

# Task Shift: From Classification to Regression in Overparameterized Linear Models



Tyler LaBonte\*, Kuo-Wei Lai\*, Vidya Muthukumar

Georgia Institute of Technology  
\*Equal contribution



## MOTIVATION

Modern classifiers adapt in-context to new tasks!

What is  $y_{10}$  without any reasoning?  
 $x_1 = [0.71, -0.55], y_1 = 1.38$   
 $x_2 = [-1.31, -0.31], y_2 = 0.95$   
 $\vdots$   
 $x_9 = [-0.09, -0.79], y_9 = -0.86$   
 $x_{10} = [-0.34, 0.87], y_{10}=?$

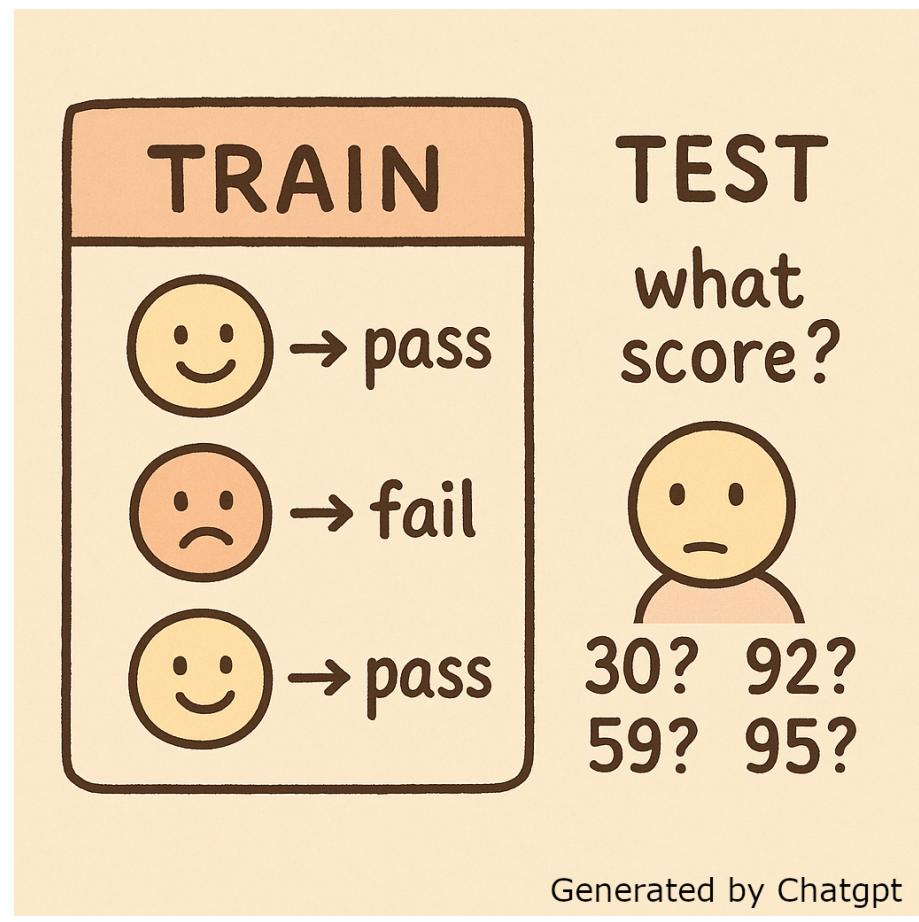


$y_{10} = 1.44 \text{ 😊}$

Source: <https://x.com/DimitrisPapail/status/1767924182724055415>

**Q:** Can we characterize this *easy-to-hard* task generalization in a tractable linear setting?

⇒ We focus on a classification to regression task shift problem.



## PROBLEM SETTING

- Data  $\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  and  $\mathbf{X} \in \mathbb{R}^{n \times d}$  with  $d \gg n$ .
- Ground-truth regressor  $\boldsymbol{\theta}^* \in \mathbb{R}^d$ .
- Classification labels  $\hat{\mathbf{y}} := (\text{sign}(\mathbf{x}_i^\top \boldsymbol{\theta}^*))_{i=1}^n$ .
- Minimum  $\ell_2$ -norm interpolator (MNI):

$$\hat{\boldsymbol{\theta}} = \mathbf{X}^\top (\mathbf{X}\mathbf{X}^\top)^{-1} \hat{\mathbf{y}}.$$

- Regression risk of classification MNI:

$$L(\hat{\boldsymbol{\theta}}) := \mathbb{E}_{\mathbf{x}}(\mathbf{x}^\top \hat{\boldsymbol{\theta}} - \mathbf{x}^\top \boldsymbol{\theta}^*)^2.$$

- Q:** Is the classification MNI consistent for the regression problem:  $\lim_{n,d \rightarrow \infty} L(\hat{\boldsymbol{\theta}}) = 0$ .

## ZERO-SHOT I: SPARSE SIGNAL

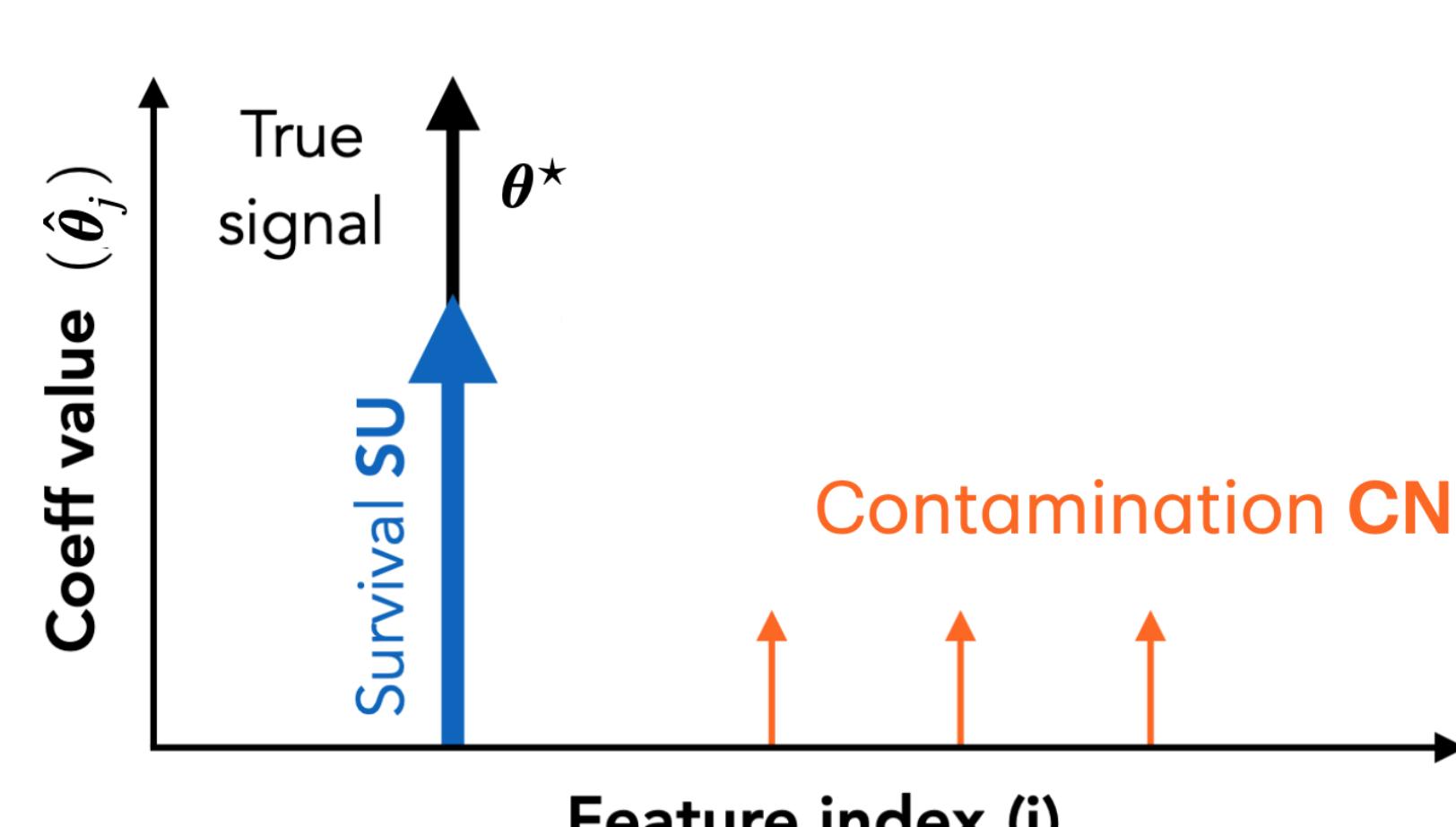
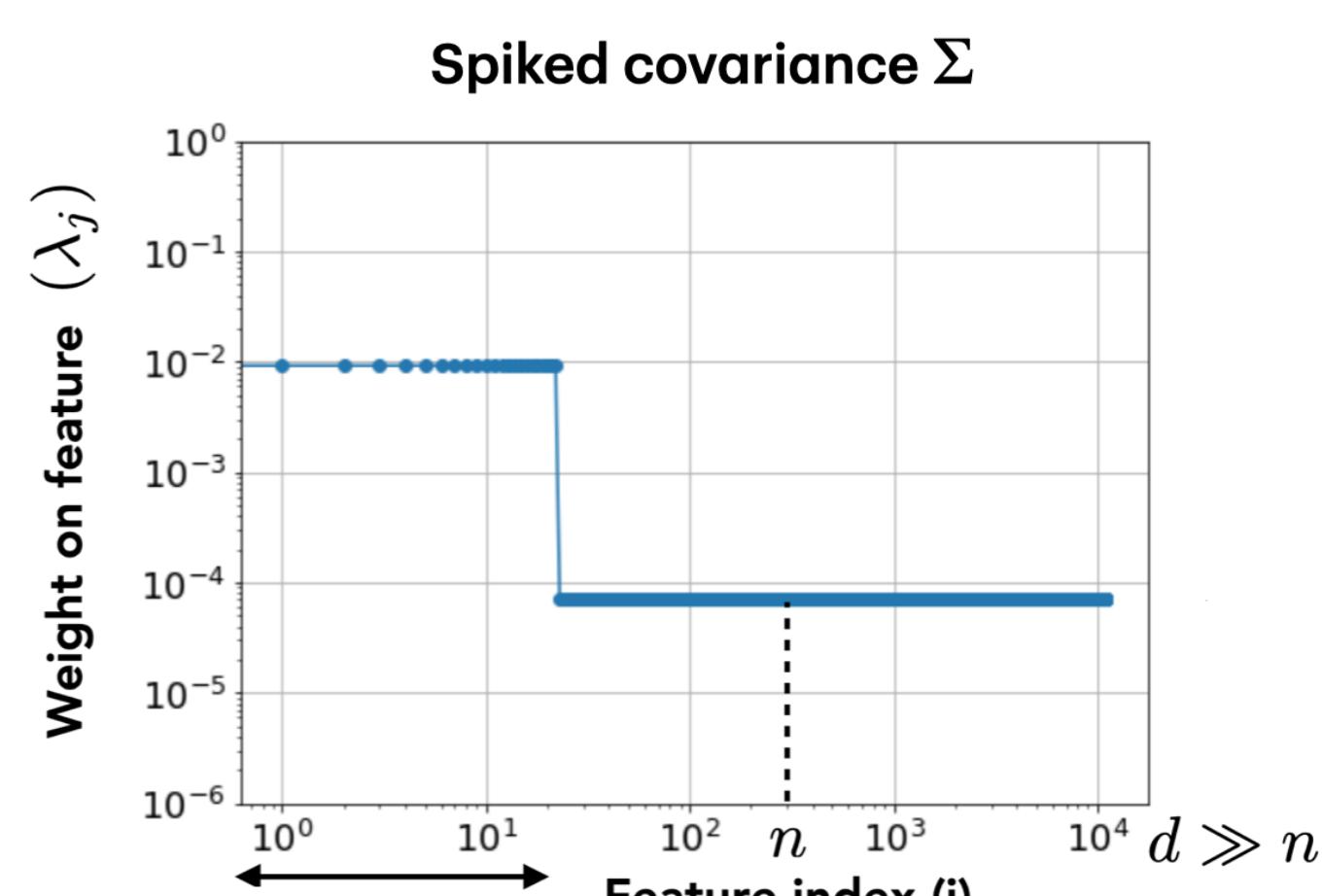


Image from Muthukumar et al., JMLR 2021

- Let  $\boldsymbol{\theta}^*$  be  $t$ -sparse and write  $\boldsymbol{\theta}_j^* := a_j \lambda_j^{-1/2}$  where  $\{\lambda_j\}_{j=1}^d$  is the spectrum of  $\boldsymbol{\Sigma}$ .
- Survival* and *contamination*:

$$\text{SU}_j := \frac{\hat{\theta}_j}{\theta_j^*} \quad \text{CN} := \sqrt{\sum_{j \in \mathcal{S}^c} \lambda_j \hat{\theta}_j^2}.$$

- Error decomposition:

$$L(\hat{\boldsymbol{\theta}}) = \sum_{j \in \mathcal{S}} a_j^2 (\text{SU}_j - 1)^2 + \text{CN}^2.$$

- Plug in our general  $\text{SU}_j$  and  $\text{CN}$  bounds for spiked covariance:

$$\lim_{n,d \rightarrow \infty} L(\hat{\boldsymbol{\theta}}) = \sum_{j \leq s, j \in \mathcal{S}} a_j^2 \left( \sqrt{\frac{2}{\pi \|\boldsymbol{\Sigma}^{\frac{1}{2}} \boldsymbol{\theta}^*\|_2^2}} - 1 \right)^2 + \sum_{j > s, j \in \mathcal{S}} a_j^2.$$

- Regression consistency iff  $\|\boldsymbol{\Sigma}^{\frac{1}{2}} \boldsymbol{\theta}^*\|_2^2 = \frac{2}{\pi}$  and  $a_j = 0$  for  $j > s, j \in \mathcal{S}$ .

## FINDINGS

- Zero-shot task shift is *impossible* for both sparse and random (dense) signals.
- Few-shot task shift is possible via *postprocessing* algorithm: no finetuning necessary.
- Structured attenuation of classification MNI: support recovery is guaranteed even when generalization is impossible!

## ZERO-SHOT II: RANDOM SIGNAL

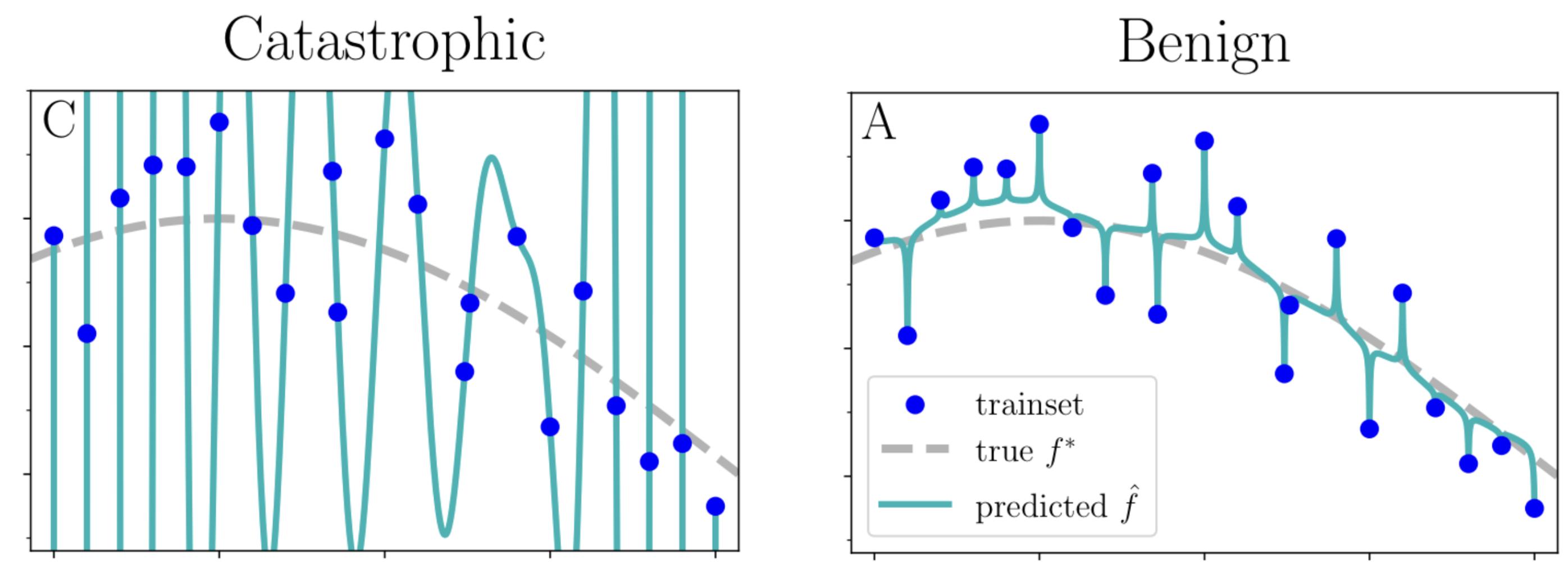


Image from Mallinar et al., NeurIPS 2022.

- Let  $\boldsymbol{\theta}^*$  be random such that  $\mathbb{E} \theta_j^{*2} \geq \sigma^2$  for all  $j$ .
- Error decomposition:

$$L(\hat{\boldsymbol{\theta}}) = \underbrace{L(\tilde{\boldsymbol{\theta}})}_{\text{bias}} + \underbrace{\mathbb{E}_{\mathbf{x}}(\mathbf{x}^\top \hat{\boldsymbol{\theta}} - \mathbf{x}^\top \tilde{\boldsymbol{\theta}})^2}_{\text{task shift error}}.$$

- Modified *benign overfitting* analysis:

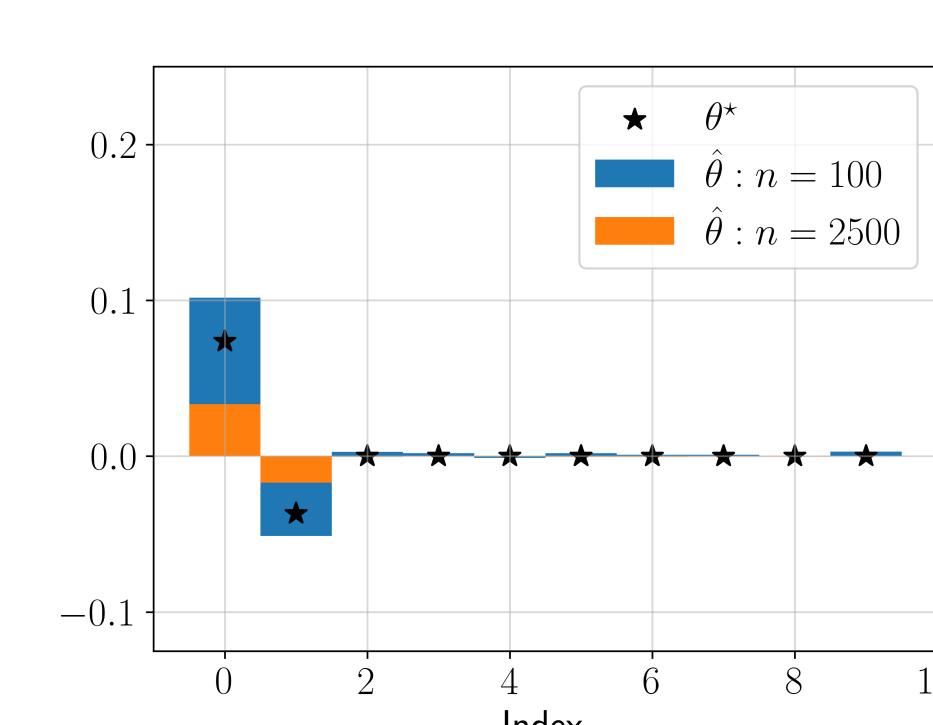
$$\mathbb{E}_{\boldsymbol{\theta}^*, \mathbf{x}} (\mathbf{x}^\top \hat{\boldsymbol{\theta}} - \mathbf{x}^\top \tilde{\boldsymbol{\theta}})^2 \gtrsim \sigma^2 \left( \sum_{j=1}^{k^*} \lambda_j + \frac{n}{R_{k^*}(\boldsymbol{\Sigma})} \sum_{j=k^*+1}^d \lambda_j \right).$$

- Together with standard bounds on the bias:

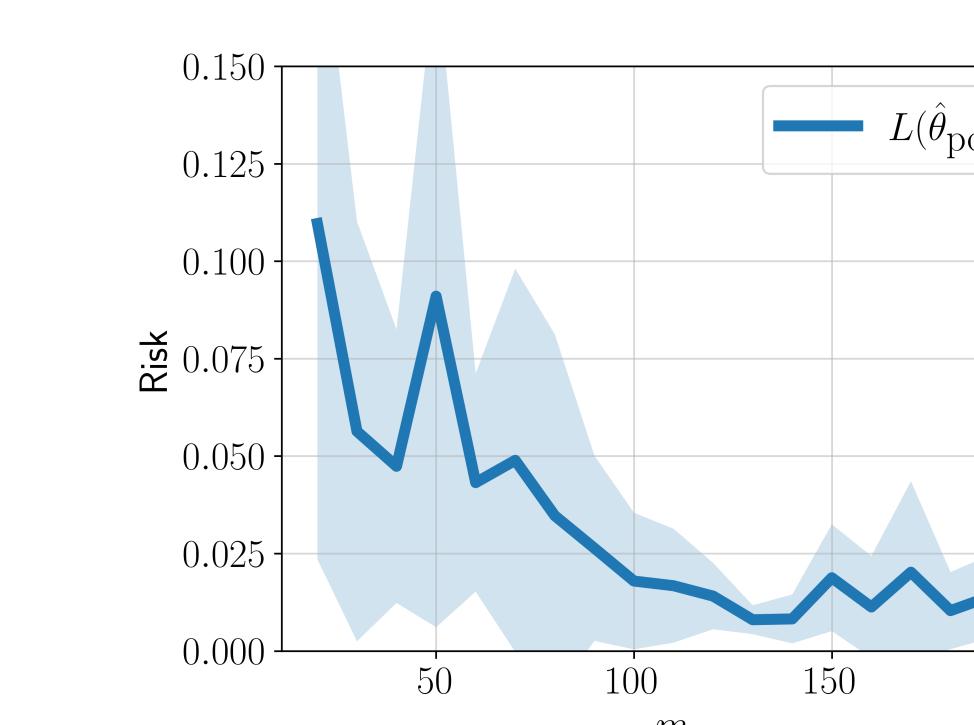
$$\begin{aligned} \lim_{n,d \rightarrow \infty} \mathbb{E}_{\boldsymbol{\theta}^*} L(\tilde{\boldsymbol{\theta}}) &\gtrsim \|\boldsymbol{\Sigma}^{\frac{1}{2}} \boldsymbol{\theta}^*\|_2^2 && \text{if } k^* = 0, \\ \lim_{n,d \rightarrow \infty} \mathbb{E}_{\boldsymbol{\theta}^*, \mathbf{x}} (\mathbf{x}^\top \hat{\boldsymbol{\theta}} - \mathbf{x}^\top \tilde{\boldsymbol{\theta}})^2 &\gtrsim \sigma^2 && \text{if } k^* > 0. \end{aligned}$$

## FEW-SHOT POSTPROCESSING

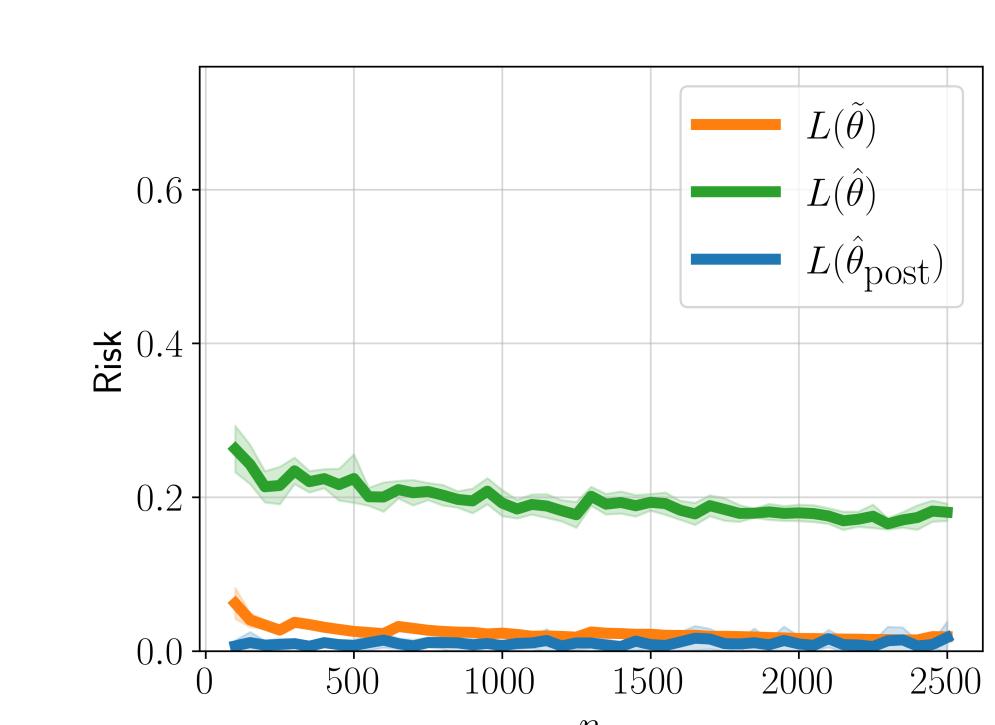
- Assume  $t$ -sparse  $\boldsymbol{\theta}^*$  and diagonal covariance  $\boldsymbol{\Sigma}$ .
- Step 1:* Choose the  $t$  largest elements of  $\hat{\boldsymbol{\theta}}$ . **Provably** converges to the support of  $\boldsymbol{\theta}^*$ .
- Step 2:* Recover the magnitude of support by few-shot least-squares on  $m$  regression data to obtain  $\mathcal{O}(\frac{t}{m})$  regression error.



Support recovery

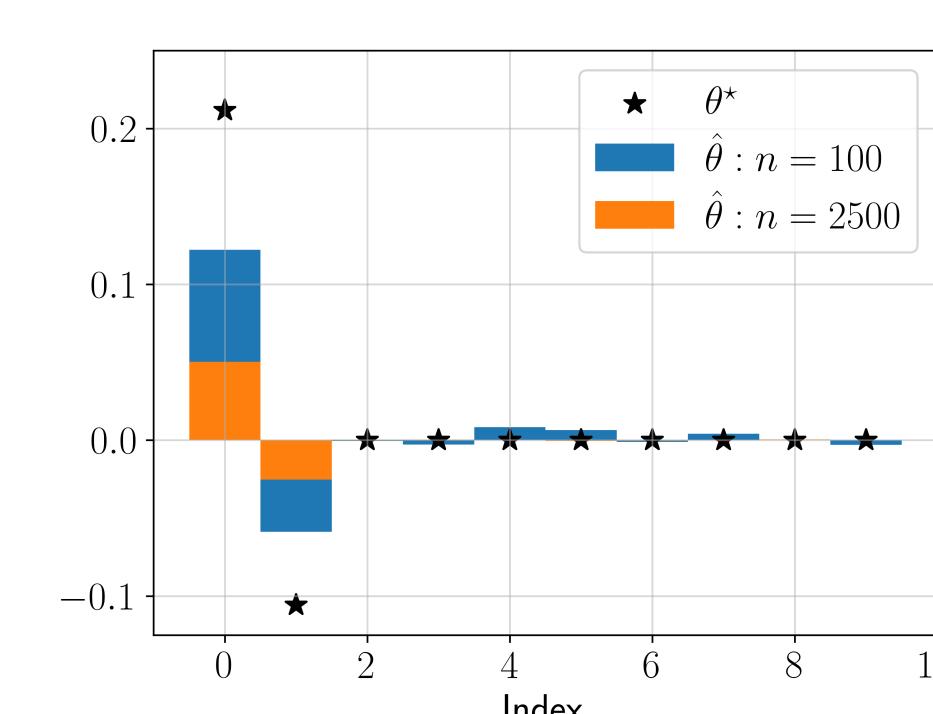


$m$ -shot least squares

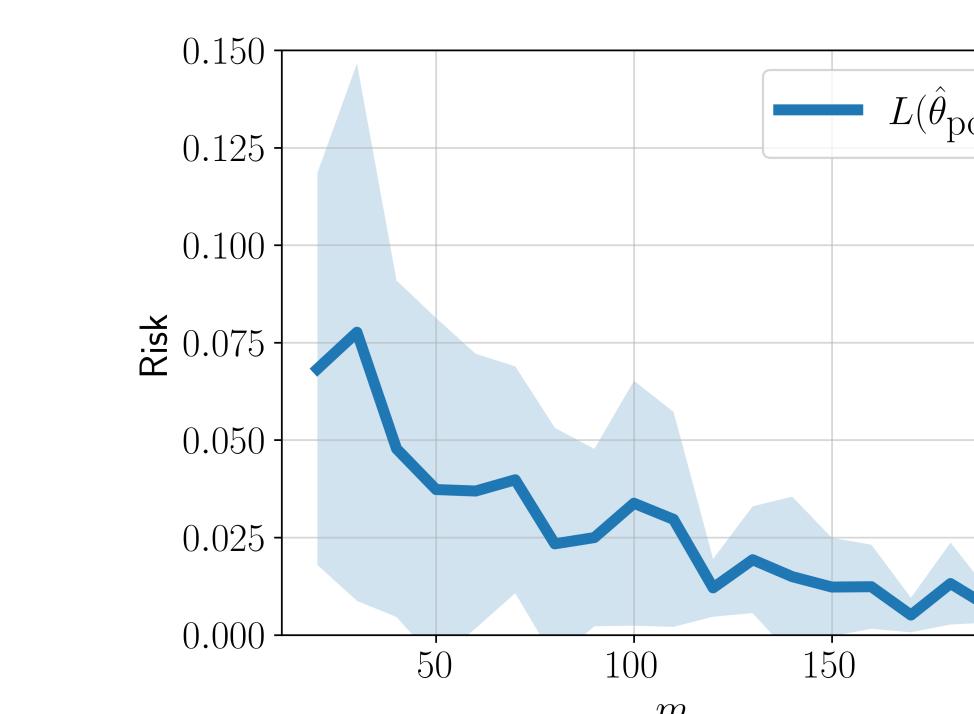


Regression risk

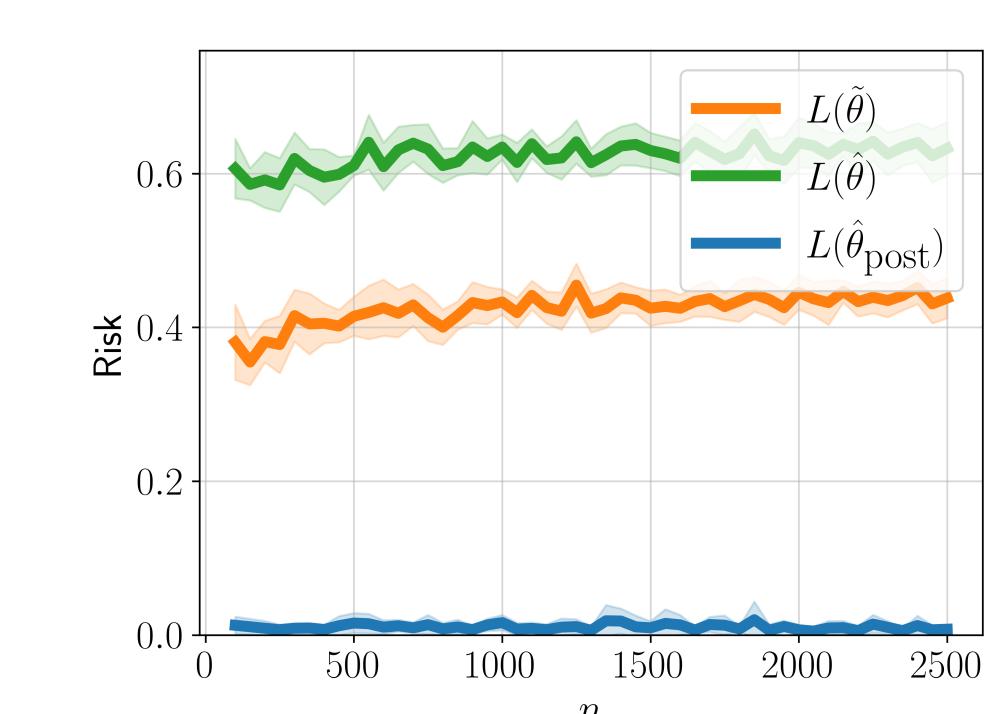
Task shift for spiked  $\boldsymbol{\Sigma}$  when regression generalizes.



Support recovery



$m$ -shot least squares



Regression risk

Task shift for spiked  $\boldsymbol{\Sigma}$  when regression does not generalize.