

# Near-Optimal Model Discrimination with Non-Disclosure Properties

`arxiv.org/abs/2012.02901`

**Jointly with**

Mohamed Ndaoud (ESSEC)

Adel Javanmard (USC)    Meisam Razaviyayn (USC)

- General problem formulation
- Linear models
- Extensions

## General problem formulation

# 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.

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

with strictly convex loss  $\ell(\cdot, z) : \mathbb{R}^d \rightarrow \mathbb{R}$ , 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  $\ell(\cdot, z)$ , and observes **two** i.i.d. samples:

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

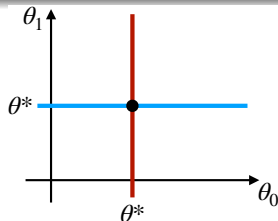
- **Task:** distinguish between the two hypotheses

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

# Model discrimination task

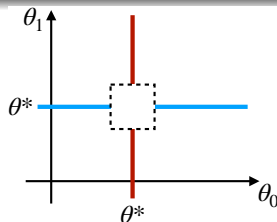
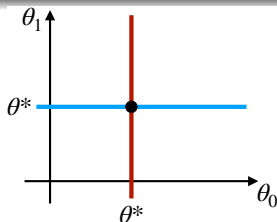
- **Classical setup:** both  $\theta_0, \theta_1$  known; one sample  $Z \sim \mathbb{P}_\theta^{\otimes n}$  observed.  
*Which  $\theta \in \{\theta_0, \theta_1\}$  corresponds to the sample?*  
*Two **simple** hypotheses about  $\theta$ .*
- **Our setup:** we observe both samples but only one model  $\theta^* \in \{\theta_0, \theta_1\}$ .  
*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  $\ell(\cdot, z)$ , and observes **two** i.i.d. samples:  
$$Z^0 = (z_1^0, \dots, z_n^0) \sim \mathbb{P}_0^{\otimes n}, \quad Z^1 = (z_1^1, \dots, z_n^1) \sim \mathbb{P}_1^{\otimes n}.$$
- **Task:** distinguish between the two hypotheses about  $(\theta_0, \theta_1) \in \mathbb{R}^{2d}$ :  
$$\mathcal{H}_0 : (\theta_0, \theta_1) \in \{\theta^*\} \times \{\theta \neq \theta^*\} \quad \text{vs.} \quad \mathcal{H}_1 : (\theta_0, \theta_1) \in \{\theta \neq \theta^*\} \times \{\theta^*\}$$

# Separation and sample complexity



$$\mathcal{H}_0 : (\theta_0, \theta_1) \in \{\theta^*\} \times \{\theta \neq \theta^*\} \quad \text{vs.} \quad \mathcal{H}_1 : (\theta_0, \theta_1) \in \{\theta \neq \theta^*\} \times \{\theta^*\}$$

# Separation and sample complexity



$$\mathcal{H}_0 : (\theta_0, \theta_1) \in \{\theta^*\} \times \{\theta \neq \theta^*\} \quad \text{vs.} \quad \mathcal{H}_1 : (\theta_0, \theta_1) \in \{\theta \neq \theta^*\} \times \{\theta^*\}$$

- **Separate**  $\theta_0$  and  $\theta_1$  to exclude the degenerate case  $\theta_0 = \theta_1 = \theta^*$ .

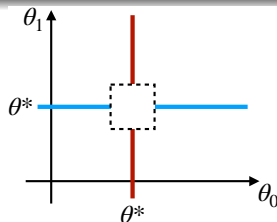
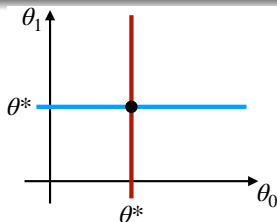
$$\mathcal{H}_0 : (\theta_0, \theta_1) \in \{\theta^*\} \times \bar{\Theta}_0 \quad \text{vs.} \quad \mathcal{H}_1 : (\theta_0, \theta_1) \in \bar{\Theta}_1 \times \{\theta^*\}$$

- “Prediction-wise” separation:

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

- Implicitly choose  $\bar{\Theta}_0, \bar{\Theta}_1$  according to this separation assumption.

# Separation and sample complexity



$$\mathcal{H}_0 : (\theta_0, \theta_1) \in \{\theta^*\} \times \{\theta \neq \theta^*\} \quad \text{vs.} \quad \mathcal{H}_1 : (\theta_0, \theta_1) \in \{\theta \neq \theta^*\} \times \{\theta^*\}$$

- **Separate**  $\theta_0$  and  $\theta_1$  to exclude the degenerate case  $\theta_0 = \theta_1 = \theta^*$ .

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

- “Prediction-wise” separation:

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

- Implicitly choose  $\bar{\Theta}_0, \bar{\Theta}_1$  according to this separation assumption.

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  $\Delta_0, \Delta_1, d$ .



## Well-specified linear models

Well-specified linear model:  $z = (x, y) \in \mathbb{R}^{d+1}$ ,  $\ell(\theta, z) = \frac{1}{2}(x^\top \theta - y)^2$ , and

$\mathbb{P}_k : x \sim \mathcal{N}(0, \Sigma_k)$ ,  $y = x^\top \theta_k + \xi$  with  $\xi \sim \mathcal{N}(0, 1)$  for  $k \in \{0, 1\}$ .

- 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:  $\hat{\Sigma}_k := \frac{1}{n} X_k^\top X_k$ .
- Population and empirical ranks:  $r_k = \text{rank}(\Sigma_k)$  and  $\hat{r}_k = \text{rank}(\hat{\Sigma}_k)$ .
- Separations and their empirical counterparts:

$$\Delta_k = \frac{1}{2} \|\theta_1 - \theta_0\|_{\Sigma_k}^2 = \frac{1}{2} \|\Sigma_k^{1/2}(\theta_1 - \theta_0)\|^2,$$
$$\hat{\Delta}_k = \frac{1}{2} \|\theta_1 - \theta_0\|_{\hat{\Sigma}_k}^2 = \frac{1}{2n} \|X_k(\theta_1 - \theta_0)\|^2.$$

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

$$\mathbb{1} \left\{ \|Y_0 - X_0\theta^*\|^2 - n \geq \|Y_1 - X_1\theta^*\|^2 - n \right\}.$$

Let  $\xi_k = Y_k - X_k\theta_k \sim \mathcal{N}(0, I_n)$  be the noises. Under  $\mathcal{H}_0 : \theta^* = \theta_0$ , one has

$$\text{LHS} = \|\xi_0\|^2 - n,$$

$$\text{RHS} = \|\xi_1\|^2 - n - 2 \langle \xi_1, X_1(\theta_0 - \theta_1) \rangle + \|X_1(\theta_1 - \theta_0)\|^2.$$

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

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

is the empirical counterpart of  $\Delta_1 = \|\theta_1 - \theta_0\|_{\Sigma_1}^2$ .

- This motivates the basic test: type-I error  $\iff$  “fluctuations  $\geq n\Delta_1$ .”

$$\mathbb{1} \left\{ \|Y_0 - X_0 \theta^*\|^2 - n \geq \|Y_1 - X_1 \theta^*\|^2 - n \right\}.$$

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

- Basic tail inequalities for Gaussian and  $\chi^2$  laws:

$$\mathbb{P}[\mathcal{N}(0, 1) \geq u] \leq \exp(-u^2), \quad \mathbb{P}[|\chi_s^2 - s| \geq v] \lesssim \exp(-c \min\{v, v^2/s\}).$$

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

$$\begin{aligned} \mathbb{P}_I &= \mathbb{P}[\text{fluctuations} \geq n\hat{\Delta}_1] \\ &\leq \mathbb{P}\left[\chi_n^2 - n \geq \frac{n\hat{\Delta}_1}{3}\right] + \mathbb{P}\left[n - \chi_n^2 \geq \frac{n\hat{\Delta}_1}{3}\right] + \mathbb{P}\left[\mathcal{N}(0, n\hat{\Delta}_1) \geq \frac{n\hat{\Delta}_1}{6}\right] \\ &\lesssim \exp\left(-\frac{cn^2\hat{\Delta}_1^2}{n}\right) + \exp(-cn\hat{\Delta}_1). \end{aligned}$$

- 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.

**Idea:** reduce  $\chi^2$  fluctuations by projecting the residuals on signal spaces.

## Test for linear model

$$\hat{T} = \mathbb{1} \left\{ \|\Pi_{X_0}[Y_0 - X_0\theta^*]\|^2 - \hat{r}_0 \geq \|\Pi_{X_1}[Y_1 - X_1\theta^*]\|^2 - \hat{r}_1 \right\},$$

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

- Recall that  $\hat{r}_k := \text{rank}(\hat{\Sigma}_k)$  and  $\hat{\Sigma} = \frac{1}{n}X^\top X$ , hence indeed

$$\dim(\text{col}(X)) = \text{Tr}(\Pi_X) = \text{Tr}[(X^\top X)^\dagger X^\top X] = \text{rank}(X^\top X) = \text{rank}(\hat{\Sigma}).$$

## Test for linear model

$$\hat{T} = \mathbb{1} \left\{ \|\Pi_{X_0}[Y_0 - X_0\theta^*]\|^2 - \hat{r}_0 \geq \|\Pi_{X_1}[Y_1 - X_1\theta^*]\|^2 - \hat{r}_1 \right\},$$

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

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

$$\text{LHS}|X_0 \sim \chi_{\hat{r}_0}^2 - \hat{r}_0, \quad \text{RHS}|X_1 \sim \chi_{\hat{r}_1}^2 - \hat{r}_1 + 2\mathcal{N}(0, n\hat{\Delta}_1) + n\hat{\Delta}_1.$$

- Smaller  $\chi^2$  fluctuations since  $\hat{r}_k \stackrel{\text{a.s.}}{\leq} \min\{r_k, n\} \leq n$ . Type-I error prob.:

$$\begin{aligned} & \mathbb{P} \left[ \chi_{\hat{r}_0}^2 - \hat{r}_0 \geq \frac{n\hat{\Delta}_1}{3} \right] + \mathbb{P} \left[ \hat{r}_1 - \chi_{\hat{r}_1}^2 \geq \frac{n\hat{\Delta}_1}{3} \right] + \mathbb{P} \left[ \mathcal{N}(0, n\hat{\Delta}_1) \geq \frac{n\hat{\Delta}_1}{6} \right] \\ & \lesssim \exp \left( -\frac{cn^2\hat{\Delta}_1^2}{\hat{r}_0} \right) + \exp \left( -\frac{cn^2\hat{\Delta}_1^2}{\hat{r}_1} \right) + \exp(-cn\hat{\Delta}_1). \end{aligned}$$

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

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

## Error probability bound

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

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

## Error probability bound

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

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

## Sample complexity bound

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

$$n \gtrsim \min\left\{\frac{1}{\Delta^2}, \frac{\sqrt{r_{\max}}}{\Delta}\right\}.$$



## Error probability bound

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

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

## Sample complexity bound

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

$$n \gtrsim \min\left\{\frac{1}{\Delta^2}, \frac{\sqrt{r_{\max}}}{\Delta}\right\}.$$

**Proof:**

1. Prove:  $n\Delta \gtrsim \min\left\{\frac{1}{\Delta}, \sqrt{r_{\max}}\right\} \iff n\Delta \min\left\{1, \frac{n\Delta}{\min\{n, r_{\max}\}}\right\} \gtrsim 1$  if  $\Delta \lesssim 1$ .
2. The second condition reads  $n\Delta \gtrsim \max\left\{1, \min\left\{\frac{1}{\Delta}, \frac{r_{\max}}{n\Delta}\right\}\right\}$ , or equivalently  $n\Delta \gtrsim \min\left\{\frac{1}{\Delta}, \max\left\{1, \frac{r_{\max}}{n\Delta}\right\}\right\}$  by using  $\Delta \lesssim 1$  and treating all possible cases.
3. It remains to verify that  $n\Delta \gtrsim \sqrt{r_{\max}}$  if and only if  $n\Delta \gtrsim \max\left\{1, \frac{r_{\max}}{n\Delta}\right\}$ .  $\square$

## Basic test

$$\mathbb{1} \left\{ \|Y_0 - X_0 \theta^*\|^2 - n \geq \|Y_1 - X_1 \theta^*\|^2 - n \right\},$$

$$\text{Sample complexity: } n = O\left(\frac{1}{\Delta^2}\right).$$

## Improved test

$$\mathbb{1} \left\{ \|\Pi_{X_0}[Y_0 - X_0 \theta^*]\|^2 - \hat{r}_0 \geq \|\Pi_{X_1}[Y_1 - X_1 \theta^*]\|^2 - \hat{r}_1 \right\}.$$

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

Note:  $\hat{r}_k \stackrel{\text{a.s.}}{=} \min\{r_k, n\}$  and  $\Pi_{X_k}$  projects on  $\text{col}(X_k) \subset \mathbb{R}^n$  of dimension  $\hat{r}_k$ . Thus, the tests coincide when  $n \leq \min\{r_0, r_1\}$ . In fact, a **phase transition**:

- **Well-separated:**  $\Delta \gtrsim \frac{1}{\sqrt{r_{\max}}}$ . Sample complexity  $n = O(1/\Delta^2) \lesssim r_{\max}$ .
- **Ill-separated:**  $\Delta \ll \frac{1}{\sqrt{r_{\max}}}$ . Sample complexity  $\gg r_{\max} \Rightarrow$  projections.

# Interpretation via least-squares

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

$$\hat{\Sigma}_0 \hat{\theta}_0 = \frac{1}{n} X_0^\top Y_0, \quad \hat{\Sigma}_1 \hat{\theta}_1 = \frac{1}{n} X_1^\top Y_1.$$

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

$$\begin{aligned} \|\Pi_X[Y - X\theta^*]\|^2 &= (Y - X\theta^*)^\top \Pi_X (Y - X\theta^*) \\ &= (X^\top Y - X^\top X\theta^*)^\top (X^\top X)^\dagger (X^\top Y - X^\top X\theta^*) \\ &= n^2 (\hat{\Sigma}(\hat{\theta} - \theta^*))^\top (X^\top X)^\dagger \hat{\Sigma}(\hat{\theta} - \theta^*) \\ &= n(\hat{\theta} - \theta^*)^\top \hat{\Sigma} \hat{\Sigma}^\dagger \hat{\Sigma}(\hat{\theta} - \theta^*) = n(\hat{\theta} - \theta^*)^\top \hat{\Sigma}(\hat{\theta} - \theta^*) \\ &= n \|\hat{\theta} - \theta^*\|_{\hat{\Sigma}}^2. \end{aligned}$$

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

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

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

- 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  $\frac{r_{\min}}{\Delta}$ .
- Thus, when  $r_0 \approx r_1$ , **recovery is way harder than discrimination!**

## Non-disclosure property

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

- In fact, our tests access  $\theta^*$  through “scalar” statistic  $\|\Pi_X[Y - X\theta^*]\|^2$  that carries only  $O(1)$  Fisher information about  $\theta^*$ .
- Hence, we also guarantee **non-disclosure of  $\theta^*$**  (up to accuracy  $\Delta$ ).

## Lower bound: key ideas

We need to prove two bounds:

$$\inf_{\hat{T}} \sup_{\|\theta_1 - \theta_0\|_{r'}^2 \geq \Delta} P_I(\hat{T}) + P_{II}(\hat{T}) \gtrsim \max \left\{ \exp(-cn\Delta), \exp\left(-c \frac{n^2 \Delta^2}{\min\{n, r\}}\right) \right\}.$$

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

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

Likelihood-ratio (LR) test

$$T_{\text{LR}} = \mathbb{1}\{\|Y_0 - X_0\theta^*\|^2 + \|Y_1 - X_1\bar{\theta}\|^2 \geq \|Y_0 - X_0\bar{\theta}\|^2 + \|Y_1 - X_1\theta^*\|^2\}$$

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

$$\begin{aligned} & \mathbb{P}_{\tilde{\mathcal{H}}_0}[T_{\text{LR}} = 1 | X_0, X_1] \\ &= \mathbb{P}[\|Y_0 - X_0\theta_0\|^2 + \|Y_1 - X_1\theta_1\|^2 \geq \|Y_0 - X_0\theta_1\|^2 + \|Y_1 - X_1\theta_0\|^2 | X_0, X_1] \\ &= \mathbb{P}[2\langle \xi_0, X_0(\theta_0 - \theta_1) \rangle + 2\langle \xi_1, X_1(\theta_0 - \theta_1) \rangle \geq \|X_0(\theta_0 - \theta_1)\|^2 + \|X_1(\theta_0 - \theta_1)\|^2] \\ &\geq \mathbb{P}[2\mathcal{N}(0, n\hat{\Delta}_0) + 2\mathcal{N}(0, n\hat{\Delta}_1) \geq n\hat{\Delta}_0 + n\hat{\Delta}_1] \\ &\geq \mathbb{P}[\mathcal{N}(0, n\hat{\Delta}_0) \geq n\hat{\Delta}_0/2] \cdot \mathbb{P}[\mathcal{N}(0, n\hat{\Delta}_1) \geq n\hat{\Delta}_1/2] \gtrsim \exp(-cn \max\{\hat{\Delta}_0, \hat{\Delta}_1\}). \end{aligned}$$

Then  $\max\{\hat{\Delta}_0, \hat{\Delta}_1\} \lesssim \Delta$  with fixed probability by Markov's inequality.

We need to prove two bounds:

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

**Second bound** captures dependence on the rank. Bayesian approach:

- Put a Gaussian prior  $\pi$  on  $\bar{\theta}$  such that  $\pi\{\|\bar{\theta} - \theta^*\|_{I_r}^2 \leq \Delta\}$  is very small.
- This allows to lower-bound the maximal risk by the Bayes risk.
- The Bayes risk can be lower-bounded by the Neyman-Pearson lemma—through surprisingly tedious calculations.  $\square$

## In our paper:

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

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

where  $\rho_0, \rho_1$  are “effective model ranks”.

## Open questions:

- Closing the gap for linear models
- General nonasymptotic result
- Mixtures
- New insights on two-sample testing?

**Thank you!**



## General asymptotics

# General setup: Newton decrement test

**Linear model:**  $\mathbb{1} \{ \|\boldsymbol{\Pi}_{X_0}[Y_0 - X_0\theta^*]\|^2 - \hat{r}_0 \geq \|\boldsymbol{\Pi}_{X_1}[Y_1 - X_1\theta^*]\|^2 - \hat{r}_1 \}.$

## General setup:

- Empirical risk  $\hat{L}_k(\theta) = \frac{1}{n} \sum_{i=1}^n \ell(\theta, z_i^{(k)})$  has gradient  $\nabla \hat{L}_k(\theta)$  and Hessian  $\hat{\mathbf{H}}_k(\theta)$ :

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

- Let  $\mathbf{G}_k(\theta) := \text{Cov}_{\mathbb{P}_k}[\nabla \ell_z(\theta)]$ . **For well-specified models:**

$$\mathbf{G}_k(\theta_k) = \mathbf{H}_k(\theta_k).$$

- Standardized Fisher matrix:  $\mathbf{J}_k(\theta) := \mathbf{H}_k(\theta)^{-\dagger/2} \mathbf{G}_k(\theta) \mathbf{H}_k(\theta)^{-\dagger/2}.$
- Effective rank  $\rho_k := \text{Tr}[\mathbf{J}_k(\theta_k)]$ . **For well-specified models:**  $\rho_k = r_k.$

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

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

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

$$\mathbb{1} \{ \|\mathbf{\Pi}_{X_0}[Y_0 - X_0\theta^*]\|^2 - \hat{r}_0 \geq \|\mathbf{\Pi}_{X_1}[Y_1 - X_1\theta^*]\|^2 - \hat{r}_1 \}.$$

- Replace  $\|\mathbf{\Pi}_{X_k}[Y_k - X_k\theta^*]\|^2$  with  $n\|\hat{\mathbf{H}}_k(\theta_k)^{\dagger/2}\nabla\hat{L}_k(\theta^*)\|^2$ .
- When  $n \rightarrow \infty$ ,

$$\mathbb{E}_k[n\|\hat{\mathbf{H}}_k(\theta_k)^{\dagger/2}\nabla\hat{L}_k(\theta_k)\|^2] \rightarrow \rho_k = \text{Tr}[\mathbf{J}_k(\theta_k)].$$

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

$$\text{Tr}[\mathbf{J}_k(\theta^*)] = n_k \mathbb{E}_k [\|\mathbf{H}_k(\theta^*)^{\dagger/2}(\nabla\hat{L}_k(\theta^*) - \nabla L_k(\theta^*))\|^2],$$

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

$$\hat{\mathbf{T}}_k = \frac{1}{2}n_k\|\mathbf{H}_k(\theta^*)^{\dagger/2}(\nabla\hat{L}_k(\theta^*) - \hat{\nabla}L'_k(\theta^*))\|^2.$$

$$\hat{\mathbf{T}} = \mathbb{1} \{ n_0\|\hat{\mathbf{H}}_0(\theta^*)^{\dagger/2}\nabla\hat{L}_0(\theta^*)\|^2 - \hat{\mathbf{T}}_0 \geq n_1\|\hat{\mathbf{H}}_1(\theta^*)^{\dagger/2}\nabla\hat{L}_1(\theta^*)\|^2 - \hat{\mathbf{T}}_1 \}.$$

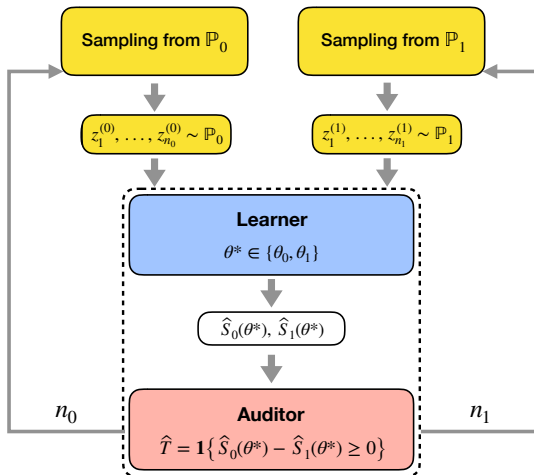
**Theorem.** Denoting  $\rho_{\max} := \max\{\rho_0, \rho_1\}$ , we have that

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

# Applications

# Applications: generic testing protocol

**Key observation:**  $\theta^*$  does not have to be known to run the test, and it cannot be inferred from  $Z_0, Z_1$  when they are small.



- We want to protect  $\theta^*$  and  $\bar{\theta} = \theta_0 + \theta_1 - \theta^*$  from inference via  $Z_0, Z_1$ .

# Application #1: testing for data deletion

## Testing for data deletion

- Company FAANG<sup>1</sup> trained a prediction model  $\theta^*$  on a large dataset  $\mathbb{P}^*$  pertaining to many users.
- Some users ask their data to be removed—and  $\theta^*$  retrained accordingly.
- FAANG should comply—and would like to demonstrate the compliance.
- Model  $\theta^*$  is proprietary, hence FAANG would like to avoid disclosing it.

Given a subsample  $Z^* \sim \mathbb{P}^*$  of FAANG's dataset and the pool  $Q$  of deletion queries, we can check that FAANG indeed retrained the model excluding  $Q$ .

- Let  $\mathbb{P}_0, \mathbb{P}_1$  correspond to hypotheses  $\mathcal{H}_0 : \mathbb{P}^* = \mathbb{P}_0$  (“clean data”) and

$$\mathcal{H}_1 : \mathbb{P}^* = \mathbb{P}_1 := (1 - \delta)\mathbb{P}_0 + \delta Q,$$

where  $\mathbb{P}_0$  is “clean” data, and  $\delta \in (0, 1)$  is the share of deletion queries.

- FAANG (“Learner”) gives to the tester (“Auditor”) access to  $Z_0 \sim \mathbb{P}_0$ .
- Having  $Z_0$ ,  $Q$ , and  $\delta$ , Auditor can generate  $Z_1 = (1 - \delta)Z_0 + \delta Q \sim \mathbb{P}_1$ .

## Application #2: testing fair representation of subpopulations

- Let  $\mathbb{P}_{\text{dem}}$  and  $\mathbb{P}_{\text{rep}}$  be two populations: Democrats and Republicans.
- We want them to be equally represented in the dataset  $\mathbb{P}^*$ .

- Define the hypotheses:

$$\mathcal{H}_0 : \mathbb{P}^* = \mathbb{P}_0 := \frac{1}{2}\mathbb{P}_{\text{dem}} + \frac{1}{2}\mathbb{P}_{\text{rep}},$$

$$\mathcal{H}_1 : \mathbb{P}^* = \mathbb{P}_1 := (\frac{1}{2} + \delta)\mathbb{P}_{\text{dem}} + (\frac{1}{2} - \delta)\mathbb{P}_{\text{rep}} \text{ for some } \delta \in (-\frac{1}{2}, \frac{1}{2}).$$

- Knowing  $\delta$ , we can easily implement the sampling oracles for  $\mathbb{P}_0$  and  $\mathbb{P}_1$  can be implemented using those for  $\mathbb{P}_{\text{dem}}$  and  $\mathbb{P}_{\text{rep}}$ .