho_1\}}\right\}\right)$$

where  $\rho_0, \rho_1$  are "effective model ranks".

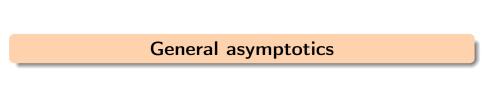
#### Open questions:

- General nonasymptotic result;
- Full optimality;
- Adaptation.

### Thank you!

And check our paper:

arxiv.org/abs/2012.02901



# General setup: Newton decrement test

Linear model: 
$$\mathbb{1}\left\{\|\mathbf{\Pi}_{X^0}[Y^0 - X^0\theta^*]\|^2 - \widehat{r}_0 \geqslant \|\mathbf{\Pi}_{X^1}[Y^1 - X^1\theta^*]\|^2 - \widehat{r}_1\right\}$$
.

#### General setup:

• Empirical risk  $\widehat{L}_k(\theta) = \frac{1}{n} \sum_{i=1}^n \ell_{z_i^k}(\theta)$  has gradient  $\nabla \widehat{L}_k(\theta)$  and Hessian  $\widehat{\boldsymbol{H}}_k(\theta)$ :

$$\widehat{\boldsymbol{H}}_k(\theta) := \nabla^2 \widehat{L}_k(\theta), \quad \boldsymbol{H}_k(\theta) := \nabla^2 L_k(\theta).$$

• Let  $G_k(\theta) := Cov_{\mathbb{P}_k}[\nabla \ell_z(\theta)]$ . For well-specified models:

$$G_k(\theta_k) = H_k(\theta_k).$$

- Standardized Fisher matrix:  $J_k(\theta) := H_k(\theta)^{-\dagger/2} G_k(\theta) H_k(\theta)^{-\dagger/2}$ .
- Effective rank  $\rho_k := \text{Tr}[J_k(\theta_k)]$ . For well-specified models:  $\rho_k = r_k$ .

In linear regression  $\nabla \widehat{L}(\theta) = \frac{1}{n} X^{\top} (Y - X\theta)$  and  $\nabla^2 \widehat{L}(\theta) \equiv \frac{1}{n} X^{\top} X$ , hence

$$\|\mathbf{\Pi}_{X}[Y - X\theta^{*}]\|^{2} = \|(X^{\top}X)^{\dagger/2}X^{\top}(Y - X\theta^{*})\|^{2} = n\|\widehat{\mathbf{H}}(\theta^{*})^{\dagger/2}\nabla\widehat{L}(\theta^{*})\|^{2}.$$

# General setup: Newton decrement test (cont'd)

$$\mathbb{1}\left\{\|\Pi_{X^0}[Y^0-X^0\theta^*]\|^2-\widehat{r}_0\geqslant \|\Pi_{X^1}[Y^1-X^1\theta^*]\|^2-\widehat{r}_1\right\}.$$

- Replace  $\|\mathbf{\Pi}_{X^k}[Y^k X^k\theta^*]\|^2$  with  $n\|\widehat{\mathbf{H}}_k(\theta_k)^{\dagger/2}\nabla\widehat{L}_k(\theta^*)\|^2$ .
- When  $n \to \infty$ ,

$$\mathbb{E}_{k}[n\|\widehat{\boldsymbol{H}}_{k}(\theta_{k})^{\dagger/2}\nabla\widehat{L}_{k}(\theta_{k})\|^{2}] \to \rho_{k} = \mathsf{Tr}[\boldsymbol{J}_{k}(\theta_{k})].$$

• Cannot use  $\rho_k$ 's as one of them uses  $\bar{\theta}$  which is unknown. Instead use

$$\operatorname{Tr}[\boldsymbol{J}_{k}(\boldsymbol{\theta}^{*})] = n_{k} \mathbb{E}_{k} \left[ \left\| \boldsymbol{H}_{k}(\boldsymbol{\theta}^{*})^{\dagger/2} \left[ \nabla \widehat{L}_{k}(\boldsymbol{\theta}^{*}) - \nabla L_{k}(\boldsymbol{\theta}^{*}) \right] \right\|^{2} \right],$$

or, more precisely, its asymptotically (as  $n \to \infty$ ) unbiased estimate:

$$\widehat{\mathsf{T}}_k = \tfrac{1}{2} n_k \big\| \boldsymbol{H}_k(\boldsymbol{\theta}^*)^{\dagger/2} \big[ \nabla \widehat{\boldsymbol{L}}_k(\boldsymbol{\theta}^*) - \widehat{\nabla} \boldsymbol{L}_k'(\boldsymbol{\theta}^*) \big] \big\|^2.$$

$$\widehat{\mathcal{T}} = \mathbb{1}\big\{ \textit{n}_0 \| \widehat{\boldsymbol{H}}_0(\boldsymbol{\theta}^*)^{\dagger/2} \nabla \widehat{\textit{L}}_0(\boldsymbol{\theta}^*) \|^2 - \widehat{T}_0 \geqslant \textit{n}_1 \| \widehat{\boldsymbol{H}}_1(\boldsymbol{\theta}^*)^{\dagger/2} \nabla \widehat{\textit{L}}_1(\boldsymbol{\theta}^*) \|^2 - \widehat{T}_1 \big\}.$$

**Theorem.** Denoting  $\rho_{max} := \max\{\rho_0, \rho_1\}$ , we have that

$$\lim_{n\to\infty}[\max\{P_I,P_{II}\}]\lesssim \exp\left(-c\min\left\{n\Delta,\frac{n^2\Delta^2}{\rho_{\max}}\right\}\right).$$