

The All-or-Nothing Phenomenon in Sparse Tensor PCA

Jonathan Niles-Weed^{1,2} & Ilias Zadik²

¹Courant Institute of Mathematical Sciences, NYU & ²Center for Data Science, NYU

High-dimensional statistical models undergo phase transitions

small changes in the parameters of the model imply large changes in the model (e.g. large changes in terms of recovery error).

A canonical example

Principal Component Analysis [BBP'05];

For $\beta \sim \text{Unif}(S^{p-1})$ and $W_{i,j} = W_{j,i} \sim N(0, 1)$ (Wigner),

$$Y = \sqrt{\lambda p} \beta \beta^\top + W.$$

$$\lim_{p \rightarrow +\infty} \mathbb{E} \left[\|\beta - \mathbb{E}[\beta|Y]\|_2^2 \right] = \begin{cases} < 1, & \text{if } \lambda > 1 \\ 1, & \text{if } \lambda < 1. \end{cases}$$

Under (sublinear) sparsity: Much sharper transitions

Sparse Linear Regression

For $\beta \sim \text{Unif}(\{v \in \{0, \frac{1}{\sqrt{k}}\}^p, \|v\|_0 = k\})$, $X_{i,j}, W_i \sim N(0, 1)$,

$$Y = X\beta + \sigma W \in \mathbb{R}^n.$$

[Gamarnik, Zadik '17], [Reeves, Xu, Zadik '19] “The All-or-Nothing Phenomenon”

For some critical sample size n^* and $k = o(\sqrt{p})$, if $\text{SNR} = k/\sigma^2 \rightarrow +\infty$,

$$\lim_{p \rightarrow +\infty} \mathbb{E} \left[\|\beta - \mathbb{E}[\beta|Y, X]\|_2^2 \right] = \begin{cases} 0, & \text{if } n > n^* \\ 1, & \text{if } n < n^*. \end{cases}$$

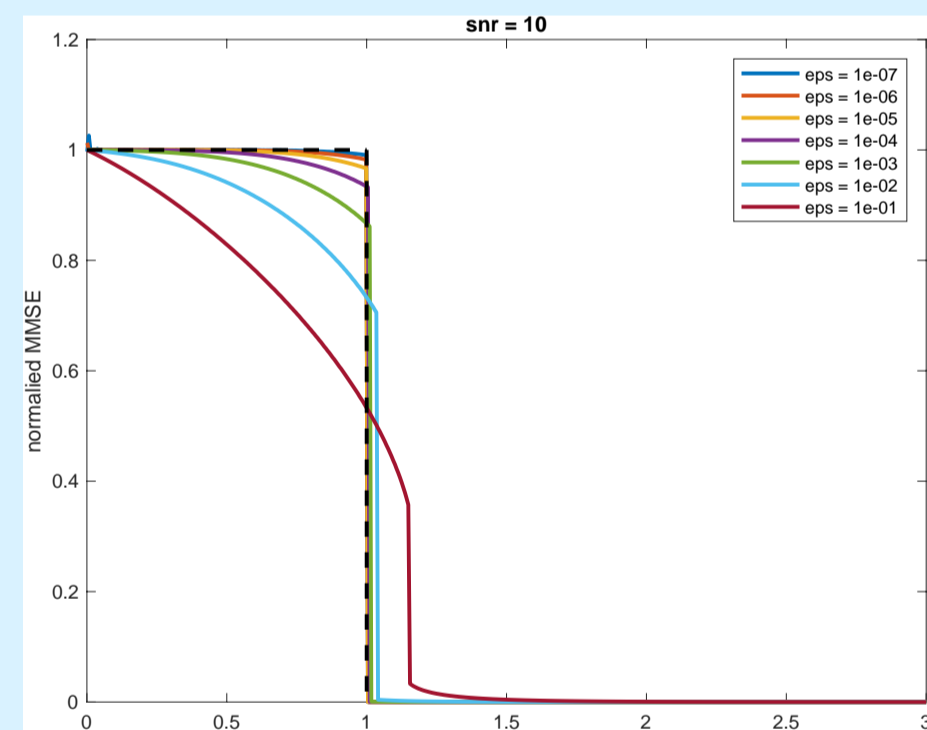


Figure from [Reeves, Xu, Zadik '19], MMSE vs n/n^* as $k/p = \epsilon \rightarrow 0$.

Sparse (tensor) PCA

For $\beta \sim \text{Unif}(\{v \in \{0, \frac{1}{\sqrt{k}}\}^p, \|v\|_0 = k\})$, $W_{i,j} \sim N(0, 1)$ and $d \geq 2$,

$$Y = \sqrt{2\lambda k \log \frac{p}{k}} \beta^{\otimes d} + W.$$

[Barbier, Macris, Rush '20] If $d = 2$ (the matrix case), $p^{\frac{12}{13}} \ll k = o(p)$ then

$$\lim_{p \rightarrow +\infty} \mathbb{E} \left[\|\beta - \mathbb{E}[\beta|Y]\|_2^2 \right] = \begin{cases} 0, & \text{if } \lambda > 1 \\ 1, & \text{if } \lambda < 1. \end{cases}$$

Bernoulli Group Testing

For $\beta \sim \text{Unif}(\{v \in \{0, \frac{1}{\sqrt{k}}\}^p, \|v\|_0 = k\})$, $X_{i,j} \sim \text{Bern}(\ln 2/k)$

$$Y = (\mathbf{1}(\langle X, v \rangle \geq 1))_{i=1, \dots, n} \in \{0, 1\}^n.$$

[Truong, Aldridge, Scarlett '20] For some critical sample size n^* and $k = p^{o(1)}$,

$$\lim_{p \rightarrow +\infty} \mathbb{E} \left[\|\beta - \mathbb{E}[\beta|Y, X]\|_2^2 \right] = \begin{cases} 0, & \text{if } n > n^* \\ 1, & \text{if } n < n^*. \end{cases}$$

Main Contributions

- Systematic study of a simple but generic inference model (the Gaussian Additive Model)
- A simple necessary and sufficient condition for all-or-nothing to hold.
- Sparse tensor PCA exhibits all-or-nothing for all $d \geq 2$ and $k = o(p)$.

The Gaussian Additive Model

For arbitrary discrete $S = S_p \subset S^{p-1} = \{v \in \mathbb{R}^p : \|v\|_2 = 1\}$, let $\beta \sim P_p := \text{Unif}(S_p)$ and $W_i \sim N(0, 1)$,

$$Y_\lambda = \sqrt{\lambda} \beta + W.$$

Subcases:

- Sparse Tensor PCA: $S = \{v^{\otimes d} : v \in \{0, \frac{1}{\sqrt{k}}\}^p, \|v\|_0 = k\}$.
- Submatrix Localization: $S = \{(\mathbf{1}_{\sigma(i)=\sigma(j)} - \frac{1}{k})_{i,j \in [n]} : \sigma : [n] \rightarrow [k] \text{ “balanced”}\}$. “Gaussian version of stochastic block model”

Definition 1 A sequence of priors $(P_p)_{p \in \mathbb{N}}$ satisfies the all-or-nothing phenomenon holds if for some critical $\lambda_c = \lambda_c(S, p)$:

$$\lim_{p \rightarrow \infty} \mathbb{E} \|\beta - \mathbb{E}[\beta|Y_\lambda]\|_2^2 = \begin{cases} 0 & \text{if } \lambda > \lambda_c \\ 1 & \text{if } \lambda < \lambda_c. \end{cases}$$

An equivalent condition

Definition 2 We denote by D the Kullback-Leibler divergence; given two random variables Z_1, Z_2 with the law of Z_1, P_1 , absolutely continuous to the law of Z_2, P_2 ,

$$D(Z_1, Z_2) = \mathbb{E} \left[\frac{dP_1(Z_2)}{dP_2(Z_2)} \log \left(\frac{dP_1(Z_2)}{dP_2(Z_2)} \right) \right].$$

Theorem 1 Suppose $|S| = |S_p| \rightarrow +\infty$ and S is “sufficiently spread out” as $p \rightarrow +\infty$. A sequence $\{P_p\}$ satisfies the all-or-nothing phenomenon for some λ_c if and only if

$$\lim_{p \rightarrow +\infty} \frac{1}{\log |S_p|} D(Y_{2 \log |S_p|}, Y_0) = 0$$

and $\lambda_c = (1 + o(1))2 \log |S_p|$.

*Sufficient spread is necessary for “all” recovery to happen.

Key Intuition

- I-MMSE formula:

$$\frac{d}{d\lambda} D(Y_\lambda, Y_0) = \frac{1}{2} - \frac{1}{2} \mathbb{E} \|\beta - \mathbb{E}[\beta|Y_\lambda]\|_2^2$$

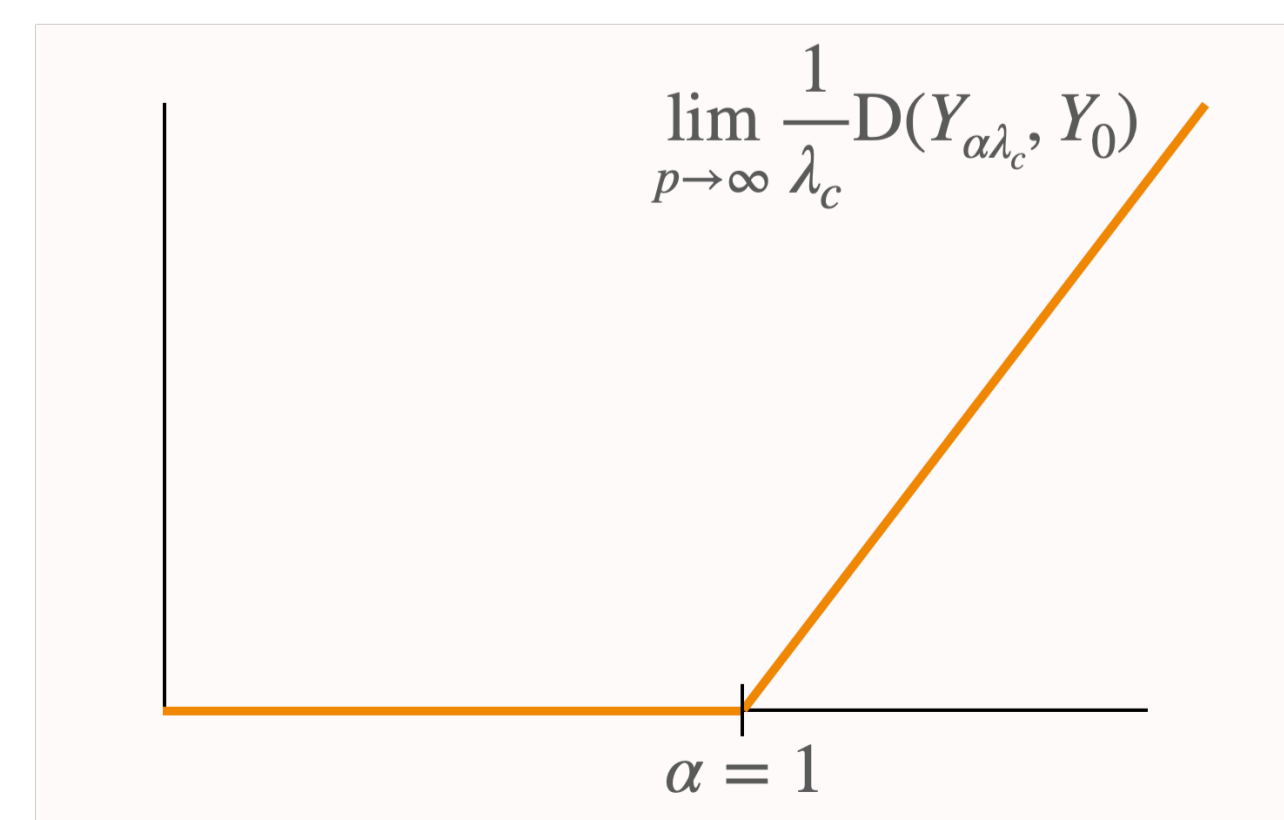
An abrupt change in the slope of KL!

- All-or-nothing phenomenon at λ_c if and only if

$$\lim_{p \rightarrow \infty} \frac{1}{\lambda_c} D(Y_{\alpha \lambda_c}, Y_0) = \frac{1}{2}(\alpha - 1)_+ \quad \forall \alpha \geq 0.$$

- Since limiting KL is Lipschitz, convex, and nonnegative, evaluating the limit at a $\alpha = 1$ and $\alpha \rightarrow \infty$ gives information about the whole function.

- For “large” $\alpha > 1$, $\frac{1}{\lambda_c} D(Y_{\alpha \lambda_c}, Y_0) = \frac{\alpha}{2} - \frac{I(\beta; Y_{\alpha \lambda_c})}{\lambda_c} \approx \frac{\alpha}{2} - \frac{\log |S|}{\lambda_c}$. Hence $\lambda_c \approx 2 \log |S|$.



“Near-orthogonality” implies all-or-nothing

Corollary 1 (Informal) Suppose $S = S_p$ consist of “nearly-orthogonal” vectors, rigorously for any $t > 0$ and large p ,

$$P_p^{\otimes 2}(\langle X, X' \rangle \geq t) \leq |S|^{-\frac{2t}{t+1}}$$

then all-or-nothing phenomenon always holds at $\lambda_c = 2 \log |S_p|$.

Proof by conditional second moment method [Banks et al. '18, Perry et al. '20].

Easy implication: Any uniform distribution over discrete subsets of orthogonal vectors satisfies the all-or-nothing phenomenon!

Application: Sparse Tensor PCA

Theorem 2 For any $d \geq 2$ and $k = o(p)$, the sparse tensor PCA model

$$Y = \sqrt{2\lambda k \log \left(\frac{p}{k} \right)} \beta^{\otimes d} + W, \quad \beta \sim \tilde{P}_p$$

with $\tilde{P}_p = \text{Unif}(\{v \in \{0, \frac{1}{\sqrt{k}}\}^p : \|v\|_0 = k\})$ exhibits the all-or-nothing phenomenon:

$$\lim_{p \rightarrow \infty} \mathbb{E} \|\beta^{\otimes d} - \mathbb{E}[\beta^{\otimes d}|Y]\|_2^2 = \begin{cases} 1 & \text{if } \lambda < 1 \\ 0 & \text{if } \lambda > 1. \end{cases}$$

Proof idea

- If β and β' are independent draws from $\tilde{P}_p = \text{Unif}(\{v \in \{0, \frac{1}{\sqrt{k}}\}^p : \|v\|_0 = k\})$, then the overlap $\langle \beta, \beta' \rangle$ is a rescaled Hypergeometric random variable.

- When $k = o(p)$, this distribution satisfies $\tilde{P}_p^{\otimes 2}(\langle \beta, \beta' \rangle \geq t) \leq |S|^{-t+o(1)}$.
- Therefore, for any $d \geq 2$,

$$\tilde{P}_p^{\otimes 2}(\langle \beta^{\otimes d}, (\beta')^{\otimes d} \rangle \geq t) \leq |S|^{-t/d+o(1)} \leq |S|^{-\frac{2t}{t+1}}.$$

Result follows from Corollary 1.

Open questions

- How does Theorem 1 extend to more general models (e.g. GLMs)?
- All-or-nothing for polynomial-time estimators?

References

- [BBP05] Jinho Baik, Gérard Ben Arous, and Sandrine Péché. Phase transition of the largest eigenvalue for nonnull complex sample covariance matrices. *The Annals of Probability*, 33(5):1643–1697, 2005.
- [Barbier, Macris, Rush '20] Jean Barbier, Nicolas Macris, and Cynthia Rush. All-or-nothing statistical and computational phase transitions in sparse spiked matrix estimation. arXiv:2006.07971, 2020.
- [Banks et al. '18] Jess Banks, Cristopher Moore, Roman Vershynin, Nicolas Verzelen, and Jiaming Xu. Information-theoretic bounds and phase transitions in clustering, sparse PCA, and submatrix localization. *IEEE Trans. Inform. Theory*, 64(7):4872–4994, 2018.
- [Gamarnik, Zadik '17] David Gamarnik and Ilias Zadik. High dimensional linear regression with binary coefficients: Mean squared error and a phase transition. *Conference on Learning Theory (COLT)*, 2017.
- [Perry et al. '20] Amelia Perry, Alexander S. Wein, and Afonso S. Bandeira. Statistical limits of spiked tensor models. *Ann. Inst. Henri Poincaré Probab. Stat.*, 56(1):230–264, 2020.
- [Reeves, Xu, Zadik '19] Galen Reeves, Jiaming Xu, and Ilias Zadik. The all-or-nothing phenomenon in sparse linear regression. *Conference on Learning Theory (COLT)*, 2019.
- [Truong, Aldridge, Scarlett '20] Lan V. Truong, Matthew Aldridge and Jonathan Scarlett On the All-Or-Nothing Behavior of Bernoulli Group Testing arXiv preprint, 2020.