



Randomized tests for high-dimensional regression

Yue Li¹, Ilmun Kim², Yuting Wei¹
Carnegie Mellon University¹; University of Cambridge²



BACKGROUND AND PROBLEM

- Global testing problem.** $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ observations generated i.i.d. from a linear model

$$y_i = \langle x_i, \beta \rangle + \sigma z_i, \quad x_i \sim \mathcal{N}(\mathbf{0}, \Sigma)$$

for some unknown vector $\beta \in \mathbb{R}^p$.

Goal : $H_0 : \beta = \mathbf{0}$ versus $H_1 : \beta \neq \mathbf{0}$.

Setting: n/p of **constant** order and β non-sparse.

- Challenges.** Classical F -test does **not** work when $p \geq n$!
 - Impose specific structure assumptions [ZC11, CGZ18, JJ14, JM14a, ACCP11].
 - Can we find an *adaptive* and *general* approach?
- Solution: **random projection/sketching**
 - Widely studied in reducing computational cost and preserving privacy [BM01, LKR05, Sar06, PW17].
 - Statistical behaviors have been less studied.
- Close to our work: Kernel regression [YPW17], two-sample test [LJW11].

CONTRIBUTIONS

- Propose a **sketched F -test** which does not restrain the size of n, p .
- Provide a systematic way of selecting the projection dimension based on the underlying intrinsic dimension.
- Characterize situations where our test enjoys better power than existing competitors.

EXAMPLES: CHOICE OF k ($k \leq r$)

With SVD $\Sigma = U \Lambda U^\top$, define $\tilde{\beta} = U^\top \beta$. Then

- α -polynomial decay:** $\lambda_j \propto j^{-\alpha}$ with $\alpha > 1$ and homogeneous $\tilde{\beta}_i$. We have

$$r \lesssim (\log p)^{\frac{1}{\alpha-1}}.$$

- γ -exponential decay:** $\lambda_j \propto \exp(-j^\gamma)$ with $\gamma > 0$ and homogeneous $\tilde{\beta}_i$. We have

$$r \lesssim (\log \log p)^{\frac{1}{\gamma}}.$$

- structured coefficient:** $0 < c_1 \leq \tilde{\beta}_i \sqrt{i} \leq c_2$ and $\lambda_j \propto j^{-1}$. We have

$$r \lesssim (\log p)^3.$$

ALGORITHM

Algorithm **Sketched F -test**

Input: data matrix $\mathbf{X} \in \mathbb{R}^{n \times p}$, response vector $\mathbf{y} \in \mathbb{R}^n$, a sketching dimension $k < n$

Output: global testing result for the linear model.

Step 1: generate a sketching matrix $S_k \in \mathbb{R}^{p \times k}$ with i.i.d. $\mathcal{N}(0, 1)$ entries;

Step 2: compute the least square regression estimate $\hat{\beta}^S := (S_k^\top \mathbf{X}^\top \mathbf{X} S_k)^{-1} S_k^\top \mathbf{X}^\top \mathbf{y}$;

Step 3: calculate the sketched F -test statistic

$$F(S_k) := \frac{\mathbf{y}^\top \mathbf{X} S_k \hat{\beta}^S / k}{\|\mathbf{y} - \mathbf{X} S_k \hat{\beta}^S\|_2^2 / (n - k)};$$

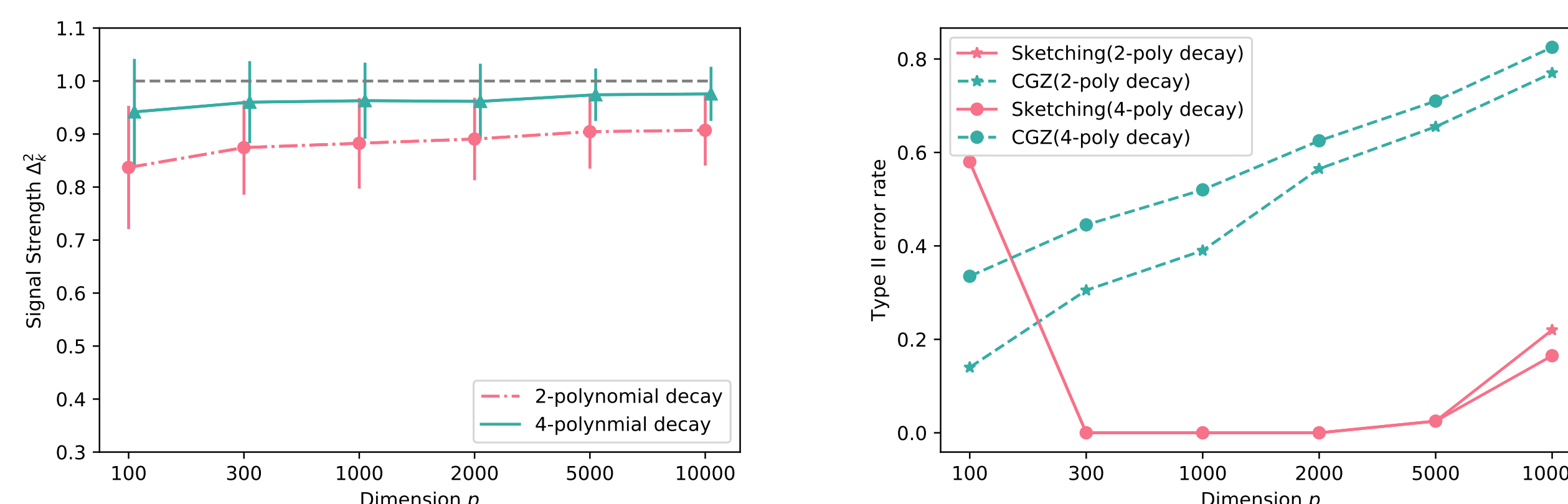
Step 4: if $F(S_k) \geq q_{\alpha, k, n-k}$, reject H_0 ; otherwise accept H_0 .

NUMERICAL RESULTS

Numerical comparisons with CGZ18, ZC11 with decaying patterns: slow-decay (log) and fast-decay (polynomial) with $k = \lfloor n/2 \rfloor$.

$\ \Sigma\ _F = 100$		$H_0: \ \beta\ _2 = 0$	$\ \beta\ _2 = 1$	$\ \beta\ _2 = 5$
slow-decay	Sketching	3.2%	1.4%	0.0%
	CGZ	6.0%	5.4%	4.6%
	ZC	2.1%	16.8%	0.6%
fast-decay	Sketching	4.0%	1.4%	2.4%
	CGZ	6.2%	10.4%	12.4%
	ZC	4.2%	14.7%	6.3%

Asymptotic Behavior. With the structure design and optimal choice of k , we plot the signal strength and error v.s. feature dimension p .



REFERENCES

- [1] Y. Yang, M. Pilanci, and M. J. Wainwright, "Randomized sketches for kernels: Fast and optimal nonparametric regression," *The Annals of Statistics*, 2017.
- [2] M. Lopes, L. Jacob, and M. J. Wainwright, "A more powerful two-sample test in high dimensions using random projection," in *Advances in Neural Information Processing Systems*, 2011, pp. 1206–1214.
- [3] H. Cui, W. Guo, and W. Zhong, "Test for high-dimensional regression coefficients using refitted cross-validation variance estimation," *The Annals of Statistics*, 2018.

CHARACTERIZATION OF POWER

The power of the proposed test is determined by

$$\Delta_k^2 := \beta^\top \Sigma S_k (S_k^\top \Sigma S_k)^{-1} S_k^\top \Sigma \beta.$$

Theorem 1. When the data and noise both follow Gaussian distributions and are independent to each other, the power of the proposed test satisfies

$$\Psi_n^S(S_k) - \Phi \left(-z_\alpha + \sqrt{\frac{(1-\rho)n}{2\rho}} \frac{\Delta_k^2}{\sigma^2} \right) \rightarrow 0,$$

where $\rho_n = k/n \rightarrow \rho$.

Comparisons. Assuming the normalized vector $\Sigma^{1/2} \beta / \|\Sigma^{1/2} \beta\|_2$ is uniformly distributed on the p -dimensional unit sphere independent of S_k , the proposed test has higher power than [CGZ18] *w.h.p.* if

$$\frac{4}{\sqrt{\rho(1-\rho)}} \frac{\text{tr}(\Sigma)}{\sqrt{\text{tr}(\Sigma^2)}} \frac{1}{\sqrt{n}} \leq 1.$$

Optimal choice of k . In this case, choose $k = \lfloor n/2 \rfloor$.

OPTIMALITY UNDER STRUCTURE DESIGN

- The model class with **intrinsic dimension** up to r is defined as, with some $\eta = o(1)$ and $\lambda_1 \geq \dots \geq \lambda_p$ being eigenvalues of Σ ,

$$\sum_{i=r+1}^p \lambda_i \leq \eta \sum_{i=1}^p \lambda_i \quad \text{and} \quad r \lambda_{r+1} \leq \eta \sum_{i=1}^p \lambda_i.$$

- When we choose k proportional to the intrinsic dimension r of the model, we can **fully preserve the signal** *w.h.p.*!

Theorem 2. Within the r -intrinsic dimensional model class, the proposed test is minimax rate optimal with radius

$$\epsilon_n^2 = \frac{r^{1/2}}{n},$$

and the upper bound is reached by choosing sketching dimension $k = O(r)$.