

COLUMN THRESHOLDING FOR SPARSE SPIKED WIGNER MODELS: IMPROVED SIGNAL STRENGTH REQUIREMENTS

Anonymous authors

Paper under double-blind review

ABSTRACT

We study the sparse spiked Wigner model, where the goal is to recover an s -sparse unit vector $\mathbf{u} \in \mathbb{R}^d$ from a noisy matrix observation $\mathbf{Y} = \beta \mathbf{u} \mathbf{u}^\top + \mathbf{W}$. While the information-theoretic threshold is $\beta = \tilde{\Omega}(\sqrt{s})$, existing polynomial-time algorithms require $\beta = \tilde{\Omega}(s)$, yielding a substantial computational-statistical gap. We propose a column thresholding method that attains the $\tilde{\Omega}(\sqrt{s})$ scaling for estimation and support recovery under the non-uniformity condition $\|\mathbf{u}\|_\infty = \Omega(1)$. Building on this initializer, we further develop a truncated power method that iteratively refines the estimate with provable linear convergence. Experiments validate our theoretical guarantees and demonstrate superior performance in estimation accuracy, support recovery, and computational efficiency.

1 INTRODUCTION

We study the sparse spiked Wigner model (Deshpande & Montanari, 2014; Lesieur et al., 2015), which addresses the problem of recovering a sparse vector \mathbf{u} from a noisy matrix $\mathbf{Y} \in \mathbb{R}^{d \times d}$:

$$\mathbf{Y} = \beta \mathbf{u} \mathbf{u}^\top + \mathbf{W}, \quad (1)$$

where $\mathbf{u} \in \mathbb{R}^d$ is an unknown s -sparse unit vector, $\beta > 0$ denotes the signal strength, and $\mathbf{W} \sim \text{GOE}(d)$ is distributed as the Gaussian orthogonal ensemble, i.e., $\mathbf{W} = \frac{1}{\sqrt{2}}(\mathbf{A} + \mathbf{A}^\top)$ with \mathbf{A} having i.i.d. $\mathcal{N}(0, 1)$ entries. This model captures fundamental inference problems involving pairwise measurements, including Gaussian variants of community detection (Deshpande et al., 2016) and $\mathbb{Z}/2$ synchronization (Javanmard et al., 2016).

Information-theoretic limits The fundamental limits for sparse PCA in this model are well-understood. For support recovery and estimation up to constant error, the information-theoretic lower bound on signal strength is $\beta = \tilde{\Omega}(\sqrt{s})$ (Banks et al., 2018; Perry et al., 2018; 2020), which is achievable through exhaustive search and other exponential-time procedures, but remains out of reach for polynomial-time algorithms.

Polynomial-time algorithms Existing polynomial-time approaches fall into two main categories, each with fundamental limitations:

Spectral methods. The vanilla spectral algorithm computes the leading eigenpair of \mathbf{Y} , while spectral projection (Brennan et al., 2018) additionally projects this eigenvector onto the set of unit s -sparse vectors. Both methods incur $O(d^3)$ computational cost and require signal strength $\beta = \tilde{\Omega}(\sqrt{d})$ for successful recovery (Baik et al., 2005; Péché, 2006; Féral & Péché, 2007; Paul, 2007; Capitaine et al., 2009; Benaych-Georges & Nadakuditi, 2011; Brennan et al., 2018). This requirement is tight—spectral methods provably fail when $\beta = \tilde{O}(\sqrt{d})$ (Montanari et al., 2015).

Thresholding methods. Diagonal thresholding (Johnstone & Lu, 2009) identifies the support by selecting the s largest diagonal entries of \mathbf{Y} , then estimates \mathbf{u} via the leading eigenvector of the corresponding $s \times s$ submatrix, achieving $O(d \log d + s^3)$ complexity. Covariance thresholding first applies soft (Deshpande & Montanari, 2016) or hard (Krauthgamer et al., 2015) thresholding to all entries of \mathbf{Y} before eigendecomposition. While computationally more efficient than spectral methods when $s \ll d$, these approaches uniformly require $\beta = \tilde{\Omega}(s)$ for successful recovery (Hopkins et al., 2017; Brennan et al., 2018; Choo & d’Orsi, 2021).

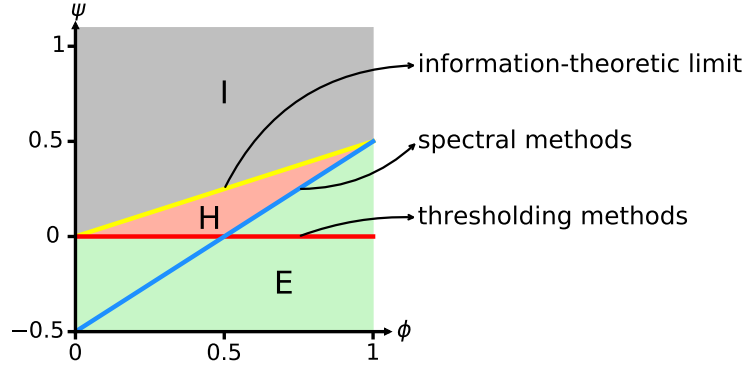


Figure 1: Recovery regimes in the sparse spiked Wigner model (Brennan et al., 2018), with $s = \tilde{\Theta}(d^\phi)$ and $s/\beta = \tilde{\Theta}(d^\psi)$. Regions I, H, and E correspond to the information-theoretically impossible, computationally hard, and polynomial-time tractable regimes, respectively. The yellow line indicates the information-theoretic threshold, while blue and orange lines show computational boundaries achieved by spectral and thresholding methods, respectively.

Both the spectral and thresholding methods discussed above fall under a common spectral paradigm. In this framework, matrix perturbation theory shows that the empirical matrix concentrates around its expectation, which can be viewed as a low-rank “signal” perturbation of a simple baseline (such as the identity); classical eigenvector perturbation bounds then imply that the leading eigenvector aligns with the true spike once the signal is sufficiently strong. This spectral approach is often used as an initialization, followed by a refinement stage, in problems such as phase retrieval (Candes et al., 2015), matrix sensing (Tu et al., 2016), and blind deconvolution (Ma et al., 2018), and such two-stage algorithms are now standard for tackling nonconvex optimization problems (Chi et al., 2019). In contrast, for the sparse spiked Wigner model, the analogous top-eigenvector method has been analyzed and is known to require the suboptimal signal strength $\tilde{\Omega}(\sqrt{d})$.

The computational-statistical gap and three regimes A significant gap exists between the information-theoretic threshold of $\tilde{\Omega}(\sqrt{s})$ and the $\tilde{\Omega}(s)$ signal strength required by polynomial-time algorithms. Brennan et al. (2018) characterized this phenomenon through three regimes. Parameterizing $s = \tilde{\Theta}(d^\phi)$ and $s/\beta = \tilde{\Theta}(d^\psi)$ with $\phi \in [0, 1]$, these regimes are:

- **Regime I (Impossible):** When $\beta = \tilde{O}(\sqrt{s})$ (i.e., $\psi > \phi/2$), recovery is information-theoretically impossible for any algorithm (Banks et al., 2018; Perry et al., 2018; 2020; Barbier et al., 2016; Lelarge & Miolane, 2017).
- **Regime H (Hard):** When $0 < \psi \leq \phi/2$ and $\psi > \phi - 1/2$, the problem is information-theoretically solvable but conjectured to be computationally intractable in polynomial time under the planted clique hypothesis (Brennan et al., 2018).
- **Regime E (Easy):** When $\beta = \tilde{\Omega}(s)$ or $\beta = \tilde{\Omega}(\sqrt{d})$ (i.e., $\psi \leq 0$ or $\psi \leq \phi - 1/2$), polynomial-time algorithms exist. Thresholding methods succeed when $\beta = \tilde{\Omega}(s)$ ($\psi \leq 0$) (Hopkins et al., 2017; Brennan et al., 2018; Choo & d’Orsi, 2021), while spectral methods succeed when $\beta = \tilde{\Omega}(\sqrt{d})$ ($\psi \leq \phi - 1/2$) (Benaych-Georges & Nadakuditi, 2011; Brennan et al., 2018).

Assuming the planted clique conjecture holds, the reductions from the planted clique problem show that any polynomial-time algorithms cannot recover a *uniform* spiked vector in the hard regime of Figure 1, where the signal strength lies between \sqrt{s} and s . Thus, if one restricts to uniform amplitudes, the statistical-computational gap in this regime is believed to be fundamental and cannot be closed (assuming the conjecture).

By contrast, the computational complexity and the associated phase transitions are much less understood when we consider different *classes of spikes* beyond the uniform case. Our goal is not to claim progress in the classical hard regime for uniform spikes, but rather to clarify what can be achieved once we move beyond the uniform setting. We identify a specific class of spikes, defined by an ℓ_∞ lower bound on \mathbf{u} , in which the uniform vector is ruled out and recovery at the \sqrt{s} signal strength becomes possible.

Under this ℓ_∞ condition, we prove that the column thresholding method succeeds at signal strength $\beta = \Omega(\sqrt{s})$, thereby providing a polynomial-time algorithm that operates in the hard regime for this class of non-uniform spikes, where the planted clique lower bound does not apply.

Therefore, our paper does not resolve the planted-clique–hard regime for uniform spikes. Instead, it identifies a different class of spikes, characterized by the ℓ_∞ condition, for which recovery at the \sqrt{s} signal strength is provably achievable in polynomial time. We have revised the manuscript to make this distinction clearer.

Our contributions In this paper, we propose two algorithms for sparse spike recovery: a polynomial-time column-thresholding method and a truncated power method (TPM) that uses the column-thresholding output as an initialization to further refine the estimate. The column-thresholding procedure still fits within the general spectral framework, but it is based on a different statistic: rather than aggregating information through all diagonal entries or the global top eigenvector, we first select a data-driven column by locating the largest observed diagonal entry, and then apply entrywise thresholding to that column. This construction yields a stronger separation between in-support and out-of-support indices, which in turn leads to an improved scaling in the signal strength required for successful recovery. Our main contributions are:

- We prove that column thresholding achieves the $\tilde{\Omega}(\sqrt{s})$ signal-strength scaling for both estimation and support recovery, under the assumption that $\|\mathbf{u}\|_\infty = \Omega(1)$. This assumption is not merely technical: it is essential for attaining the $\tilde{\Omega}(\sqrt{s})$ rate and explicitly rules out the uniform spike case in which planted-clique–based hardness results apply. Our work does not resolve the planted-clique–hard regime for uniform spikes. Rather, by imposing this ℓ_∞ condition and analyzing the column-thresholding algorithm, we identify a concrete class of non-uniform spikes, lying outside the reach of existing planted-clique reductions, for which recovery at signal strength $\tilde{\Omega}(\sqrt{s})$ is provably achievable in polynomial time.

Additionally, the condition $\|\mathbf{u}\|_\infty = \Omega(1)$ naturally covers power-law decaying signals (Jagatap & Hegde, 2019; Chen et al., 2015). When this condition fails, our bound degrades to $\tilde{\Omega}(s)$, matching existing thresholding methods. Conversely, to our knowledge, existing thresholding methods cannot achieve the $\tilde{\Omega}(\sqrt{s})$ rate even when $\|\mathbf{u}\|_\infty = \Omega(1)$ holds.

- We demonstrate that using column thresholding as initialization for truncated power iteration yields a two-stage algorithm with both rigorous guarantees and strong practical performance. While Yuan & Zhang (2013) established convergence theory conditional on a correlation condition, they did not provide a concrete initialization procedure. Our work fills this gap in the sparse spiked Wigner model by explicitly constructing an initialization that satisfies their theoretical requirements at the optimal signal level, enabling the refinement framework to operate in this previously inaccessible regime.
- Experiments validate our theory and demonstrate strong empirical performance. The column–thresholding method matches the predicted signal–strength scaling, while TPM achieves superior estimation accuracy and exact support recovery compared to baseline methods, all with competitive computational efficiency.

It is also natural to extend our methods and analysis to related models. For instance, in the symmetric two-cluster sparse Gaussian mixture model (Pesce et al., 2022; Löffler et al., 2022), the expected sample covariance exhibits the same structural form as in the sparse spiked Wigner model. This analogy allows both diagonal thresholding and our column-thresholding procedure to be used for support estimation of the sparse cluster mean, after which standard eigenvector-based methods can be applied to recover the cluster mean itself.

Notations: We use $f(n) = O(g(n))$ when $f(n) \leq c_1 g(n)$, $f(n) = \Omega(g(n))$ when $f(n) \geq c_2 g(n)$, and $f(n) = \Theta(g(n))$ when both hold, for some constants $c_1, c_2 > 0$. We use $\tilde{O}, \tilde{\Omega}, \tilde{\Theta}$ to denote the logarithm-suppressing variants of O, Ω, Θ that hide polylogarithmic factors in d . For vector \mathbf{a} , a_i denotes the i -th element, $\|\mathbf{a}\|_0$ counts nonzero entries, and $\|\mathbf{a}\|_2, \|\mathbf{a}\|_\infty$ denote the ℓ_2 and ℓ_∞ norms. Given set \mathcal{R} , $\mathbf{a}_{\mathcal{R}}$ zeros out elements indexed by \mathcal{R}^c . For matrix $\mathbf{A} \in \mathbb{R}^{m \times q}$, A_{ij} is the (i, j) -th element. With sets \mathcal{R} and \mathcal{C} , $\mathbf{A}_{\mathcal{R}, \mathcal{C}}$ retains rows in \mathcal{R} and columns in \mathcal{C} , zeroing others. Special cases: $\mathbf{A}_{:, \mathcal{C}} = \mathbf{A}_{\mathcal{R}, \mathcal{C}}$ when $|\mathcal{R}| = m$, and $\mathbf{A}_{\mathcal{R}} = \mathbf{A}_{\mathcal{R}, \mathcal{C}}$ when $\mathcal{C} = \mathcal{R}$.

2 COLUMN THRESHOLDING

We present a novel column thresholding algorithm for the spiked Wigner model that achieves the information-theoretically optimal signal strength requirement. Our method exploits the key insight that column entries of the observation matrix provide stronger statistical separation than diagonal entries, enabling recovery with signal strength $\beta = \tilde{\Omega}(\sqrt{s})$ rather than the $\tilde{\Omega}(s)$ required by existing polynomial-time methods. After developing the algorithm and analyzing its computational complexity, we establish theoretical guarantees for both estimation accuracy and support recovery.

2.1 ALGORITHM

Diagonal thresholding (Johnstone & Lu, 2009) is a well-studied algorithm for the spiked Wigner model that offers low computational cost but requires signal strength $\beta = \tilde{\Omega}(s)$ for consistent estimation (Hopkins et al., 2017; Choo & d’Orsi, 2021). This requirement significantly exceeds the information-theoretic lower bound of $\beta = \tilde{\Omega}(\sqrt{s})$ (Banks et al., 2018; Perry et al., 2018; 2020). We propose a novel thresholding algorithm that closes this gap.

To understand the limitations of diagonal thresholding, we analyze its signal strength requirements. The algorithm estimates the support of \mathbf{u} by selecting indices corresponding to the s largest diagonal entries of \mathbf{Y} , then computes the leading eigenvector of the resulting submatrix. This approach exploits the expected diagonal structure:

$$\mathbb{E}[\mathbf{Y}]_{ii} = \begin{cases} \beta |u_i|^2, & i \in \mathcal{T}, \\ 0, & i \in \mathcal{T}^c, \end{cases} \quad (2)$$

where \mathcal{T} is the support of \mathbf{u} . The statistical gap between in-support and out-of-support entries is

$$g_{\text{diag}} := \min_{i \in \mathcal{T}} \mathbb{E}[\mathbf{Y}]_{ii} - \max_{i \in \mathcal{T}^c} \mathbb{E}[\mathbf{Y}]_{ii} = \beta \cdot \min_{i \in \mathcal{T}} |u_i|^2. \quad (3)$$

The following proposition shows when diagonal thresholding successfully identifies the support:

Proposition 2.1. *If $|W_{jj}| \leq \frac{1}{2} g_{\text{diag}}$ holds for all $j \in [d]$, then $Y_{ii} > Y_{i'i'}$ for all $i \in \mathcal{T}$ and $i' \in \mathcal{T}^c$.*

The proof of Proposition 2.1 is provided in Appendix A.3. Since $W_{jj} \sim \mathcal{N}(0, 2)$ independently, the condition holds with high probability when g_{diag} is sufficiently large. In that case, the diagonal entries Y_{ii} for $i \in \mathcal{T}$ exceed those for $i \in \mathcal{T}^c$, so diagonal thresholding recovers the support of \mathbf{u} by selecting the largest s diagonal entries of \mathbf{Y} . The success probability is governed by the gap g_{diag} : a larger gap g_{diag} yields a higher probability of correctly estimating the support. Alternatively, when g_{diag}/β is large, achieving any target gap sufficient for successful recovery requires less β .

Our key insight is to leverage column entries instead of diagonal entries to achieve better separation between in-support and out-of-support indices. Our approach is based on the observation that, when $l \in \mathcal{T}$, the expected column structure is

$$\mathbb{E}[\mathbf{Y}]_{il} = \begin{cases} \beta u_i u_l, & i \in \mathcal{T}, \\ 0, & i \in \mathcal{T}^c. \end{cases} \quad (4)$$

The resultant gap becomes:

$$g_{\text{col}} := \min_{i \in \mathcal{T}} |\mathbb{E}[\mathbf{Y}]_{il}| - \max_{i \in \mathcal{T}^c} |\mathbb{E}[\mathbf{Y}]_{il}| = \beta |u_l| \min_{i \in \mathcal{T}} |u_i|. \quad (5)$$

Crucially, $g_{\text{col}} = \beta |u_l| \min_{i \in \mathcal{T}} |u_i| \geq \beta \min_{i \in \mathcal{T}} |u_i|^2 = g_{\text{diag}}$ whenever $l \in \mathcal{T}$, providing enhanced separation that enables recovery with weaker signal strength requirement. To maximize the gap g_{col} , l should ideally be the index of the largest absolute element of \mathbf{u} , which aligns with the index of the largest diagonal entry of $\mathbb{E}[\mathbf{Y}]$, as shown in (2). However, since we only have the noisy matrix \mathbf{Y} , we choose l as the index of the largest diagonal entry of \mathbf{Y} , denoted by i_0 .

Algorithm 1 implements our column thresholding approach in two steps: (1) estimate the support $\hat{\mathcal{T}}$ using the s largest entries of the i_0 -th column, where $i_0 = \arg \max_i Y_{ii}$; (2) reconstruct the spike vector using the leading eigenvector of the $s \times s$ submatrix $\mathbf{Y}_{\hat{\mathcal{T}}}$, formed by restricting to rows and columns indexed by $\hat{\mathcal{T}}$. For computational efficiency, Algorithm 2 presents a variant that directly normalizes the selected column entries.

Algorithm 1 Column Thresholding

```

1: Input: Matrix  $\mathbf{Y} \in \mathbb{R}^{d \times d}$ , sparsity level  $s$ 
2: Output: Estimated sparse unit vector  $\hat{\mathbf{u}} \in \mathbb{R}^d$ 
3:  $i_0 \leftarrow \arg \max_{i \in [d]} Y_{ii}$  ▷ Find index of largest diagonal entry
4:  $\hat{\mathcal{T}} \leftarrow$  indices of  $s$  largest entries in  $|\mathbf{Y}_{:,i_0}|$  ▷ Estimate support
5:  $\mathbf{Y}_{\hat{\mathcal{T}}} \leftarrow$  submatrix of  $\mathbf{Y}$  with rows and columns in  $\hat{\mathcal{T}}$ 
6:  $\mathbf{v} \leftarrow$  leading eigenvector of  $\mathbf{Y}_{\hat{\mathcal{T}}}$  with  $\|\mathbf{v}\|_2 = 1$ 
7: Initialize  $\hat{\mathbf{u}} \leftarrow \mathbf{0} \in \mathbb{R}^d$  and set  $\hat{\mathbf{u}}_{\hat{\mathcal{T}}} \leftarrow \mathbf{v}$  ▷ Embed eigenvector into full space
8: return  $\hat{\mathbf{u}}$ 

```

Algorithm 2 Column Thresholding (Normalization Variant)

```

1: Input: Matrix  $\mathbf{Y} \in \mathbb{R}^{d \times d}$ , sparsity level  $s$ 
2: Output: Estimated sparse unit vector  $\hat{\mathbf{u}} \in \mathbb{R}^d$ 
3:  $i_0 \leftarrow \arg \max_{i \in [d]} Y_{ii}$  ▷ Find index of largest diagonal entry
4:  $\hat{\mathcal{T}} \leftarrow$  indices of  $s$  largest entries in  $|\mathbf{Y}_{:,i_0}|$  ▷ Estimate support
5: Set  $\hat{\mathbf{u}}_{\text{nv}} \leftarrow \mathbf{Y}_{\hat{\mathcal{T}},i_0} / \|\mathbf{Y}_{\hat{\mathcal{T}},i_0}\|_2$  ▷ Embed normalized subvector into full space
6: return  $\hat{\mathbf{u}}_{\text{nv}}$ 

```

The enhanced statistical gap in our column-based approach is the key to achieving the optimal signal strength requirement of $\beta = \tilde{\Omega}(\sqrt{s})$. As detailed in Section 3.2, this improvement stems from leveraging the correlations between entries in the selected column, which provides stronger signal concentration than the independent diagonal entries used in diagonal thresholding.

Algorithm 2 presents a computationally efficient variant that applies the same thresholding strategy for support estimation but directly normalizes the i_0 -th column rather than computing an eigenvalue decomposition. This variant trades modest estimation accuracy for reduced computational cost while preserving the optimal signal strength requirement of $\beta = \tilde{\Omega}(\sqrt{s})$ and maintaining the same theoretical guarantees for support recovery. The variant is practical when computational resources are limited or rapid support identification is prioritized over exact reconstruction.

2.2 COMPUTATIONAL COMPLEXITY

Algorithm 1 requires three main operations: finding the largest diagonal entry ($O(d)$), selecting the s largest column entries ($O(d \log d)$ via sorting or $O(d + s \log s)$ using partial sorting), and computing the leading eigenvector of an $s \times s$ matrix ($O(s^3)$). The total complexity is $O(d \log d + s^3)$, which reduces to $O(d \log d)$ when $s = O((d \log d)^{1/3})$ —a regime covering many practical sparse recovery scenarios. The normalization variant (Algorithm 2) eliminates the eigendecomposition step, achieving $O(d \log d)$ complexity uniformly.

For comparison, diagonal thresholding (Johnstone & Lu, 2009) follows a similar computational pattern: it finds the s largest diagonal entries ($O(d \log d)$) and computes the leading eigenvector of the resulting $s \times s$ submatrix ($O(s^3)$), yielding the same $O(d \log d + s^3)$ complexity. Spectral methods, the vanilla spectral algorithm and spectral projection (Brennan et al., 2018), compute the leading eigenvector of the full $d \times d$ matrix \mathbf{Y} , requiring $O(d^3)$ operations.

Column thresholding achieves fundamental improvement over existing approaches. While diagonal thresholding and spectral methods both require suboptimal signal strength $\beta = \tilde{\Omega}(s)$ for consistent recovery, our algorithm achieves the information-theoretically optimal requirement of $\beta = \tilde{\Omega}(\sqrt{s})$. Moreover, we maintain the computational efficiency of diagonal thresholding and offer substantial speedup over spectral methods. This unique combination of computational efficiency and statistical optimality makes our approach particularly valuable for modern high-dimensional applications.

2.3 THEORETICAL ANALYSIS

We establish theoretical guarantees showing that column thresholding achieves the information-theoretically optimal signal strength requirement of $\beta = \tilde{\Omega}(\sqrt{s})$ under mild conditions. We analyze estimation accuracy and support recovery separately.

2.3.1 ESTIMATION ERROR

We analyze estimation accuracy using the following distance metric accounting for sign ambiguity:

$$\text{dist}(\mathbf{u}, \hat{\mathbf{u}}) := \min \{ \|\mathbf{u} - \hat{\mathbf{u}}\|_2, \|\mathbf{u} + \hat{\mathbf{u}}\|_2 \}. \quad (6)$$

This metric is standard in PCA and phase retrieval problems where the sign of the recovered vector is inherently ambiguous.

Theorem 2.2. *Let $\mathbf{u} \in \mathbb{R}^d$ be an s -sparse unit vector and $\mathbf{Y} = \beta \mathbf{u} \mathbf{u}^\top + \mathbf{W}$, where $\beta > 0$ and $\mathbf{W} \in \mathbb{R}^{d \times d}$ is distributed as $\text{GOE}(d)$. For any target accuracy $\zeta \in (0, 1]$, if the signal strength satisfies*

$$\beta \geq C_1 \zeta^{-1} \|\mathbf{u}\|_\infty^{-1} \sqrt{s \log d}$$

for some universal constant $C_1 > 0$, then with probability at least $1 - 1.6d^{-1}$, the output $\hat{\mathbf{u}}$ of Algorithm 1 satisfies $\text{dist}(\mathbf{u}, \hat{\mathbf{u}}) \leq \zeta$.

This theorem, proved in Appendix A.4, shows that column thresholding achieves signal strength scaling of $\Omega(\sqrt{s \log d})$ for constant estimation error when $\|\mathbf{u}\|_\infty = \Omega(1)$, matching the information-theoretic optimum of $\tilde{\Omega}(\sqrt{s})$. This bridges the computational-statistical gap between the information-theoretic lower bound and the $\tilde{\Omega}(s)$ requirement of existing polynomial-time algorithms, including diagonal thresholding and spectral methods.

Attaining the optimal signal strength $\tilde{\Omega}(\sqrt{s})$ in polynomial time requires additional assumptions. Indeed, computational hardness results based on the planted clique conjecture indicate that no polynomial-time algorithm can achieve the optimal signal strength without additional structural assumptions (Brennan et al., 2018). This infinity norm condition is mild and naturally satisfied in many applications. For instance, when the nonzero entries of \mathbf{u} follow a power-law decay—a common model in compressive sensing (Donoho, 2006; Candès et al., 2006)—the infinity norm requirement is automatically satisfied. Similar phenomena arise in sparse phase retrieval, where power-law signals enable optimal recovery (Jagatap & Hegde, 2019).

Theorem 2.3. *Under the same model assumptions as Theorem 2.2, if the signal strength satisfies*

$$\beta \geq C_2 \zeta^{-2} \|\mathbf{u}\|_\infty^{-1} \sqrt{s \log d}$$

for some universal constant $C_2 > 0$, then with probability at least $1 - 1.4d^{-1}$, the output $\hat{\mathbf{u}}_{\text{nv}}$ of Algorithm 2 satisfies $\text{dist}(\mathbf{u}, \hat{\mathbf{u}}_{\text{nv}}) \leq \zeta$.

The proof of Theorem 2.3 is provided in Appendix A.5. While the normalization variant (Algorithm 2) requires a stronger dependence of β on the accuracy parameter ζ (quadratic rather than linear), it maintains the optimal scaling with respect to s and d . This highlights a key insight: the statistical efficiency is determined by the column thresholding step for support estimation, not by the reconstruction method (eigendecomposition versus direct normalization). The reconstruction only affects the constant factors in the accuracy guarantee, confirming that our column-based support estimation is the fundamental innovation enabling optimal signal strength requirements.

2.3.2 SUPPORT RECOVERY

Beyond estimation accuracy, exact support recovery is essential for interpretability and many downstream applications. We now establish conditions under which our algorithms correctly identify the support $\mathcal{T} := \{i : u_i \neq 0\}$ of the spike vector.

Theorem 2.4. *Let $\mathbf{u} \in \mathbb{R}^d$ be an s -sparse unit vector satisfying $|u_i| \geq \theta/\sqrt{s}$ for all $i \in \mathcal{T}$ and some constant $\theta > 0$. Under the spiked Wigner model with signal strength*

$$\beta \geq C_3 \theta^{-1} \|\mathbf{u}\|_\infty^{-1} \sqrt{s \log d}$$

for some universal constant $C_3 > 0$, both Algorithms 1 and 2 recover the support exactly (i.e., $\hat{\mathcal{T}} = \mathcal{T}$) with probability at least $1 - 1.3d^{-1}$.

The proof of Theorem 2.4 is shown in Appendix A.6. Our support recovery guarantee attains the signal strength scaling of $\Omega(\sqrt{s \log d})$, matching the information-theoretic limits $\tilde{\Omega}(\sqrt{s})$. The minimum magnitude assumption $|u_i| \geq \theta/\sqrt{s}$ ensures that all nonzero entries are sufficiently strong to be distinguished from noise, as detailed in the proof in Appendix A.6. This condition is standard in the sparse recovery literature and naturally holds for many structured signals. Since Algorithms 1 and 2 differ only in their estimation procedures while using identical support recovery methods, they achieve the same support recovery guarantees.

3 TRUNCATED POWER METHOD

Column thresholding meets the optimal signal-strength requirement but can benefit from iterative refinement to improve accuracy. Its one-shot estimate, though supported by strong theory, may not attain the smallest achievable error. In this section, we show how the truncated power method iteratively refines the initial estimate, yielding improved estimation accuracy.

In the spiked Wigner model, recovering the sparse spike \mathbf{u} from the noisy observation \mathbf{Y} in (1) naturally leads to the sparse PCA formulation:

$$\underset{\mathbf{w}}{\text{maximize}} \mathbf{w}^\top \mathbf{Y} \mathbf{w}, \quad \text{subject to } \|\mathbf{w}\|_2 = 1, \|\mathbf{w}\|_0 \leq k. \quad (7)$$

where k is a sparsity parameter. Since \mathbf{u} is the leading eigenvector of $\mathbb{E}[\mathbf{Y}] = \beta \mathbf{u} \mathbf{u}^\top$, the solution to (7) provides a natural estimator for \mathbf{u} .

The truncated power method (Yuan & Zhang, 2013) is an iterative algorithm designed to solve the sparse PCA problem (7). Starting from an initial vector \mathbf{u}^0 , it alternates between power iteration and hard thresholding:

$$\mathbf{u}^t = \mathcal{P}_{\mathbb{S}^{d-1}}(\mathcal{H}_k(\mathbf{Y} \mathbf{u}^{t-1})), \quad (8)$$

where $\mathcal{P}_{\mathbb{S}^{d-1}} : \mathbb{R}^d \setminus \{\mathbf{0}\} \rightarrow \mathbb{S}^{d-1}$ defined by $\mathcal{P}_{\mathbb{S}^{d-1}}(\mathbf{z}) = \mathbf{z}/\|\mathbf{z}\|_2$, and $\mathcal{H}_k : \mathbb{R}^d \rightarrow \mathbb{R}^d$ is the hard thresholding operator that retains the k largest entries (in absolute value) and zeros out the rest. The parameter k denotes the sparsity level used in (7), whereas s denotes the sparsity of the true spike. We assume k is of the same order as the true sparsity s throughout, and set $k = s$ in all experiments.

For computational efficiency, we exploit the sparsity structure. Let $\mathcal{T}^t = \text{supp}(\mathbf{u}^t)$ denote the support of iterate t . Since $|\mathcal{T}^t| \leq k$, we can rewrite the matrix-vector multiplication as:

$$\mathbf{u}^t = \mathcal{P}_{\mathbb{S}^{d-1}}(\mathcal{H}_k(\mathbf{Y}_{:, \mathcal{T}^{t-1}} \mathbf{u}_{\mathcal{T}^{t-1}}^{t-1})), \quad (9)$$

where $\mathbf{Y}_{:, \mathcal{T}}$ denotes the submatrix of \mathbf{Y} with columns indexed by \mathcal{T} . This reduces the per-iteration complexity from $O(d^2)$ to $O(ds)$.

The sparse PCA problem (7) is highly nonconvex due to the cardinality constraint, resulting in a landscape riddled with local maxima. Like all iterative methods for such problems, the truncated power method’s performance hinges on initialization quality—poor starting points can trap the algorithm in suboptimal local maxima or prevent convergence entirely. The choice of initialization thus becomes crucial for achieving good performance.

Yuan & Zhang (2013) provided a sharp characterization of when the truncated power method succeeds: geometric convergence to a near-optimal solution is guaranteed when the initial vector \mathbf{u}^0 has sufficient correlation with the truth:

$$|\langle \mathbf{u}^0, \mathbf{u} \rangle| \geq c \quad (10)$$

for some constant $c > 0$. However, obtaining such initialization is the key challenge—random initialization typically fails in high dimensions, while existing polynomial-time algorithms require at least $\beta = \tilde{\Omega}(s)$. The column thresholding algorithm provides a simple and effective approach that operates under weaker signal strength conditions in this setting. As shown in Section 2.3, it produces an initialization with two key properties:

- Near-perfect correlation: $|\langle \hat{\mathbf{u}}, \mathbf{u} \rangle| \geq 1 - \frac{\zeta^2}{2}$ for arbitrarily small $\zeta > 0$.
- Optimal signal strength: It succeeds under $\beta = \Omega(\sqrt{s \log d})$, matching the information-theoretic limit of $\tilde{\Omega}(\sqrt{s})$.

With this initialization, the truncated power method attains near-optimal estimation accuracy under minimal signal strength requirements—a combination unattainable by either method alone. We formalize the resulting convergence guarantee in Section 3.2, and summarize the full two-stage procedure in Algorithm 3.

3.1 COMPUTATIONAL COMPLEXITY

The column thresholding initialization requires $O(d \log d + s^3)$ operations, as detailed in Section 2.2. Each truncated power iteration involves two main steps: (i) a sparse matrix-vector multiplication costing $O(ds)$ operations, and (ii) sorting the resulting vector requiring $O(d \log d)$ operations. Thus, each iteration costs $O(ds + d \log d)$ operations. Since the method converges in $O(\log(1/\epsilon))$ iterations to achieve ϵ -accuracy (Section 3.2), the total refinement cost is $O((ds + d \log d) \log(1/\epsilon))$.

Algorithm 3 Truncated Power Method (TPM)

```

1: Input: Matrix  $\mathbf{Y} \in \mathbb{R}^{d \times d}$ , sparsity  $s$ , parameter  $k$ 
2: Output: Estimated sparse unit vector  $\mathbf{u}^t \in \mathbb{R}^d$ 
3:  $\mathbf{u}^0 \leftarrow \text{Column Thresholding}(\mathbf{Y}, s)$  ▷ Column thresholding initialization
4: for  $t = 1, 2, \dots$  do
5:    $\mathcal{T}^{t-1} \leftarrow \text{supp}(\mathbf{u}^{t-1})$ 
6:    $\mathbf{z}^t \leftarrow \mathbf{Y}_{:, \mathcal{T}^{t-1}} \mathbf{u}_{\mathcal{T}^{t-1}}$  ▷ Sparse matrix-vector product
7:    $\mathbf{u}^t \leftarrow \mathcal{P}_{\mathbb{S}^{d-1}}(\mathcal{H}_k(\mathbf{z}^t))$ 
8: end for
9: return  $\mathbf{u}^t$ 

```

3.2 THEORETICAL RESULTS

This section establishes the theoretical convergence guarantee for truncated power method, showing that under the information-theoretically optimal signal scaling the algorithm achieves geometric contraction of the optimization error down to an irreducible statistical floor.

Theorem 3.1. *Let $\mathbf{u} \in \mathbb{R}^d$ be an s -sparse unit vector and $\mathbf{Y} = \beta \mathbf{u} \mathbf{u}^\top + \mathbf{W}$, where $\beta > 0$ and $\mathbf{W} \sim \text{GOE}(d)$. Fix any $\zeta \in (0, 1)$. There exist universal constants $C_4, C_5 > 0$ such that, if*

$$\beta \geq C_4 \max \{ \|\mathbf{u}\|_\infty^{-1}, \zeta^{-1} \} \sqrt{s \log d},$$

then, with probability at least $1 - 1.5d^{-1}$, the sequence $\{\mathbf{u}^t\}_{t \geq 1}$ produced by Algorithm 3, initialized with \mathbf{u}^0 from Algorithm 1 and using parameter $k = C_5 s$, satisfies

$$\text{dist}(\mathbf{u}, \mathbf{u}^t) \leq \eta^t \text{dist}(\mathbf{u}, \mathbf{u}^0) + h\zeta, \quad (11)$$

where $\eta, h \in (0, 1)$ are universal constants.

The proof is provided in Appendix A.7. Theorem 3.1 decomposes the estimation error into two parts. The first term, $\eta^t \text{dist}(\mathbf{u}, \mathbf{u}^0)$, represents an optimization error that decays geometrically with iteration count t . The second term, $h\zeta$, captures the irreducible statistical error inherent to the problem. Consequently, the truncated power method rapidly eliminates the optimization error—requiring only $O(\log \zeta^{-1})$ iterations to achieve ζ -accuracy—while operating under the information-theoretically optimal signal strength $\beta = \tilde{\Omega}(\sqrt{s})$ when $\|\mathbf{u}\|_\infty = \Omega(1)$.

4 EXPERIMENTAL RESULTS

We empirically verify that column thresholding satisfies the information-theoretic signal-strength requirement, thereby validating our theoretical guarantees. We further show that TPM outperforms existing methods in estimation accuracy, support recovery, and computational efficiency. Performance is assessed by estimation error (via the distance metric in (6)) and F-score (0–1 scale, with 1 indicating perfect recovery). All results are averaged over 200 independent runs.

4.1 EMPIRICAL VALIDATION OF OPTIMAL SIGNAL THRESHOLDS

We empirically verify that our column thresholding achieves the information-theoretic signal strength requirement for constant estimation error and exact support recovery, thereby validating Theorem 2.2 and Theorem 2.4. In each trial, we construct an s -sparse spike \mathbf{u} with one entry of magnitude 0.5 and the remaining $s - 1$ nonzeros of equal magnitude, normalized so that $\|\mathbf{u}\|_2 = 1$. This design ensures two key properties: (i) $\|\mathbf{u}\|_\infty$ is constant in d and s (as required by Theorem 2.2), and (ii) every nonzero entry satisfies $|u_i| \geq \theta/\sqrt{s}$ for a universal constant θ (as required by Theorem 2.4). We then generate \mathbf{Y} according to the spiked Wigner model in (1).

We conduct experiments across two complementary regimes: (i) varying the dimension $d \in \{2000, 5000, 10000, 15000, 20000\}$ with fixed sparsity $s = 20$, and (ii) varying the sparsity $s \in \{20, 50, 100, 150, 200\}$ with fixed dimension $d = 10000$. When performance is plotted against the scaled signal strength $\beta/\sqrt{s \log d}$, the curves from different (d, s) collapse, indicating that the phase transition depends only on this scaled quantity, confirming the theoretical scaling.

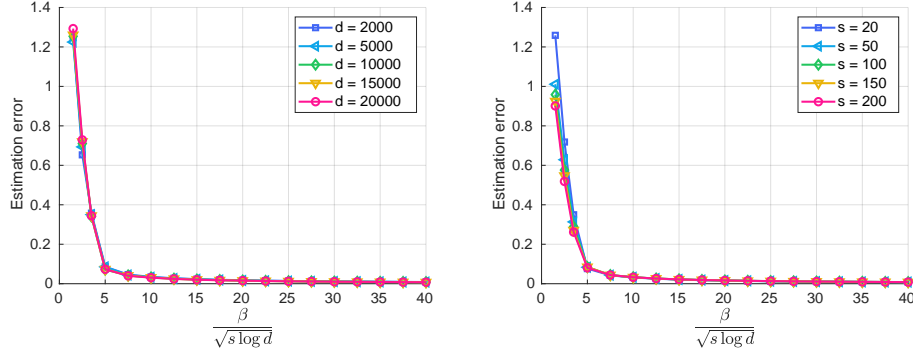


Figure 2: Estimation error versus scaled signal strength for our column thresholding under varying dimensions (left) and sparsities (right).

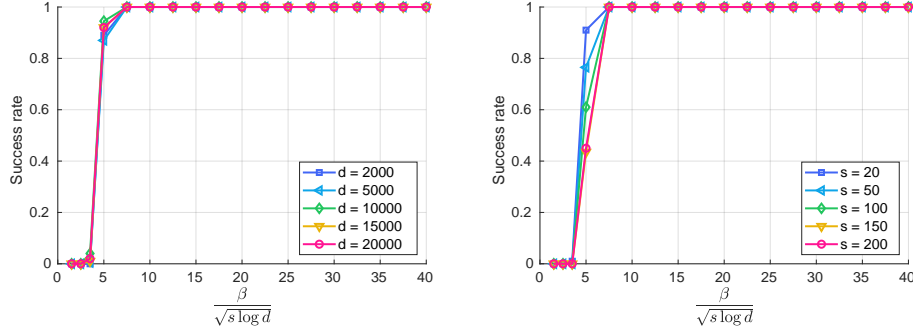


Figure 3: Success rate versus scaled signal strength for our column thresholding under varying dimensions (left) and sparsities (right).

Figure 2 shows that curves from all (d, s) settings collapse once $\beta/\sqrt{s \log d} \geq 10$, regardless of the specific values of d and s . This indicates that column thresholding succeeds when $\beta \geq C_1 \sqrt{s \log d}$ with a universal constant $C_1 \approx 10$, thereby validating the $\Omega(\sqrt{s \log d})$ signal strength requirement established in Theorem 2.2.

Figure 3 exhibits a sharp phase transition for support recovery at $\beta/\sqrt{s \log d} \approx 7.5$. Below this threshold, perfect recovery is not guaranteed; above it, the success rate reaches 1 uniformly across all tested d and s . This empirically validates the $\Omega(\sqrt{s \log d})$ signal strength requirement for successful support recovery established in Theorem 2.4.

4.2 COMPARATIVE EVALUATION: STATISTICAL AND COMPUTATIONAL PERFORMANCE

We evaluate TPM against three established approaches: diagonal thresholding (DT) (Johnstone & Lu, 2009), covariance thresholding (CT) (Krauthgamer et al., 2015), and spectral projection (SP) (Brennan et al., 2018). For all experiments, we construct the true spike \mathbf{u} with s randomly-located nonzero entries, each taking values $\pm 1/\sqrt{s}$ with equal probability. This balanced spike design, standard in the sparse PCA literature (Krauthgamer et al., 2015), ensures $\|\mathbf{u}\|_2 = 1$ while maintaining uniform entry magnitudes.

Figure 4 reports performance as the signal strength β varies, with fixed dimension $d = 2000$ and sparsity $s = 10$. Our TPM shows clear advantages, especially in weak-signal regimes. As the signal weakens, CT and SP degrade markedly, while TPM maintains strong accuracy and support recovery. TPM consistently outperforms DT across all signal strengths, validating our analysis in Section 2.1 that column thresholding yields a fundamentally larger statistical gap than diagonal thresholding.

Figure 5 evaluates computational scalability by varying the dimension d from 2000 to 10000, with fixed signal strength $\beta = 100$ and sparsity $s = 15$. Our TPM attains the lowest estimation error across all dimensions, whereas competing methods deteriorate as d grows. Notably, this accuracy comes with minimal computational overhead: TPM’s runtime scales comparably to the efficient DT method. By contrast, CT and SP incur substantially higher costs. Overall, these results show that TPM combines strong statistical performance with practical computational efficiency.

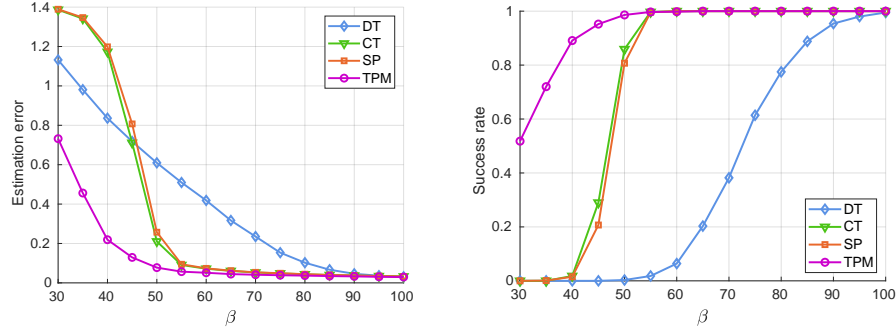
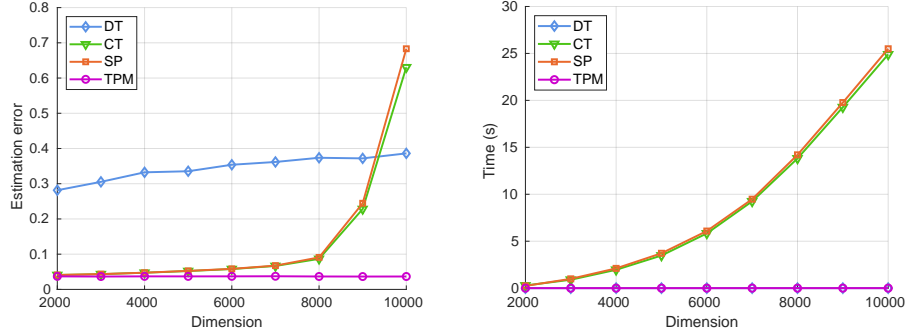
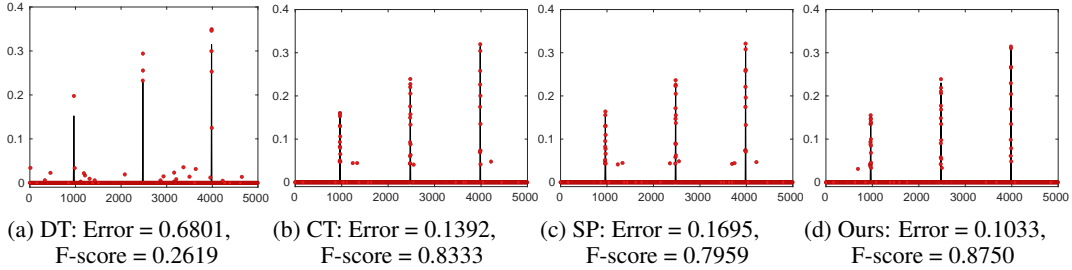
Figure 4: Estimation error (left) and support-recovery success rate (right) versus signal strength β .

Figure 5: Scalability analysis across dimensions: estimation error (left) and runtime (right).

Figure 6: Three-peak benchmark results. True spike (black curve) versus estimated spike (red markers) for four methods: (a) DT, (b) CT, (c) SP, and (d) TPM. The true signal comprises three Beta densities on $[0, 1]$ with dimension $p = 5000$ and signal strength $\beta = 100$.

4.3 THREE-PEAK BENCHMARK EVALUATION

We evaluate our method on the canonical “three-peak” experiment (Johnstone & Lu, 2009), a demanding benchmark for sparse recovery in high dimensions. The true spike \mathbf{v} is constructed as a mixture of three Beta densities on $[0, 1]$, producing three pronounced peaks separated by near-zero valleys. This setup rigorously tests an algorithm’s ability to localize multiple signal components while suppressing inter-peak noise. The experiment stresses methods’ capacity to distinguish true signal peaks from spurious activations. As shown in Figure 6, TPM faithfully recovers all three peaks, achieving lower estimation error and higher F-score than competing methods, which either misestimate the peaks or introduce false detections in the valleys.

5 CONCLUSIONS

We introduce two complementary algorithms for sparse PCA in the spiked Wigner model. Column thresholding achieves a breakthrough in computational-statistical tradeoffs: it runs in polynomial time while requiring only $\tilde{\Omega}(\sqrt{s})$ signal strength when $\|\mathbf{u}\|_\infty = \Omega(1)$ holds. Truncated power method iteratively refines the column thresholding estimate with provable linear convergence. Extensive experiments validate our theoretical guarantees and demonstrate superior performance over existing methods in estimation accuracy, support recovery, and computational efficiency.

REFERENCES

- 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.
- 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 Transactions on Information Theory*, 64(7):4872–4894, 2018.
- Richard Baraniuk, Mark Davenport, Ronald DeVore, and Michael Wakin. A simple proof of the restricted isometry property for random matrices. *Constructive Approximation*, 28:253–263, 2008.
- Jean Barbier, Mohamad Dia, Nicolas Macris, Florent Krzakala, Thibault Lesieur, and Lenka Zdeborová. Mutual information for symmetric rank-one matrix estimation: A proof of the replica formula. In *Advances in Neural Information Processing Systems*, volume 29, 2016.
- Florent Benaych-Georges and Raj Rao Nadakuditi. The eigenvalues and eigenvectors of finite, low rank perturbations of large random matrices. *Advances in Mathematics*, 227(1):494–521, 2011.
- Matthew Brennan, Guy Bresler, and Wasim Huleihel. Reducibility and computational lower bounds for problems with planted sparse structure. In *Conference On Learning Theory*, pp. 48–166, 2018.
- Emmanuel J Candès, Justin Romberg, and Terence Tao. Robust uncertainty principles: Exact signal reconstruction from highly incomplete frequency information. *IEEE Transactions on Information Theory*, 52(2):489–509, 2006.
- Emmanuel J Candes, Xiaodong Li, and Mahdi Soltanolkotabi. Phase retrieval via wirtinger flow: Theory and algorithms. *IEEE Transactions on Information Theory*, 61(4):1985–2007, 2015.
- Mireille Capitaine, Catherine Donati-Martin, and Delphine Féral. The largest eigenvalues of finite rank deformation of large Wigner matrices: Convergence and nonuniversality of the fluctuations. *The Annals of Probability*, 37(1):1–47, 2009.
- Yuxin Chen, Yuejie Chi, and Andrea J. Goldsmith. Exact and stable covariance estimation from quadratic sampling via convex programming. *IEEE Transactions on Information Theory*, 61(7):4034–4059, 2015.
- Yuejie Chi, Yue M Lu, and Yuxin Chen. Nonconvex optimization meets low-rank matrix factorization: An overview. *IEEE Transactions on Signal Processing*, 67(20):5239–5269, 2019.
- Davin Choo and Tommaso d’Orsi. The complexity of sparse tensor PCA. *Advances in Neural Information Processing Systems*, 34:7993–8005, 2021.
- Yash Deshpande and Andrea Montanari. Information-theoretically optimal sparse PCA. In *2014 IEEE International Symposium on Information Theory*, pp. 2197–2201, 2014.
- Yash Deshpande and Andrea Montanari. Sparse PCA via covariance thresholding. *Journal of Machine Learning Research*, 17(141):1–41, 2016.
- Yash Deshpande, Emmanuel Abbe, and Andrea Montanari. Asymptotic mutual information for the binary stochastic block model. In *2016 IEEE International Symposium on Information Theory (ISIT)*, pp. 185–189, 2016.
- David L Donoho. Compressed sensing. *IEEE Transactions on information theory*, 52(4):1289–1306, 2006.
- Delphine Féral and Sandrine Péché. The largest eigenvalue of rank one deformation of large Wigner matrices. *Communications in Mathematical Physics*, 272(1):185–228, 2007.
- Samuel B Hopkins, Pravesh K Kothari, Aaron Potechin, Prasad Raghavendra, Tselil Schramm, and David Steurer. The power of sum-of-squares for detecting hidden structures. In *2017 IEEE 58th Annual Symposium on Foundations of Computer Science (FOCS)*, pp. 720–731, 2017.
- Roger A Horn and Charles R Johnson. *Matrix analysis*. Cambridge University Press, 2012.

- Gauri Jagatap and Chinmay Hegde. Sample-efficient algorithms for recovering structured signals from magnitude-only measurements. *IEEE Transactions on Information Theory*, 65(7):4434–4456, 2019.
- Adel Javanmard, Andrea Montanari, and Federico Ricci-Tersenghi. Phase transitions in semidefinite relaxations. *Proceedings of the National Academy of Sciences*, 113(16):E2218–E2223, 2016.
- Iain M Johnstone and Arthur Yu Lu. On consistency and sparsity for principal components analysis in high dimensions. *Journal of the American Statistical Association*, 104(486):682–693, 2009.
- Robert Krauthgamer, Boaz Nadler, and Dan Vilenchik. Do semidefinite relaxations solve sparse PCA up to the information limit? *The Annals of Statistics*, 43(3):1300–1322, 2015.
- Marc Lelarge and Léo Miolane. Fundamental limits of symmetric low-rank matrix estimation. In *Conference on Learning Theory*, pp. 1297–1301. PMLR, 2017.
- Thibault Lesieur, Florent Krzakala, and Lenka Zdeborová. Phase transitions in sparse PCA. In *2015 IEEE International Symposium on Information Theory (ISIT)*, pp. 1635–1639, 2015.
- Matthias Löffler, Alexander S Wein, and Afonso S Bandeira. Computationally efficient sparse clustering. *Information and Inference: A Journal of the IMA*, 11(4):1255–1286, 2022.
- Cong Ma, Kaizheng Wang, Yuejie Chi, and Yuxin Chen. Implicit regularization in nonconvex statistical estimation: Gradient descent converges linearly for phase retrieval and matrix completion. In *International Conference on Machine Learning*, pp. 3345–3354, 2018.
- Andrea Montanari, Daniel Reichman, and Ofer Zeitouni. On the limitation of spectral methods: From the Gaussian hidden clique problem to rank-one perturbations of Gaussian tensors. In *Advances in Neural Information Processing Systems*, volume 28, 2015.
- Debashis Paul. Asymptotics of sample eigenstructure for a large dimensional spiked covariance model. *Statistica Sinica*, 17(4):1617–1642, 2007.
- S. Péché. The largest eigenvalue of small rank perturbations of Hermitian random matrices. *Probability Theory and Related Fields*, 134(1):127–173, 2006.
- Amelia Perry, Alexander S. Wein, Afonso S. Bandeira, and Ankur Moitra. Optimality and sub-optimality of PCA I: Spiked random matrix models. *The Annals of Statistics*, 46(5):2416–2451, 2018.
- Amelia Perry, Alexander S. Wein, and Afonso S. Bandeira. Statistical limits of spiked tensor models. *Annales de l’Institut Henri Poincaré, Probabilités et Statistiques*, 56(1):230–264, 2020.
- Luca Pesce, Bruno Loureiro, Florent Krzakala, and Lenka Zdeborová. Subspace clustering in high-dimensions: Phase transitions & statistical-to-computational gap. *Advances in Neural Information Processing Systems*, 35:27087–27099, 2022.
- Stephen Tu, Ross Boczar, Max Simchowitz, Mahdi Soltanolkotabi, and Ben Recht. Low-rank solutions of linear matrix equations via procrustes flow. In *International conference on machine learning*, pp. 964–973, 2016.
- Roman Vershynin. *High-dimensional probability: An introduction with applications in data science*, volume 47. Cambridge University Press, 2018.
- Xiao-Tong Yuan and Tong Zhang. Truncated power method for sparse eigenvalue problems. *Journal of Machine Learning Research*, 14(4), 2013.

Appendix A contains the proofs of our theoretical results, while Appendix B provides additional experimental results. We clarify that our usage of large language models (LLMs) is limited strictly to polish writing.

A PROOFS

In Appendix A, we prove the proposition and theorems introduced in Sections 2 and 3. First, we present some auxiliary lemmas in Appendix A.1 and show some technical lemmas in Appendix A.2. Next, we prove Proposition 2.1 in Appendix A.3. Subsequently, we present the proofs for Theorem 2.2, Theorem 2.3 and Theorem 2.4 in Appendix A.4, Appendix A.5 and Appendix A.6, respectively. Finally, we show the proof for Theorem 3.1 in Appendix A.7.

Throughout Appendix A, we define the largest and smallest ℓ -sparse eigenvalue of a symmetric matrix $\mathbf{B} \in \mathbb{R}^{m \times m}$ by

$$\lambda_{\max}(\mathbf{B}, \ell) = \max_{\mathbf{w} \in \mathbb{R}^m, \|\mathbf{w}\|_2=1, \|\mathbf{w}\|_0=\ell} \mathbf{w}^T \mathbf{B} \mathbf{w}, \quad \lambda_{\min}(\mathbf{B}, \ell) = \min_{\mathbf{w} \in \mathbb{R}^m, \|\mathbf{w}\|_2=1, \|\mathbf{w}\|_0=\ell} \mathbf{w}^T \mathbf{B} \mathbf{w},$$

respectively. Then we define the maximum spectral norm of all $\ell \times \ell$ submatrices of \mathbf{B} by

$$\rho(\mathbf{B}, \ell) = \max \{|\lambda_{\max}(\mathbf{B}, \ell)|, |\lambda_{\min}(\mathbf{B}, \ell)|\}. \quad (12)$$

A.1 AUXILIARY LEMMAS

The following two lemmas are used to prove the convergence of truncated power method in Algorithm 3, which will be used for the proof of Theorem 3.1 in Appendix A.7.

Lemma A.1 ((Yuan & Zhang, 2013)). *Let \mathbf{z} be the eigenvector with the largest eigenvalue (in absolute value) of a symmetric matrix \mathbf{B} , and let $\kappa < 1$ be the ratio of the second to the largest eigenvalue in absolute values. Given any \mathbf{y} such that $\|\mathbf{y}\|_2 = 1$, let $\mathbf{y}' = \mathbf{B}\mathbf{y}/\|\mathbf{B}\mathbf{y}\|_2$, then*

$$|\mathbf{z}^\top \mathbf{y}'| \geq |\mathbf{z}^\top \mathbf{y}| \left(1 + \frac{1}{2} (1 - \kappa^2) (1 - |\mathbf{z}^\top \mathbf{y}|^2) \right).$$

Lemma A.2 ((Yuan & Zhang, 2013)). *Consider \mathbf{y} with $\|\mathbf{y}\|_0 = \ell$. Consider \mathbf{z} and let $\mathcal{F} = \text{supp}(\mathbf{z}, \ell')$ be the set of indices with the ℓ' largest absolute values in \mathbf{z} . If $\|\mathbf{y}\|_2 = \|\mathbf{z}\|_2 = 1$, then*

$$|\mathbf{y}^\top \mathbf{z}_{\mathcal{F}}| \geq |\mathbf{y}^\top \mathbf{z}| - \sqrt{\ell/\ell'} \min \left\{ \sqrt{1 - |\mathbf{y}^\top \mathbf{z}|^2}, (1 + \sqrt{\ell/\ell'}) (1 - |\mathbf{y}^\top \mathbf{z}|^2) \right\}.$$

A.2 TECHNICAL LEMMAS

In Appendix A.2, we show some technical lemmas that will be used for the proofs of Theorem 2.2 and Theorem 3.1. The first lemma bounds the quantity $\rho(\mathbf{W}, \ell)$ defined in (12).

Lemma A.3. *For any $r \in (0, 1)$,*

$$\mathbb{P} \{ \rho(\mathbf{W}, \ell) \leq 3r\beta \} \geq 1 - \frac{2}{\sqrt{\pi}r\beta} \left(\frac{9ed}{\ell} \right)^\ell \exp \left(- \frac{r^2 \beta^2}{4} \right). \quad (13)$$

Proof. Denote the set of ℓ -sparse vectors in \mathbb{R}^d by $\mathbb{T}_\ell^d := \{ \mathbf{w} : \|\mathbf{w}\|_2 = 1, \|\mathbf{w}\|_0 = \ell \}$. For any $\delta \in (0, 1)$, there exists a set $\mathcal{N}_\delta \subset \mathbb{T}_\ell^d$ such that for any $\mathbf{w} \in \mathbb{T}_\ell^d$, there exists $\mathbf{w}_\delta \in \mathcal{N}_\delta$ such that $\text{supp}(\mathbf{w}) = \text{supp}(\mathbf{w}_\delta)$ and $\|\mathbf{w} - \mathbf{w}_\delta\|_2 \leq \delta$ and $|\mathcal{N}_\delta| \leq \binom{d}{\ell} \left(\frac{3}{\delta} \right)^\ell \leq \left(\frac{3ed}{\delta\ell} \right)^\ell$ (Baraniuk et al., 2008).

From (12), we obtain

$$\rho(\mathbf{W}, \ell) = \max_{\substack{\mathbf{y}, \mathbf{z} \in \mathbb{T}_\ell^d, \\ \text{supp}(\mathbf{y}) = \text{supp}(\mathbf{z})}} \mathbf{y}^\top \mathbf{W} \mathbf{z} =: \mathbf{y}_*^\top \mathbf{W} \mathbf{z}_*.$$

From the definition of \mathcal{N}_δ , there exists $\mathbf{y}_\delta, \mathbf{z}_\delta \in \mathcal{N}_\delta$ such that $\text{supp}(\mathbf{y}_\delta) = \text{supp}(\mathbf{y}_*) = \text{supp}(\mathbf{z}_*) = \text{supp}(\mathbf{z}_\delta)$, $\|\mathbf{y}_* - \mathbf{y}_\delta\|_2 \leq \delta$ and $\|\mathbf{z}_* - \mathbf{z}_\delta\|_2 \leq \delta$. Then we have

$$\mathbf{y}_*^\top \mathbf{W} \mathbf{z}_* = \mathbf{y}_*^\top \mathbf{W} (\mathbf{z}_* - \mathbf{z}_\delta) + (\mathbf{y}_* - \mathbf{y}_\delta)^\top \mathbf{W} \mathbf{z}_\delta + \mathbf{y}_\delta^\top \mathbf{W} \mathbf{z}_\delta \leq 2\delta \mathbf{y}_*^\top \mathbf{W} \mathbf{z}_* + \mathbf{y}_\delta^\top \mathbf{W} \mathbf{z}_\delta,$$

which implies that

$$\rho(\mathbf{W}, \ell) \leq (1 - 2\delta)^{-1} \mathbf{y}_\delta^\top \mathbf{W} \mathbf{z}_\delta \leq (1 - 2\delta)^{-1} \max_{\substack{\mathbf{y}, \mathbf{z} \in \mathcal{N}_\delta, \\ \text{supp}(\mathbf{y}) = \text{supp}(\mathbf{z})}} \mathbf{y}^\top \mathbf{W} \mathbf{z}, \quad (14)$$

where the inequalities hold when $1 - 2\delta > 0$.

Now for any $(\mathbf{y}, \mathbf{z}) \in \mathcal{N}_\delta$, we bound $|\mathbf{y}^\top \mathbf{W} \mathbf{z}|$ as follows. Since $\mathbf{W} = \frac{1}{\sqrt{2}}(\mathbf{A} + \mathbf{A}^\top)$ with some random matrix $\mathbf{A} \sim \mathcal{N}(0, 1)^{\otimes d \times d}$, we obtain

$$\mathbf{y}^\top \mathbf{W} \mathbf{z} \sim \mathcal{N}(0, 1 + |\mathbf{y}^\top \mathbf{z}|^2).$$

Therefore, using the tail of a Gaussian variable (Vershynin, 2018), for any $r \in (0, 1)$, it holds that

$$\mathbb{P}\{|\mathbf{y}^\top \mathbf{W} \mathbf{z}| \geq r\beta\} \leq \sqrt{\frac{2}{\pi}} \frac{\sqrt{1 + |\mathbf{y}^\top \mathbf{z}|^2}}{r\beta} \exp\left(-\frac{r^2\beta^2}{2(1 + |\mathbf{y}^\top \mathbf{z}|^2)}\right) \leq \frac{2}{\sqrt{\pi}r\beta} \exp\left(-\frac{r^2\beta^2}{4}\right),$$

where the last inequality we use $\|\mathbf{y}\|_2 = \|\mathbf{z}\|_2 = 1$.

By taking union bounds for all $(\mathbf{y}, \mathbf{z}) \in \mathcal{N}_\delta$, we obtain that

$$\mathbb{P}\left\{\max_{\substack{\mathbf{y}, \mathbf{z} \in \mathcal{N}_\delta, \\ \text{supp}(\mathbf{y}) = \text{supp}(\mathbf{z})}} |\mathbf{y}^\top \mathbf{W} \mathbf{z}| \leq r\beta\right\} \geq 1 - \frac{2}{\sqrt{\pi}r\beta} \left(\frac{3ed}{\delta\ell}\right)^\ell \exp\left(-\frac{r^2\beta^2}{4}\right).$$

Setting $\delta = \frac{1}{3}$ together with (14) leads to (13). \square

The next lemma bounds the error between \mathbf{u} and the ℓ -sparse largest eigenvector of \mathbf{Y} .

Lemma A.4. *Let $\Lambda \subset [d]$ be such that $\Lambda \cup \mathcal{T} \neq \emptyset$ and $|\Lambda| = \ell$. Let \mathbf{w} be the largest eigenvector of \mathbf{Y}_Λ with $\|\mathbf{w}\|_2 = 1$. If $\rho(\mathbf{W}, \ell) < \frac{\beta}{2}\|\mathbf{u}_\Lambda\|_2^2$, then we have*

$$\text{dist}(\mathbf{w}, \mathbf{u}_\Lambda)^2 \leq \|\mathbf{u}_\Lambda\|_2^2 + 1 - 2 \frac{\|\mathbf{u}_\Lambda\|_2}{\sqrt{1 + \frac{\rho(\mathbf{W}, \ell)^2}{\left(\beta\|\mathbf{u}_\Lambda\|_2^2 - 2\rho(\mathbf{W}, \ell)\right)^2}}}.$$

Proof. Denote $\bar{\lambda}$ the largest eigenvalue of \mathbf{Y}_Λ , i.e. $\bar{\lambda} = \lambda_1(\mathbf{Y}_\Lambda)$. Recall that $\mathbf{Y} = \beta\mathbf{u}\mathbf{u}^\top + \mathbf{W}$ and $\mathbb{E}[\mathbf{Y}] = \beta\mathbf{u}\mathbf{u}^\top$. Using Weyl's inequality (Horn & Johnson, 2012), it holds that

$$\bar{\lambda} \geq \lambda_1(\mathbb{E}[\mathbf{Y}_\Lambda]) + \lambda_n(\mathbf{W}_\Lambda) \geq \beta\|\mathbf{u}_\Lambda\|_2^2 - \rho(\mathbf{W}, \ell), \quad (15)$$

where the last inequality holds by $\lambda_n(\mathbf{W}_\Lambda) \geq -\rho(\mathbf{W}, \ell)$ from the definition. Similarly, we have for all $i \geq 2$,

$$\begin{aligned} |\lambda_i(\mathbf{Y}_\Lambda)| &\leq |\lambda_i(\mathbb{E}[\mathbf{Y}_\Lambda])| + |\lambda_i(\mathbf{Y}_\Lambda) - \lambda_i(\mathbb{E}[\mathbf{Y}_\Lambda])| \\ &= \max\{|\lambda_1(\mathbf{W}_\Lambda)|, |\lambda_n(\mathbf{W}_\Lambda)|\} \leq \rho(\mathbf{W}, \ell). \end{aligned} \quad (16)$$

Notice $\|\mathbf{w}\|_2 = 1$ but $\|\mathbf{u}_\Lambda\|_2 \leq 1$. We divide \mathbf{w} as

$$\mathbf{w} = a_1 \frac{\mathbf{u}_\Lambda}{\|\mathbf{u}_\Lambda\|_2} + a_2 \mathbf{y}$$

with $\mathbf{u}_\Lambda^\top \mathbf{y} = 0$, $\|\mathbf{y}\|_2 = 1$ and $a_1^2 + a_2^2 = 1$. Then we have $\text{supp}(\mathbf{y}) \subset \Lambda$, and

$$\bar{\lambda} a_1 \frac{\mathbf{u}_\Lambda}{\|\mathbf{u}_\Lambda\|_2} + \bar{\lambda} a_2 \mathbf{y} = \bar{\lambda} \mathbf{w} = \mathbf{Y}_\Lambda \mathbf{w} = a_1 \frac{\mathbf{Y}_\Lambda \mathbf{u}_\Lambda}{\|\mathbf{u}_\Lambda\|_2} + a_2 \mathbf{Y}_\Lambda \mathbf{y}.$$

By taking the inner product with \mathbf{y} , we obtain

$$\bar{\lambda} a_2 = a_1 \frac{\mathbf{y}^\top \mathbf{Y}_\Lambda \mathbf{u}_\Lambda}{\|\mathbf{u}_\Lambda\|_2} + a_2 \mathbf{y}^\top \mathbf{Y}_\Lambda \mathbf{y}.$$

Since \mathbf{u}_Λ is the eigenvector of $\mathbb{E}[\mathbf{Y}_\Lambda]$ and $\mathbf{u}_\Lambda^\top \mathbf{y} = 0$, we have $\mathbf{y}^\top \mathbb{E}[\mathbf{Y}_\Lambda] \mathbf{u}_\Lambda = 0$. This leads to

$$|a_2| = |a_1| \frac{\left| \mathbf{y}^\top (\mathbf{W}_\Lambda + \mathbb{E}[\mathbf{Y}_\Lambda]) \frac{\mathbf{u}_\Lambda}{\|\mathbf{u}_\Lambda\|_2} \right|}{|\bar{\lambda} - \mathbf{y}^\top \mathbf{Y}_\Lambda \mathbf{y}|} = |a_1| \frac{\left| \mathbf{y}^\top \mathbf{W}_\Lambda \frac{\mathbf{u}_\Lambda}{\|\mathbf{u}_\Lambda\|_2} \right|}{|\bar{\lambda} - \mathbf{y}^\top \mathbf{Y}_\Lambda \mathbf{y}|}.$$

Since $\text{supp}(\mathbf{y}) \subset \Lambda$, we have $\left| \mathbf{y}^\top \mathbf{W}_\Lambda \frac{\mathbf{u}_\Lambda}{\|\mathbf{u}_\Lambda\|_2} \right| \leq \rho(\mathbf{W}, \ell)$. Moreover, since \mathbf{y} is perpendicular to \mathbf{u}_Λ , from (16) we have

$$|\mathbf{y}^\top \mathbf{Y}_\Lambda \mathbf{y}| \leq \max_{i \geq 2} |\lambda_i(\mathbf{Y}_\Lambda)| \leq \rho(\mathbf{W}, \ell).$$

So from (15) and $\rho(\mathbf{W}, \ell) < \frac{\beta}{2} \|\mathbf{u}_\Lambda\|_2^2$, we have

$$\frac{|a_2|}{|a_1|} = \frac{\left| \mathbf{y}^\top \mathbf{W}_\Lambda \frac{\mathbf{u}_\Lambda}{\|\mathbf{u}_\Lambda\|_2} \right|}{|\bar{\lambda} - \mathbf{y}^\top \mathbf{Y}_\Lambda \mathbf{y}|} \leq \frac{\rho(\mathbf{W}, \ell)}{\beta \|\mathbf{u}_\Lambda\|_2^2 - 2\rho(\mathbf{W}, \ell)}.$$

Then, since $a_1^2 + a_2^2 = 1$, we have

$$a_1^2 \geq \frac{1}{1 + \frac{\rho(\mathbf{W}, \ell)^2}{\left(\beta \|\mathbf{u}_\Lambda\|_2^2 - 2\rho(\mathbf{W}, \ell) \right)^2}},$$

which implies that

$$\begin{aligned} \text{dist}(\mathbf{w}, \mathbf{u}_\Lambda)^2 &= \min \{ \|\mathbf{u}_\Lambda - \mathbf{w}\|_2^2, \|\mathbf{u}_\Lambda + \mathbf{w}\|_2^2 \} \\ &= \|\mathbf{u}_\Lambda\|_2^2 + 1 - 2|a_1| \cdot \|\mathbf{u}_\Lambda\|_2 \\ &\leq \|\mathbf{u}_\Lambda\|_2^2 + 1 - 2 \frac{\|\mathbf{u}_\Lambda\|_2}{\sqrt{1 + \frac{\rho(\mathbf{W}, \ell)^2}{\left(\beta \|\mathbf{u}_\Lambda\|_2^2 - 2\rho(\mathbf{W}, \ell) \right)^2}}}. \end{aligned}$$

□

A.3 PROOF OF PROPOSITION 2.1

Proof of Proposition 2.1. Recall that $\mathbf{Y} = \lambda \mathbf{u} \mathbf{u}^\top + \mathbf{W}$ and $\mathbb{E}[\mathbf{Y}] = \lambda \mathbf{u} \mathbf{u}^\top$. For any $i \in \mathcal{T}$ and any $i' \in \mathcal{T}^c$, using (2)(3), we obtain

$$\begin{aligned} Y_{ii} &\geq (\mathbb{E}[\mathbf{Y}])_{ii} - |(\mathbb{E}[\mathbf{Y}])_{ii} - Y_{ii}| \geq \beta |u_i|^2 - \frac{1}{2} g_{\text{diag}}, \\ Y_{i'i'} &\leq |(\mathbb{E}[\mathbf{Y}])_{i'i'}| + |(\mathbb{E}[\mathbf{Y}])_{i'i'} - Y_{i'i'}| \leq \frac{1}{2} g_{\text{diag}}. \end{aligned}$$

Following from the fact that $g_{\text{diag}} = \beta \cdot \min_{i \in \mathcal{T}} |u_i|^2$, one has $Y_{ii} \geq Y_{i'i'}$. □

A.4 PROOF OF THEOREM 2.2

In Appendix A.4, we prove Theorem 2.2 in three steps. First, we prove that the index i_0 chosen in Algorithm 1 satisfies $|u_{i_0}| \geq \frac{\|\mathbf{u}\|_\infty}{2}$ with high probability. Second, we show that $\hat{\mathcal{T}}$ chosen in Algorithm 1 contains the indices of most of the larger nonzero entries of \mathbf{u} with high probability. Finally, we put everything together.

Step 1: Estimating $|u_{i_0}|$. Recall that $i_0 = \arg \max_{i \in [d]} Y_{ii}$.

Lemma A.5. *If β satisfies*

$$\beta \geq \frac{32\sqrt{2}}{3} \|\mathbf{u}\|_\infty^{-2} \sqrt{\log d},$$

with probability exceeding $1 - \frac{1}{4\sqrt{2\pi \log 2}} d^{-1}$, $|u_{i_0}| \geq \frac{\|\mathbf{u}\|_\infty}{2}$.

Proof. From (1)(4), for any $i \in [d]$, we obtain $Y_{ii} = \beta |u_i|^2 + W_{ii}$ and $\mathbb{E}[Y_{ii}] = \beta |u_i|^2$.

Firstly, we consider $Y_{i_* i_*}$, where i_* satisfies $|u_{i_*}| = \|\mathbf{u}\|_\infty$. Since $W_{i_* i_*} \sim \mathcal{N}(0, 2)$, using the tail of a Gaussian variable (Vershynin, 2018), it holds that, for any $\epsilon_1 > 0$,

$$\mathbb{P}\{Y_{i_* i_*} - \beta \|\mathbf{u}\|_\infty^2 \leq -\epsilon_1\} = \mathbb{P}\{W_{i_* i_*} \leq -\epsilon_1\} \leq \frac{1}{\sqrt{\pi}\epsilon_1} \exp\left(-\frac{\epsilon_1^2}{4}\right). \quad (17)$$

Secondly, we consider $\mathcal{T}_1 := \{i \in [d] : |u_i| < \frac{\|\mathbf{u}\|_\infty}{2}\}$. Since $W_{ii} \sim \mathcal{N}(0, 2)$, taking union bound and using the tail of a Gaussian variable (Vershynin, 2018), we have, for any $\epsilon_2 > 0$,

$$\begin{aligned} \mathbb{P}\left\{\max_{i \in \mathcal{T}_1} (Y_{ii} - \beta |u_i|^2) \geq \epsilon_2\right\} &\leq (d-1) \mathbb{P}\{W_{ii} \geq \epsilon_2 \text{ for some } i \in \mathcal{T}_1\} \\ &\leq \frac{d-1}{\sqrt{\pi}\epsilon_2} \exp\left(-\frac{\epsilon_2^2}{4}\right). \end{aligned} \quad (18)$$

Now we combine (17)(18) and set $\epsilon_1 = \epsilon_2 = \frac{3}{8}\beta \|\mathbf{u}\|_\infty^2$. The complementary events in (17)(18) are

$$\begin{aligned} Y_{i_* i_*} &\geq \beta \|\mathbf{u}\|_\infty^2 - \frac{3}{8}\beta \|\mathbf{u}\|_\infty^2, \\ \max_{i \in \mathcal{T}_1} (Y_{ii} - \beta |u_i|^2) &\leq \frac{3}{8}\beta \|\mathbf{u}\|_\infty^2, \end{aligned}$$

which leads to

$$\max_{i \in \mathcal{T}_1} Y_{ii} < \frac{3}{8}\beta \|\mathbf{u}\|_\infty^2 + \beta \left(\frac{\|\mathbf{u}\|_\infty}{2}\right)^2 = \beta \|\mathbf{u}\|_\infty^2 - \frac{3}{8}\beta \|\mathbf{u}\|_\infty^2 < Y_{i_* i_*} \leq Y_{i_0 i_0},$$

where we use the definition of \mathcal{T}_1 in the first inequality. It follows that $i_0 \notin \mathcal{T}_1$, i.e. $|u_{i_0}| \geq \frac{\|\mathbf{u}\|_\infty}{2}$. Therefore, using (17)(18), we obtain

$$\mathbb{P}\left\{|u_{i_0}| \geq \frac{\|\mathbf{u}\|_\infty}{2}\right\} \geq 1 - \frac{8d}{3\sqrt{\pi}\beta \|\mathbf{u}\|_\infty^2} \exp\left(-\frac{9\beta^2 \|\mathbf{u}\|_\infty^4}{256}\right), \quad (19)$$

which leads to the desired result with the condition of β . \square

Step 2: Estimating $\|\mathbf{u}_{\hat{\mathcal{T}}}\|_2$. For any $\zeta \in (0, 1]$, we define $\mathcal{T}_\zeta^- := \{i \in \mathcal{T} : |u_i| < \frac{\zeta}{2\sqrt{s}}\}$ and $\mathcal{T}_\zeta^+ = \mathcal{T} \setminus \mathcal{T}_\zeta^-$. Then we have $\|\mathbf{u}_{\mathcal{T}_\zeta^-}\|_2^2 < \frac{\zeta^2}{4s} \cdot s = \frac{\zeta^2}{4}$ and $\|\mathbf{u}_{\mathcal{T}_\zeta^+}\|_2^2 \geq 1 - \frac{\zeta^2}{4}$. Since $\|\mathbf{u}\|_\infty \geq \frac{1}{\sqrt{s}}$, Lemma A.5 implies that $|u_{i_0}| \geq \frac{1}{2\sqrt{s}} \geq \frac{\zeta}{2\sqrt{s}}$ with high probability, and thus $i_0 \in \mathcal{T}_\zeta^+$. The following lemma shows that $\mathcal{T}_\zeta^+ \subset \hat{\mathcal{T}}$ with high probability, where $\hat{\mathcal{T}}$ is chosen in Algorithm 1.

Lemma A.6. For any $\zeta \in (0, 1]$, if β satisfies

$$\beta \geq 16\zeta^{-1} \|\mathbf{u}\|_\infty^{-1} \sqrt{s(\log d + 2 \log s)},$$

with probability exceeding $1 - \frac{7+3\sqrt{2}}{6\sqrt{\pi \log 2}} d^{-1}$, $\mathcal{T}_\zeta^+ \subset \hat{\mathcal{T}}$.

Proof. It suffices to show that with high probability,

$$\min_{i \in \mathcal{T}_\zeta^+} |Y_{i, i_0}| > \max_{i \in \mathcal{T}^c} |Y_{i, i_0}|.$$

To prove this, first, we show that for any $l \in \mathcal{T}_2$, where $\mathcal{T}_2 := \{i \in \mathcal{T} : |u_i| \geq \frac{\|\mathbf{u}\|_\infty}{2}\}$,

$$\min_{i \in \mathcal{T}_\zeta^+} |Y_{il}| > \max_{i \in \mathcal{T}^c} |Y_{il}|,$$

which needs to bound $|Y_{il}|$ and $|Y_{il} - \mathbb{E}[Y_{il}]|$ for all $i \in \mathcal{T}^c$ and $i \in \mathcal{T}_\zeta^+$.

For any $l \in \mathcal{T}_2$, we first consider $\max_{i \in \mathcal{T}^c} |Y_{il}|$. From (4), for any $i \in \mathcal{T}^c$, we have $\mathbb{E}[Y_{il}] = 0$, and thus $Y_{il} = W_{il}$ by (1). Since $W_{il} \sim \mathcal{N}(0, 1)$, by taking union bound and using the tail of a Gaussian variable (Vershynin, 2018), it holds that, for any $\epsilon_3 > 0$,

$$\mathbb{P} \left\{ \max_{i \in \mathcal{T}^c} |Y_{il}| \geq \epsilon_3 \right\} \leq \frac{\sqrt{2}(d-s)}{\sqrt{\pi}\epsilon_3} \exp \left(-\frac{\epsilon_3^2}{2} \right). \quad (20)$$

Second, we consider $\min_{i \in \mathcal{T}_\zeta^+} |\mathbb{E}[Y_{il}]|$. From (4), $\mathbb{E}[Y_{il}] = \beta u_i u_l$ for any $i \in \mathcal{T}_\zeta^+$. Then, from the definition of \mathcal{T}_2 and \mathcal{T}_ζ^+ , we obtain

$$\min_{i \in \mathcal{T}_\zeta^+} |\mathbb{E}[Y_{il}]| \geq \frac{\zeta \beta \|\mathbf{u}\|_\infty}{4\sqrt{s}}. \quad (21)$$

Third, we estimate $\max_{i \in \mathcal{T}} |Y_{il} - \mathbb{E}[Y_{il}]|$. By (1), $Y_{il} = \beta u_i u_l + W_{il}$. Since $W_{il} \sim \mathcal{N}(0, 1)$ if $i \neq l$ or $W_{il} \sim \mathcal{N}(0, 2)$ if $i = l$, by taking union bound and using the tail of a Gaussian variable (Vershynin, 2018), we have, for any $\epsilon_4 > 0$,

$$\mathbb{P} \left\{ \max_{i \in \mathcal{T}} |Y_{il} - \mathbb{E}[Y_{il}]| \geq \epsilon_4 \right\} \leq \frac{2s}{\sqrt{\pi}\epsilon_4} \exp \left(-\frac{\epsilon_4^2}{4} \right). \quad (22)$$

Now we combine (20)(22) and set $\epsilon_3 = \epsilon_4 = \frac{\zeta \beta \|\mathbf{u}\|_\infty}{8\sqrt{s}}$. The complementary event in (20) is

$$\max_{i \in \mathcal{T}^c} |Y_{il}| \leq \frac{\zeta \beta \|\mathbf{u}\|_\infty}{8\sqrt{s}}.$$

Moreover, (21) and the complementary event in (22) lead to

$$|Y_{il}| > |\mathbb{E}[Y_{il}]| - |Y_{il} - \mathbb{E}[Y_{il}]| > \frac{\zeta \beta \|\mathbf{u}\|_\infty}{4\sqrt{s}} - \frac{\zeta \beta \|\mathbf{u}\|_\infty}{8\sqrt{s}} = \frac{\zeta \beta \|\mathbf{u}\|_\infty}{8\sqrt{s}}, \forall i \in \mathcal{T}_\zeta^+.$$

These two inequalities implies that $\min_{i \in \mathcal{T}_\zeta^+} |Y_{il}| > \max_{i \in \mathcal{T}^c} |Y_{il}|$ for any $l \in \mathcal{T}_2$.

Finally, by taking union bound and using (19)(20)(22), we obtain

$$\begin{aligned} & \mathbb{P} \left\{ \mathcal{T}_\zeta^+ \subset \hat{\mathcal{T}} \right\} \\ &= \sum_{l \in \mathcal{T}_2} \mathbb{P} \left\{ \min_{i \in \mathcal{T}_\zeta^+} |Y_{il}| > \max_{i \in \mathcal{T}^c} |Y_{il}|, i_0 = l \right\} \\ &\geq \sum_{l \in \mathcal{T}_2} (1 - \mathbb{P} \left\{ \min_{i \in \mathcal{T}_\zeta^+} |Y_{il}| \leq \max_{i \in \mathcal{T}^c} |Y_{il}| \right\} - \mathbb{P} \{i_0 \neq l\}) \\ &\geq \sum_{l \in \mathcal{T}_2} (\mathbb{P} \{i_0 = l\} - \mathbb{P} \left\{ \min_{i \in \mathcal{T}_\zeta^+} |Y_{il}| \leq \max_{i \in \mathcal{T}^c} |Y_{il}| \right\}) \\ &\geq 1 - \frac{8d}{3\sqrt{\pi}\beta \|\mathbf{u}\|_\infty^2} \exp \left(-\frac{9\beta^2 \|\mathbf{u}\|_\infty^4}{256} \right) - \frac{8\sqrt{2}ss(d-s)}{\sqrt{\pi}\zeta\beta \|\mathbf{u}\|_\infty} \exp \left(-\frac{\zeta^2\beta^2 \|\mathbf{u}\|_\infty^2}{128s} \right) \\ &\quad - \frac{16\sqrt{s}s^2}{\sqrt{\pi}\zeta\beta \|\mathbf{u}\|_\infty} \exp \left(-\frac{\zeta^2\beta^2 \|\mathbf{u}\|_\infty^2}{256s} \right). \end{aligned} \quad (23)$$

Since $\zeta \in (0, 1]$, (23) leads to the desired result with the condition of β . \square

Step 3: Putting everything together. Now we estimate $\text{dist}(\hat{\mathbf{u}}, \mathbf{u})$ and prove Theorem 2.2.

Proof of Theorem 2.2. For simplicity, we denote $\rho = \rho(\mathbf{W}, s)$. By applying Lemma A.3 with $r = \frac{1}{16}\zeta$, $\ell = s$ and Lemma A.6, if β satisfies

$$\beta \geq \max \left\{ 16\zeta^{-1} \|\mathbf{u}\|_{\infty}^{-1} \sqrt{s(\log d + 2 \log s)}, 32\zeta^{-1} \sqrt{\log d + s \log\left(\frac{9ed}{s}\right)} \right\}, \quad (24)$$

then we have

$$\mathbb{P} \left\{ \rho \leq \frac{3}{16}\zeta\beta, i_0 \in \mathcal{T}_{\zeta}^+ \subset \widehat{\mathcal{T}} \right\} \geq 1 - \frac{1}{\sqrt{\pi \log(18e)}} d^{-1} - \frac{7 + 3\sqrt{2}}{6\sqrt{\pi \log 2}} d^{-1} > 1 - 1.5558d^{-1}. \quad (25)$$

Under the event in (25), we estimate $\text{dist}(\hat{\mathbf{u}}, \mathbf{u})$. Since $\text{supp}(\hat{\mathbf{u}}) = \widehat{\mathcal{T}}$, we have

$$\text{dist}(\hat{\mathbf{u}}, \mathbf{u})^2 = \text{dist}(\hat{\mathbf{u}}, \mathbf{u}_{\widehat{\mathcal{T}}})^2 + \|\mathbf{u}_{\widehat{\mathcal{T}}^c}\|_2^2. \quad (26)$$

Firstly, we estimate $\|\mathbf{u}_{\widehat{\mathcal{T}}^c}\|_2^2$. Since $\widehat{\mathcal{T}}^c \subset (\mathcal{T} \setminus \mathcal{T}_{\zeta}^-)^c = \mathcal{T}_{\zeta}^- \cup \mathcal{T}^c$, we have

$$\|\mathbf{u}_{\widehat{\mathcal{T}}^c}\|_2^2 \leq \|\mathbf{u}_{\mathcal{T}_{\zeta}^-}\|_2^2 + \|\mathbf{u}_{\mathcal{T}^c}\|_2^2 < \frac{\zeta^2}{4} < \frac{1}{4}, \quad \|\mathbf{u}_{\widehat{\mathcal{T}}}\|_2^2 > 1 - \frac{\zeta^2}{4} > \frac{3}{4}.$$

Secondly, we estimate $\text{dist}(\hat{\mathbf{u}}, \mathbf{u}_{\widehat{\mathcal{T}}})^2$. Applying Lemma A.4 with $\Lambda = \widehat{\mathcal{T}}$ and $\ell = s$, we obtain

$$\text{dist}(\hat{\mathbf{u}}, \mathbf{u}_{\widehat{\mathcal{T}}})^2 \leq \|\mathbf{u}_{\widehat{\mathcal{T}}}\|_2 + 1 - 2 \frac{\|\mathbf{u}_{\widehat{\mathcal{T}}}\|_2}{\sqrt{1 + \frac{\rho^2}{(\frac{3}{4}\beta - 2\rho)^2}}} \leq \|\mathbf{u}_{\widehat{\mathcal{T}}}\|_2 + 1 - 2 \frac{\|\mathbf{u}_{\widehat{\mathcal{T}}}\|_2}{\sqrt{1 + \frac{\rho^2}{(\frac{3}{4}\beta - 2\rho)^2}}},$$

where the last inequality holds since $\|\mathbf{u}_{\widehat{\mathcal{T}}}\|_2^2 > \frac{3}{4}$. Therefore, using Lemma A.4 and $\|\mathbf{u}_{\widehat{\mathcal{T}}^c}\|_2^2 \leq \frac{\zeta^2}{4}$, we have

$$\begin{aligned} \text{dist}(\hat{\mathbf{u}}, \mathbf{u}_{\widehat{\mathcal{T}}})^2 &\leq \max \left\{ 2 - \frac{\zeta^2}{4} - \frac{2\sqrt{1 - \frac{\zeta^2}{4}}}{\sqrt{1 + \frac{\rho^2}{(\frac{3}{4}\beta - 2\rho)^2}}}, 2 - \frac{2}{\sqrt{1 + \frac{\rho^2}{(\frac{3}{4}\beta - 2\rho)^2}}} \right\} \\ &\leq \max \left\{ 2 - \frac{\zeta^2}{4} - 2 \frac{1 - \frac{\zeta^2}{4}}{1 + \frac{\rho^2}{(\frac{3}{4}\beta - 2\rho)^2}}, 2 - 2 \frac{1}{1 + \frac{\rho^2}{(\frac{3}{4}\beta - 2\rho)^2}} \right\} \\ &= \max \left\{ \frac{\frac{\zeta^2}{4}(\frac{3}{4}\beta - 2\rho)^2 + (2 - \frac{\zeta^2}{4})\rho^2}{(\frac{3}{4}\beta - 2\rho)^2 + \rho^2}, \frac{2\rho^2}{(\frac{3}{4}\beta - 2\rho)^2 + \rho^2} \right\} \\ &\leq \frac{\zeta^2}{4} + \frac{2\rho^2}{(\frac{3}{4}\beta - 2\rho)^2 + \rho^2}. \end{aligned}$$

It follows from (26) and $\rho \leq \frac{3}{16}\zeta\beta$ that

$$\text{dist}(\hat{\mathbf{u}}, \mathbf{u})^2 \leq \frac{\zeta^2}{2} + \frac{\zeta^2}{2} = \zeta^2,$$

completing the proof. \square

Remark A.7. From (24), a sufficient condition for the constant C_1 in Theorem 2.2 is

$$C_1 \geq \max \left\{ 16\sqrt{3}, 32\sqrt{2 + \log_2(9e)} \right\} = 32\sqrt{2 + \log_2(9e)}.$$

A.5 PROOF OF THEOREM 2.3

Since Algorithm 1 and Algorithm 2 have the same step for support estimation, we use some results and techniques in Section A.4 to prove Theorem 2.3. Specifically, it requires Lemma A.5 and Lemma A.6, which show that with high probability, $|u_{i_0}| \geq \frac{\|\mathbf{u}\|_{\infty}}{2}$ and $\mathcal{T}_{\zeta}^+ \subset \widehat{\mathcal{T}}$. Recall that

$$\hat{\mathbf{u}}_{\text{nv}} = \mathbf{Y}_{\widehat{\mathcal{T}}, i_0} / \|\mathbf{Y}_{\widehat{\mathcal{T}}, i_0}\|_2 \text{ and } \mathcal{T}_{\zeta}^+ = \left\{ i \in \mathcal{T} : |u_i| \geq \frac{\zeta}{2\sqrt{s}} \right\}.$$

Proof. From (6), similar to (26), we have

$$\text{dist}(\hat{\mathbf{u}}_{\text{nv}}, \mathbf{u})^2 = \text{dist}(\hat{\mathbf{u}}_{\text{nv}}, \mathbf{u}_{\hat{\mathcal{T}}})^2 + \|\mathbf{u}_{\hat{\mathcal{T}}^c}\|_2^2. \quad (27)$$

Recall that $\|\mathbf{u}_{\hat{\mathcal{T}}^c}\|_2^2 \leq \frac{\zeta^2}{4}$ and $\|\mathbf{u}_{\hat{\mathcal{T}}}\|_2^2 \geq 1 - \frac{\zeta^2}{4}$ from the proof of Theorem 2.2, hence we only need to estimate $\text{dist}(\hat{\mathbf{u}}_{\text{nv}}, \mathbf{u}_{\hat{\mathcal{T}}})^2$.

From the definition of $\hat{\mathbf{u}}_{\text{nv}}$, we obtain

$$\text{dist}(\hat{\mathbf{u}}_{\text{nv}}, \mathbf{u}_{\hat{\mathcal{T}}})^2 = \text{dist}\left(\frac{\mathbf{Y}_{\hat{\mathcal{T}}, i_0}}{\|\mathbf{Y}_{\hat{\mathcal{T}}, i_0}\|_2}, \mathbf{u}_{\hat{\mathcal{T}}}\right)^2,$$

To handle the randomness of i_0 . Similarly to Lemma A.6, we estimate

$$\text{dist}\left(\frac{\mathbf{Y}_{\hat{\mathcal{T}}, l}}{\|\mathbf{Y}_{\hat{\mathcal{T}}, l}\|_2}, \mathbf{u}_{\hat{\mathcal{T}}}\right)^2 \quad (28)$$

for any $l \in \mathcal{T}_2 = \{i \in \mathcal{T} : |u_i| \geq \frac{\|\mathbf{u}\|_\infty}{2}\}$. Without loss of generality, we assume $u_l > 0$ and consider

$$\left\| \frac{\mathbf{Y}_{\hat{\mathcal{T}}, l}}{\|\mathbf{Y}_{\hat{\mathcal{T}}, l}\|_2} - \mathbf{u}_{\hat{\mathcal{T}}} \right\|_2^2, \quad (29)$$

which is an upper bound of (28).

We begin with a simplification of (29):

$$\begin{aligned} \left\| \frac{\mathbf{Y}_{\hat{\mathcal{T}}, l}}{\|\mathbf{Y}_{\hat{\mathcal{T}}, l}\|_2} - \mathbf{u}_{\hat{\mathcal{T}}} \right\|_2^2 &= \sum_{i \in \hat{\mathcal{T}}} \left| \frac{Y_{il}}{\|\mathbf{Y}_{\hat{\mathcal{T}}, l}\|_2} - u_i \right|^2 \\ &= \sum_{i \in \hat{\mathcal{T}}} \frac{|Y_{il} - \beta u_i u_l + u_i(\beta u_l - \|\mathbf{Y}_{\hat{\mathcal{T}}, l}\|_2)|^2}{\|\mathbf{Y}_{\hat{\mathcal{T}}, l}\|_2^2} \\ &\leq 2 \sum_{i \in \hat{\mathcal{T}}} \frac{|Y_{il} - \beta u_i u_l|^2 + |u_i|^2 |\beta u_l - \|\mathbf{Y}_{\hat{\mathcal{T}}, l}\|_2|^2}{\|\mathbf{Y}_{\hat{\mathcal{T}}, l}\|_2^2} \\ &\leq 2s \frac{\max_{i \in [d]} |Y_{il} - \mathbb{E}[Y_{il}]|^2}{\|\mathbf{Y}_{\hat{\mathcal{T}}, l}\|_2^2} + 2 \frac{|\beta u_l - \|\mathbf{Y}_{\hat{\mathcal{T}}, l}\|_2|^2}{\|\mathbf{Y}_{\hat{\mathcal{T}}, l}\|_2^2}, \end{aligned} \quad (30)$$

where we use $(a_1 + a_2)^2 \leq 2(a_1^2 + a_2^2)$ in the first inequality and use $\|\mathbf{u}_{\hat{\mathcal{T}}}\|_2^2 \leq 1$ in the last inequality.

To estimate (30), we consider the following event:

$$\left\{ \max_{i \in \mathcal{T}^c} |Y_{il}| \leq \epsilon_5 = \frac{\varsigma_2 \zeta^2 \beta \|\mathbf{u}\|_\infty}{\sqrt{s}}, \max_{i \in \mathcal{T}} |Y_{il} - \mathbb{E}[Y_{il}]| \leq \epsilon_5, \mathcal{T}_\zeta^+ \subset \hat{\mathcal{T}} \right\}, \quad (31)$$

where $\varsigma_2 > 0$ is a constant close to 0. This event is related to (20)(22)(23).

Under the event in (31), we first estimate the lower bound of $\|\mathbf{Y}_{\hat{\mathcal{T}}, l}\|_2^2$. It holds that

$$\begin{aligned} \|\mathbf{Y}_{\hat{\mathcal{T}}, l}\|_2^2 &\geq \sum_{i \in \mathcal{T}_\zeta^+} (\beta |u_i u_l| - |\beta u_i u_l - Y_{il}|)^2 \\ &\geq \sum_{i \in \mathcal{T}_\zeta^+} (\beta |u_i u_l| - \frac{\varsigma_2 \zeta^2 \beta \|\mathbf{u}\|_\infty}{\sqrt{s}})^2 \\ &\geq \sum_{i \in \mathcal{T}_\zeta^+} (\beta |u_i u_l| - 4\varsigma_2 \beta |u_i u_l|)^2 \\ &\geq \frac{3}{4} (1 - 4\varsigma_2)^2 \beta^2 |u_l|^2, \end{aligned} \quad (32)$$

where $\mathcal{T}_\zeta^+ \subset \widehat{\mathcal{T}}$ and triangle inequality are used in the first inequality, the third inequality holds by $l \in \mathcal{T}_2$ and $i \in \mathcal{T}_\zeta^+$ and the last inequality holds by $\|\mathbf{u}_{\mathcal{T}_\zeta^+}\|_2 \geq 1 - \frac{\zeta^2}{4} \geq \frac{3}{4}$.

Second, we estimate $\left| \beta u_l - \|\mathbf{Y}_{\widehat{\mathcal{T}},l}\|_2 \right|^2$. We obtain

$$\begin{aligned} \left| \beta u_l - \|\mathbf{Y}_{\widehat{\mathcal{T}},l}\|_2 \right|^2 &= \beta^2 |u_l|^2 - 2\beta u_l \|\mathbf{Y}_{\widehat{\mathcal{T}},l}\|_2 + \|\mathbf{Y}_{\widehat{\mathcal{T}},l}\|_2^2 \\ &\leq \beta^2 |u_l|^2 - 2\beta u_l \sqrt{\sum_{i \in \widehat{\mathcal{T}}} (\beta |u_i u_l| - \epsilon_5)^2} + \sum_{i \in \widehat{\mathcal{T}}} (\beta |u_i u_l| + \epsilon_5)^2 \\ &\leq \beta^2 |u_l|^2 - 2\beta u_l \sqrt{\beta^2 |u_l|^2 \|\mathbf{u}_{\widehat{\mathcal{T}}}\|_2^2 - 2\epsilon_5 \beta u_l \|\mathbf{u}_{\widehat{\mathcal{T}}}\|_1 + s\epsilon_5^2} \\ &\quad + \beta^2 |u_l|^2 \|\mathbf{u}_{\widehat{\mathcal{T}}}\|_2^2 + 2\epsilon_5 \beta u_l \|\mathbf{u}_{\widehat{\mathcal{T}}}\|_1 + s\epsilon_5^2, \end{aligned}$$

where the first inequality holds similar to (32). Thus we have

$$\begin{aligned} &\left| \beta u_l - \|\mathbf{Y}_{\widehat{\mathcal{T}},l}\|_2 \right|^2 \\ &\leq 2\beta u_l \left(\beta u_l \frac{\|\mathbf{u}_{\widehat{\mathcal{T}}}\|_2^2 + 1}{2} + \epsilon_5 \|\mathbf{u}_{\widehat{\mathcal{T}}}\|_1 - \sqrt{\beta^2 |u_l|^2 \|\mathbf{u}_{\widehat{\mathcal{T}}}\|_2^2 - 2\epsilon_5 \beta u_l \|\mathbf{u}_{\widehat{\mathcal{T}}}\|_1 + s\epsilon_5^2} \right) + s\epsilon_5^2, \end{aligned} \quad (33)$$

To complete this estimation, we compute

$$\begin{aligned} &\beta u_l \frac{\|\mathbf{u}_{\widehat{\mathcal{T}}}\|_2^2 + 1}{2} + \epsilon_5 \|\mathbf{u}_{\widehat{\mathcal{T}}}\|_1 - \sqrt{\beta^2 |u_l|^2 \|\mathbf{u}_{\widehat{\mathcal{T}}}\|_2^2 - 2\epsilon_5 \beta u_l \|\mathbf{u}_{\widehat{\mathcal{T}}}\|_1 + s\epsilon_5^2} \\ &= \frac{(\beta u_l \frac{\|\mathbf{u}_{\widehat{\mathcal{T}}}\|_2^2 + 1}{2} + \epsilon_5 \|\mathbf{u}_{\widehat{\mathcal{T}}}\|_1)^2 - (\beta^2 |u_l|^2 \|\mathbf{u}_{\widehat{\mathcal{T}}}\|_2^2 - 2\epsilon_5 \beta u_l \|\mathbf{u}_{\widehat{\mathcal{T}}}\|_1 + s\epsilon_5^2)}{\beta u_l + \epsilon_5 \|\mathbf{u}_{\widehat{\mathcal{T}}}\|_1 + \sqrt{\beta^2 |u_l|^2 \|\mathbf{u}_{\widehat{\mathcal{T}}}\|_2^2 - 2\epsilon_5 \beta u_l \|\mathbf{u}_{\widehat{\mathcal{T}}}\|_1 + s\epsilon_5^2}} \\ &\leq \frac{1}{\beta u_l} \left(\beta^2 |u_l|^2 \frac{(1 - \|\mathbf{u}_{\widehat{\mathcal{T}}}\|_2^2)^2}{4} + 4\epsilon_5 \beta u_l \|\mathbf{u}_{\widehat{\mathcal{T}}}\|_1 + (\|\mathbf{u}_{\widehat{\mathcal{T}}}\|_1^2 - s)\epsilon_5^2 \right), \\ &\leq \beta u_l \frac{\zeta^4}{64} + 4\epsilon_5 \sqrt{s}, \end{aligned}$$

where we use $\|\mathbf{u}_{\widehat{\mathcal{T}}}\|_1 \leq \sqrt{s}$ and $1 - \frac{\zeta^2}{4} \leq \|\mathbf{u}_{\widehat{\mathcal{T}}}\|_2^2 \leq 1$ in the last inequality. It follows that (33) can be simplified as

$$\left| \beta u_l - \|\mathbf{Y}_{\widehat{\mathcal{T}},l}\|_2 \right|^2 \leq \beta^2 |u_l|^2 \frac{\zeta^4}{32} + 8\epsilon_5 \beta u_l \sqrt{s} + s\epsilon_5^2. \quad (34)$$

Therefore, under the event in (31), combining (28)(30)(32)(34), we have

$$\begin{aligned} &\text{dist}\left(\frac{\mathbf{Y}_{\widehat{\mathcal{T}},l}}{\|\mathbf{Y}_{\widehat{\mathcal{T}},l}\|_2}, \mathbf{u}_{\widehat{\mathcal{T}}}\right)^2 \\ &\leq \frac{2s \frac{\zeta^4 \beta^2 \|\mathbf{u}_{\widehat{\mathcal{T}}}\|_\infty^2}{s}}{\frac{3}{4}(1 - 4\zeta_2)^2 \beta^2 |u_l|^2} + 2 \frac{\beta^2 |u_l|^2 \frac{\zeta^4}{32} + 8 \frac{\zeta_2 \zeta^2 \beta \|\mathbf{u}_{\widehat{\mathcal{T}}}\|_\infty}{\sqrt{s}} \beta u_l \sqrt{s} + s \frac{\zeta^2 \zeta^4 \beta^2 \|\mathbf{u}_{\widehat{\mathcal{T}}}\|_\infty^2}{s}}{\frac{3}{4}(1 - 4\zeta_2)^2 \beta^2 |u_l|^2} \\ &\leq \left(\frac{32\zeta_2^2}{3(1 - 4\zeta_2)^2} + \frac{1}{12(1 - 4\zeta_2)^2} + \frac{128\zeta_2}{3(1 - 4\zeta_2)^2} + \frac{32\zeta_2^2}{3(1 - 4\zeta_2)^2} \right) \zeta^2 \\ &\leq \frac{3}{4} \zeta^2, \end{aligned} \quad (35)$$

where in the second inequality we use $|u_l| \geq \frac{\|\mathbf{u}_{\widehat{\mathcal{T}}}\|_\infty}{2}$ and $0 < \zeta \leq 1$, and in the last equality we set $\zeta_2 = 1/69$. Then, according to the event in (31), using (20)(22) with $\epsilon_3 = \epsilon_4 = \frac{\zeta^2 \beta \|\mathbf{u}_{\widehat{\mathcal{T}}}\|_\infty}{69\sqrt{s}}$ and (23),

it holds that

$$\begin{aligned} \mathbb{P} \left\{ \text{dist} \left(\frac{\mathbf{Y}_{\hat{\mathcal{T}},l}}{\|\mathbf{Y}_{\hat{\mathcal{T}},l}\|_2}, \mathbf{u}_{\hat{\mathcal{T}}} \right) \leq \zeta \right\} &\geq 1 - \frac{69\sqrt{2}s(d-s)}{\sqrt{\pi}\zeta^2\beta\|\mathbf{u}\|_\infty} \exp \left(-\frac{\zeta^4\beta^2\|\mathbf{u}\|_\infty^2}{9522s} \right) - \frac{138\sqrt{ss}}{\sqrt{\pi}\zeta\beta\|\mathbf{u}\|_\infty} \exp \left(-\frac{\zeta^4\beta^2\|\mathbf{u}\|_\infty^2}{19044s} \right) \\ &\quad - \frac{8d}{3\sqrt{\pi}\beta\|\mathbf{u}\|_\infty^2} \exp \left(-\frac{9\beta^2\|\mathbf{u}\|_\infty^4}{256} \right) - \frac{8\sqrt{2}ss(d-s)}{\sqrt{\pi}\zeta\beta\|\mathbf{u}\|_\infty} \exp \left(-\frac{\zeta^2\beta^2\|\mathbf{u}\|_\infty^2}{128s} \right) \\ &\quad - \frac{16\sqrt{ss}^2}{\sqrt{\pi}\zeta\beta\|\mathbf{u}\|_\infty} \exp \left(-\frac{\zeta^2\beta^2\|\mathbf{u}\|_\infty^2}{256s} \right). \end{aligned} \quad (36)$$

Finally, similar to (23), taking the union bound and using (19)(27)(36), we obtain

$$\begin{aligned} &\mathbb{P} \left\{ \text{dist}(\hat{\mathbf{u}}_{\text{nv}}, \mathbf{u}_{\hat{\mathcal{T}}}) \leq \zeta \right\} \\ &= \sum_{l \in \mathcal{T}_2} \mathbb{P} \left\{ \text{dist} \left(\frac{\mathbf{Y}_{\hat{\mathcal{T}},l}}{\|\mathbf{Y}_{\hat{\mathcal{T}},l}\|_2}, \mathbf{u}_{\hat{\mathcal{T}}} \right) \leq \zeta, i_0 = l \right\} \\ &\geq \sum_{l \in \mathcal{T}_2} (1 - \mathbb{P} \left\{ \text{dist} \left(\frac{\mathbf{Y}_{\hat{\mathcal{T}},l}}{\|\mathbf{Y}_{\hat{\mathcal{T}},l}\|_2}, \mathbf{u}_{\hat{\mathcal{T}}} \right) > \zeta \right\} - \mathbb{P} \{i_0 \neq l\}) \\ &\geq \sum_{l \in \mathcal{T}_2} (\mathbb{P} \{i_0 = l\} - \mathbb{P} \left\{ \text{dist} \left(\frac{\mathbf{Y}_{\hat{\mathcal{T}},l}}{\|\mathbf{Y}_{\hat{\mathcal{T}},l}\|_2}, \mathbf{u}_{\hat{\mathcal{T}}} \right) > \zeta \right\}) \\ &\geq 1 - \frac{69\sqrt{2}ss(d-s)}{\sqrt{\pi}\zeta^2\beta\|\mathbf{u}\|_\infty} \exp \left(-\frac{\zeta^4\beta^2\|\mathbf{u}\|_\infty^2}{9522s} \right) - \frac{138\sqrt{ss}^2}{\sqrt{\pi}\zeta\beta\|\mathbf{u}\|_\infty} \exp \left(-\frac{\zeta^4\beta^2\|\mathbf{u}\|_\infty^2}{19044s} \right) \\ &\quad - \frac{8(1+s)d}{3\sqrt{\pi}\beta\|\mathbf{u}\|_\infty^2} \exp \left(-\frac{9\beta^2\|\mathbf{u}\|_\infty^4}{256} \right) - \frac{8\sqrt{2}ss^2(d-s)}{\sqrt{\pi}\zeta\beta\|\mathbf{u}\|_\infty} \exp \left(-\frac{\zeta^2\beta^2\|\mathbf{u}\|_\infty^2}{128s} \right) \\ &\quad - \frac{16\sqrt{ss}^3}{\sqrt{\pi}\zeta\beta\|\mathbf{u}\|_\infty} \exp \left(-\frac{\zeta^2\beta^2\|\mathbf{u}\|_\infty^2}{256s} \right). \end{aligned}$$

If β satisfies

$$\beta \geq 138\zeta^{-2}\|\mathbf{u}\|_\infty^{-1}\sqrt{s(\log d + 2\log s)}, \quad (37)$$

it implies that

$$\mathbb{P} \left\{ \text{dist}(\hat{\mathbf{u}}_{\text{nv}}, \mathbf{u}_{\hat{\mathcal{T}}}) \leq \zeta \right\} \geq 1 - \frac{231\sqrt{2} + 470}{414\sqrt{\pi}\log 2} d^{-1} > 1 - 1.3041d^{-1}.$$

□

Remark A.8. From (37), a sufficient condition for the constant C_2 in Theorem 2.3 is

$$C_2 \geq 138\sqrt{3}.$$

A.6 PROOF OF THEOREM 2.4

Proof. Recall that $\mathcal{T}_\zeta^+ = \{i \in \mathcal{T} : |u_i| \geq \frac{\zeta}{2\sqrt{s}}\}$. From the assumption of \mathbf{u} , we have $\mathcal{T} = \mathcal{T}_{2\theta}^+$.

Therefore, using (23) and $|\hat{\mathcal{T}}| = |\mathcal{T}| = s$, we obtain

$$\begin{aligned} \mathbb{P} \left\{ \hat{\mathcal{T}} = \mathcal{T} \right\} &\geq 1 - \frac{8d}{3\sqrt{\pi}\beta\|\mathbf{u}\|_\infty^2} \exp \left(-\frac{9\beta^2\|\mathbf{u}\|_\infty^4}{256} \right) - \frac{4\sqrt{2}ss(d-s)}{\sqrt{\pi}\theta\beta\|\mathbf{u}\|_\infty} \exp \left(-\frac{\theta^2\beta^2\|\mathbf{u}\|_\infty^2}{32s} \right) \\ &\quad - \frac{8\sqrt{ss}^2}{\sqrt{\pi}\theta\beta\|\mathbf{u}\|_\infty} \exp \left(-\frac{\theta^2\beta^2\|\mathbf{u}\|_\infty^2}{64s} \right). \end{aligned}$$

If β satisfies

$$\beta \geq \max \left\{ \frac{32\sqrt{2}}{3}\|\mathbf{u}\|_\infty^{-2}\sqrt{\log d}, 8\theta^{-1}\|\mathbf{u}\|_\infty^{-1}\sqrt{s(\log d + 2\log s)} \right\}, \quad (38)$$

it implies that

$$\mathbb{P}\left\{\widehat{\mathcal{T}} = \mathcal{T}\right\} \geq 1 - \frac{5 + 4\sqrt{2}}{4\sqrt{2\pi\log 2}}d^{-1} > 1 - 1.2766d^{-1}.$$

□

Remark A.9. From (38), a sufficient condition for the constant C_3 in Theorem 2.4 is

$$C_3 \geq \frac{32\sqrt{2}}{3}.$$

A.7 PROOF OF THEOREM 3.1

The proof of Theorem 3.1 is organized into two parts. First, we show that \mathbf{u}^0 falls into a small constant neighborhood of \mathbf{u} . Subsequently, we prove the convergence of the truncated power method.

Proof of Theorem 3.1. We denote $\tilde{s} = s + 2k$, $\rho = \rho(\mathbf{W}, \tilde{s})$ and $\mathcal{F}_t = \text{supp}(\mathbf{u}^t)$, where $k = C_5 s$ for some absolute constant $C_5 \geq 1$. Similar to the proof of Theorem 3.1, setting $r = 0.01\zeta'$ and $\ell = \tilde{s}$ in Lemma A.3 with $\zeta' \in (0, 1)$ and $\zeta = 1$ in Lemma A.6, if β satisfies

$$\beta \geq \max \left\{ 16\|\mathbf{u}\|_\infty^{-1} \sqrt{s(\log d + 2\log s)}, 200(\zeta')^{-1} \sqrt{(1 + 2C_5)s \log \left(\frac{9ed}{(1 + 2C_5)s} \right) + \log d} \right\}, \quad (39)$$

then with the probability exceeding

$$1 - \frac{1}{\sqrt{\pi \log(432e^3)}}d^{-1} - \frac{7 + 3\sqrt{2}}{6\sqrt{\pi \log 2}}d^{-1} > 1 - 1.4571d^{-1},$$

the following event holds:

$$\{\rho \leq 0.03\zeta'\beta, \text{dist}(\mathbf{u}^0, \mathbf{u}) \leq 1\}.$$

We will continue the proof under this event.

Step 1: Estimating $|\mathbf{u}^\top \mathbf{u}^0|$. Since $1 \geq \text{dist}(\mathbf{u}, \mathbf{u}^0)^2 = 2 - 2|\mathbf{u}^\top \mathbf{u}^0|$, we have $|\mathbf{u}^\top \mathbf{u}^0| \geq 0.5$.

Step 2: Convergence of truncated power method. To prove (11), we will first show that $\text{dist}(\mathbf{u}, \mathbf{u}^t) \leq 1$ by induction.

We denote $\Lambda_t = \mathcal{F}_{t-1} \cup \mathcal{F}_t \cup \mathcal{T}$, then $|\Lambda_t| \leq s + 2k = \tilde{s}$. Also, we define

$$\mathbf{w}^t = \mathbf{Y}_{\Lambda_t} \mathbf{u}^{t-1} / \|\mathbf{Y}_{\Lambda_t} \mathbf{u}^{t-1}\|_2, \quad (40)$$

hence we have $\mathbf{u}^t = \mathbf{w}^t / \|\mathbf{w}^t\|_2$ and \mathcal{F}_t is the set of indices with the k largest absolute values in \mathbf{w}^t . Let κ be the ratio of the second largest (in absolute value) to the largest eigenvalue of \mathbf{Y}_{Λ_t} . Then, since $\mathcal{T} \subset \Lambda_t$, similar to (15)(16), we obtain

$$\kappa = \frac{\max_{i \neq 1} |\lambda_i(\mathbf{Y}_{\Lambda_t})|}{|\lambda_1(\mathbf{Y}_{\Lambda_t})|} \leq \frac{\rho}{\beta \|\mathbf{u}_{\Lambda_t}\|_2^2 - \rho} \leq \frac{0.03\beta}{\beta - 0.03\beta} = \frac{3}{97} < 1,$$

where in the second inequality we use $\rho \leq 0.03\zeta'\beta$ and $\zeta' < 1$.

Let $\bar{\mathbf{u}}$ be a unit eigenvector corresponding to the largest eigenvalue of \mathbf{Y}_{Λ_t} and satisfying $\mathbf{u}^\top \bar{\mathbf{u}} \geq 0$. hence we have $\text{dist}(\mathbf{u}, \bar{\mathbf{u}}) = \|\mathbf{u} - \bar{\mathbf{u}}\|_2$. Then, using (40) and Lemma A.1, we have

$$|\bar{\mathbf{u}}^\top \mathbf{w}^t| \geq |\bar{\mathbf{u}}^\top \mathbf{u}^{t-1}| \left(1 + \frac{1}{2}(1 - \kappa^2)(1 - |\bar{\mathbf{u}}^\top \mathbf{u}^{t-1}|^2) \right),$$

which implies that

$$1 - |\bar{\mathbf{u}}^\top \mathbf{w}^t| \leq \left(1 - |\bar{\mathbf{u}}^\top \mathbf{u}^{t-1}| \right) \left(1 - \frac{1 - \kappa^2}{2} \left(|\bar{\mathbf{u}}^\top \mathbf{u}^{t-1}| + |\bar{\mathbf{u}}^\top \mathbf{u}^{t-1}|^2 \right) \right). \quad (41)$$

Since $\mathcal{T} \subset \Lambda_t$, Lemma A.4 gives

$$\begin{aligned} \|\mathbf{u} - \bar{\mathbf{u}}\|_2^2 &= \text{dist}(\mathbf{u}, \bar{\mathbf{u}})^2 \leq 2 - 2 \frac{1}{\sqrt{1 + \frac{\rho^2}{(\beta - 2\rho)^2}}} \\ &\leq \frac{\rho^2}{(\beta - 2\rho)^2} \leq \frac{(0.03\zeta'\beta)^2}{(\beta - 0.06\beta)^2} = \frac{9(\zeta')^2}{8836}, \end{aligned} \quad (42)$$

where in the second inequality we use $1 - \frac{1}{\sqrt{1+a}} \leq \frac{a}{2}$ for $a \geq 0$, and in the last two inequalities we use $\rho \leq 0.03\zeta'\beta$ and $\zeta' < 1$. Note that the induction assumption $\text{dist}(\mathbf{u}, \mathbf{u}^{t-1}) \leq 1$ implies that $|\mathbf{u}^\top \mathbf{u}^{t-1}| \geq 0.5$, which with (42) further leads to

$$\begin{aligned} |\bar{\mathbf{u}}^\top \mathbf{u}^{t-1}| &\geq |\mathbf{u}^\top \mathbf{u}^{t-1}| - |(\mathbf{u} - \bar{\mathbf{u}})^\top \mathbf{u}^{t-1}| \\ &\geq |\mathbf{u}^\top \mathbf{u}^{t-1}| - \|\mathbf{u} - \bar{\mathbf{u}}\|_2 \|\mathbf{u}^{t-1}\|_2 \geq 0.5 - \frac{3}{94}. \end{aligned} \quad (43)$$

Plugging (43) into (41), we have

$$1 - |\bar{\mathbf{u}}^\top \mathbf{w}^t| \leq 0.6568(1 - |\bar{\mathbf{u}}^\top \mathbf{u}^{t-1}|),$$

which is equivalent to

$$\text{dist}(\bar{\mathbf{u}}, \mathbf{w}^t) \leq 0.8105 \cdot \text{dist}(\bar{\mathbf{u}}, \mathbf{u}^{t-1}), \quad (44)$$

where we use $\|\bar{\mathbf{u}}\|_2 = \|\mathbf{w}^t\|_2 = \|\mathbf{u}^{t-1}\|_2 = 1$. For unit vectors $\bar{\mathbf{u}}, \mathbf{u}^{t-1}, \mathbf{u}$, we obtain

$$\text{dist}(\bar{\mathbf{u}}, \mathbf{u}^{t-1}) \leq \text{dist}(\bar{\mathbf{u}}, \mathbf{u}) + \text{dist}(\mathbf{u}^{t-1}, \mathbf{u}). \quad (45)$$

This is because

$$\begin{aligned} \text{dist}(\bar{\mathbf{u}}, \mathbf{u}) + \text{dist}(\mathbf{u}^{t-1}, \mathbf{u}) &= \|\tau_1 \bar{\mathbf{u}} - \mathbf{u}\|_2 + \|\mathbf{u} + \tau_2 \mathbf{u}^{t-1}\|_2 \\ &\geq \|\tau_1 \bar{\mathbf{u}} + \tau_2 \mathbf{u}^{t-1}\|_2 \\ &\geq \text{dist}(\bar{\mathbf{u}}, \mathbf{u}^{t-1}), \end{aligned}$$

where $\tau_1, \tau_2 \in \{\pm 1\}$ and we use (6). Similarly, for unit vectors $\mathbf{u}, \mathbf{w}^t, \bar{\mathbf{u}}$, it holds that

$$\text{dist}(\mathbf{u}, \mathbf{w}^t) \leq \text{dist}(\mathbf{u}, \bar{\mathbf{u}}) + \text{dist}(\mathbf{w}^t, \bar{\mathbf{u}}). \quad (46)$$

Using (42)(44)(45)(46), we have

$$\text{dist}(\mathbf{u}, \mathbf{w}^t) \leq 0.8105 \cdot \text{dist}(\mathbf{u}, \mathbf{u}^{t-1}) + 0.0578\zeta'. \quad (47)$$

Since $k = C_5 s$ and \mathcal{F}_t is the set of indices with the largest k absolute values in \mathbf{w}^t , Lemma A.2 generates

$$\begin{aligned} |\mathbf{u}^\top \mathbf{w}_{\mathcal{F}_t}^t| &\geq |\mathbf{u}^\top \mathbf{w}^t| - C_5^{-1/2} \min \left\{ \sqrt{1 - |\mathbf{u}^\top \mathbf{w}^t|^2}, (1 + C_5^{-1/2}) (1 - |\mathbf{u}^\top \mathbf{w}^t|^2) \right\} \\ &\geq |\mathbf{u}^\top \mathbf{w}^t| - C_5^{-1/2} (1 + C_5^{-1/2}) (1 - |\mathbf{u}^\top \mathbf{w}^t|^2), \end{aligned}$$

which implies that

$$1 - |\mathbf{u}^\top \mathbf{w}_{\mathcal{F}_t}^t| \leq 1 - |\mathbf{u}^\top \mathbf{w}^t| + C_5^{-1/2} (1 + C_5^{-1/2}) (1 - |\mathbf{u}^\top \mathbf{w}^t|^2) \leq D_1^2 (1 - |\mathbf{u}^\top \mathbf{w}^t|),$$

where $D_1 := \sqrt{1 + 2C_5^{-1/2} (1 + C_5^{-1/2})}$. Then, since $\mathbf{u}^t = \mathbf{w}_{\mathcal{F}_t}^t / \|\mathbf{w}_{\mathcal{F}_t}^t\|_2$, we have

$$\begin{aligned} \text{dist}(\mathbf{u}, \mathbf{u}^t) &= \sqrt{2 - 2|\mathbf{u}^\top \mathbf{u}^t|} = \sqrt{2 - 2|\mathbf{u}^\top \mathbf{w}_{\mathcal{F}_t}^t| / \|\mathbf{w}_{\mathcal{F}_t}^t\|_2} \\ &\leq \sqrt{2 - 2|\mathbf{u}^\top \mathbf{w}_{\mathcal{F}_t}^t|} \leq D_1 \cdot \sqrt{2(1 - |\mathbf{u}^\top \mathbf{w}^t|)} \\ &= D_1 \cdot \text{dist}(\mathbf{u}, \mathbf{w}^t) \\ &\leq 0.8105 D_1 \cdot \text{dist}(\mathbf{u}, \mathbf{u}^{t-1}) + 0.0578 D_1 \zeta' \end{aligned} \quad (48)$$

where in the last second inequality we use (47). Since $\text{dist}(\mathbf{u}, \mathbf{u}^{t-1}) \leq 1$ and $\zeta' < 1$, the above inequality also implies that $\text{dist}(\mathbf{u}, \mathbf{u}^t) \leq 1$ with suitable constant C_5 (constant D_1). Therefore, we

complete the induction, which proves that $\text{dist}(\mathbf{u}, \mathbf{u}^t) \leq 1$ for all t . As a result, the above inequality holds for all t , which leads to

$$\begin{aligned} \text{dist}(\mathbf{u}, \mathbf{u}^t) &\leq \eta \cdot \text{dist}(\mathbf{u}, \mathbf{u}^{t-1}) + D_2 \zeta' \\ &\leq \eta^2 \cdot \text{dist}(\mathbf{u}, \mathbf{u}^{t-2}) + \eta D_2 \zeta' + D_2 \zeta' \\ &\leq \dots \\ &\leq \eta^t \cdot \text{dist}(\mathbf{u}, \mathbf{u}^0) + h \zeta', \end{aligned}$$

where $\eta := 0.8105D_1$, $D_2 := 0.0578D_1$ and $h := \frac{D_2}{1-\eta}$. This inequality is just (11). \square

Remark A.10. From (39)(48), a sufficient condition for constants C_4, C_5 in Theorem 3.1 is

$$\begin{aligned} C_4 &\geq \max \left\{ 16\sqrt{3}, 200\sqrt{1 + (1 + 2C_5)(1 + \log_2(9e))} \right\}, \\ (0.8105 + 0.0578) \sqrt{1 + 2C_5^{-1/2}(1 + C_5^{-1/2})} &< 1. \end{aligned}$$

It can be simplified as

$$\begin{aligned} C_4 &\geq 200\sqrt{(1 + 2\log_2(9e))C_5 + 2 + 2\log_2(9e)}, \\ C_5 &\geq 49.047. \end{aligned}$$

B ADDITIONAL EXPERIMENTAL RESULTS

B.1 COMPUTATIONAL EFFICIENCY AND STATISTICAL PERFORMANCE

Figure 7 provides complementary analysis of computational efficiency and support recovery, extending results in Figures 4 and 5. Runtime results (left panel) show that our TPM matches the efficiency of diagonal thresholding (DT), with only modest increase in cost at small β due to extra iterations with weaker initialization; in contrast, covariance thresholding (CT) and spectral projection (SP) are substantially more expensive across all signal strengths. Dimension-scaling results (right panel) demonstrate that our TPM achieves perfect support recovery (success rate = 1) across all tested dimensions, while competing methods underperform; in particular, DT maintains a low success rate throughout. This superior statistical performance incurs minimal computational overhead, as TPM’s runtime scales comparably to the efficient DT baseline (see Figure 5).

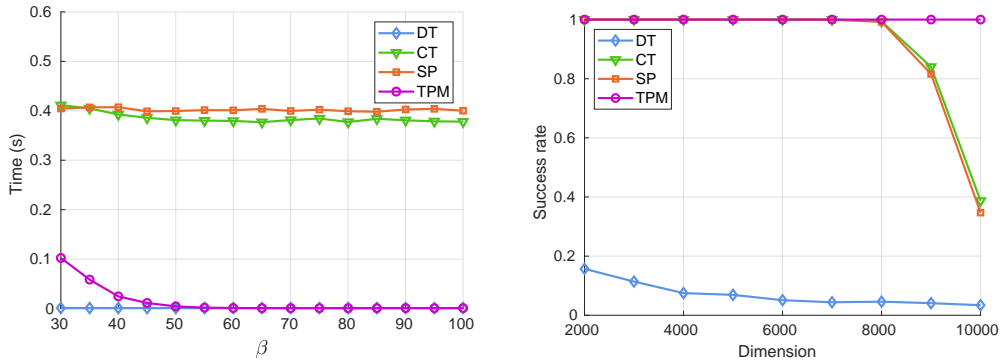


Figure 7: Left—Runtime versus signal strength β . Right—Support-recovery success rate versus dimension d . Experimental settings match those in Figures 4 and 5.

B.2 PERFORMANCE TRADE-OFFS BETWEEN COLUMN THRESHOLDING AND ITS VARIANT

We evaluate the performance trade-off between column thresholding (Algorithm 1) and its computationally efficient variant (Algorithm 2). The two methods differ only in their estimation procedures: column thresholding applies eigenvalue decomposition to selected submatrices for higher accuracy, whereas the variant adopts a simpler normalization-based estimator for speed. Both use the same thresholding rule for support identification, yielding identical support recovery. We assess this trade-off along two axes: problem dimension and sparsity level.

Table 1 shows performance across increasing dimensions. Column thresholding consistently achieves substantially lower estimation error than the variant, but requires more computation time. These findings suggest that practitioners should choose column thresholding when estimation accuracy is paramount and the variant when computational budget is tight. In addition, the runtime of column thresholding grows almost linearly with dimension, reflecting favorable scaling relative to CT and SP that require $O(d^3)$ operations.

Table 1: Performance comparison of column thresholding (Algorithm 1) and its variant (Algorithm 2) across increasing dimensions. Column thresholding achieves lower estimation error at higher computational cost, while the variant offers faster runtimes with reduced estimation accuracy; support-recovery rates are identical due to the shared thresholding strategy. Results averaged over 500 trials with $s = 25$ and $\beta = 150$.

Algorithm	Metric	$d = 1000$	$d = 3000$	$d = 5000$	$d = 7000$	$d = 9000$
Algorithm 1	Estimation error	0.0555	0.0846	0.0950	0.0993	0.1152
	Success rate	0.8660	0.7080	0.6480	0.6300	0.5420
	Runtime (s)	1.3×10^{-3}	1.4×10^{-3}	1.5×10^{-3}	1.6×10^{-3}	1.7×10^{-3}
Algorithm 2	Estimation error	0.1932	0.2087	0.2165	0.2205	0.2305
	Success rate	0.8660	0.7080	0.6480	0.6300	0.5420
	Runtime (s)	1.9×10^{-4}	2.8×10^{-4}	3.8×10^{-4}	5.5×10^{-4}	7.2×10^{-4}

Table 2 examines performance under varying sparsity. Column thresholding preserves its accuracy advantage across all sparsity levels while incurring higher runtime. As sparsity increases, its runtime grows because the eigenvalue decompositions operate on larger $s \times s$ submatrices. By contrast, the variant’s runtime remains nearly constant across sparsity levels, being driven primarily by the ambient dimension rather than sparsity.

Table 2: Performance comparison of column thresholding (Algorithm 1) and its variant (Algorithm 2) under varying sparsity levels. Column thresholding maintains superior estimation accuracy at higher computational cost, while both methods achieve identical support recovery. Results averaged over 500 trials with $d = 5000$ and $\beta = 150$.

Algorithm	Metric	$s = 10$	$s = 15$	$s = 20$	$s = 25$	$s = 30$
Algorithm 1	Estimation error	0.0198	0.0243	0.0307	0.1030	0.2549
	Success rate	1	1	0.9900	0.6000	0.0620
	Runtime (s)	1.0×10^{-3}	1.1×10^{-3}	1.1×10^{-3}	1.5×10^{-3}	1.5×10^{-3}
Algorithm 2	Estimation error	0.0738	0.1097	0.1456	0.2237	0.3619
	Success rate	1	1	0.9900	0.6000	0.0620
	Runtime (s)	4.1×10^{-4}	4.0×10^{-4}	4.0×10^{-4}	4.1×10^{-4}	4.1×10^{-4}

B.3 EMPIRICAL SIGNAL STRENGTH REQUIREMENTS UNDER UNIFORM AMPLITUDES

We examine the empirical signal strength requirement of the TPM initialized by column thresholding (Algorithm 3) in the *uniform-amplitude* setting. Our theoretical results show that, under the non-uniform ℓ_∞ condition $\|\mathbf{u}\|_\infty = \Omega(1)$, a signal strength of order $\beta = \Omega(\sqrt{s \log d})$ suffices. By contrast, the experiments here indicate that with uniform amplitudes, where $\|\mathbf{u}\|_\infty = 1/\sqrt{s}$, the algorithm empirically requires a stronger signal of order $\beta = \Omega(s^{0.6} \sqrt{\log d})$.

Figure 8 reports the estimation error as a function of the scaled signal strength $\beta/\sqrt{s \log d}$. In panel (a), when varying the dimension d , the curves collapse under this scaling. In panel (b), however, when varying the sparsity level s , the curves stabilize at different phase-transition thresholds, with larger s requiring larger values of $\beta/\sqrt{s \log d}$.

Figure 9 shows the corresponding support recovery performance. The phase-transition curves exhibit the same dependence on s : larger sparsity levels demand larger $\beta/\sqrt{s \log d}$ at transition. This parallel behavior indicates that the $\sqrt{s \log d}$ scaling is not attained for either estimation error or support identification in the uniform-amplitude case.

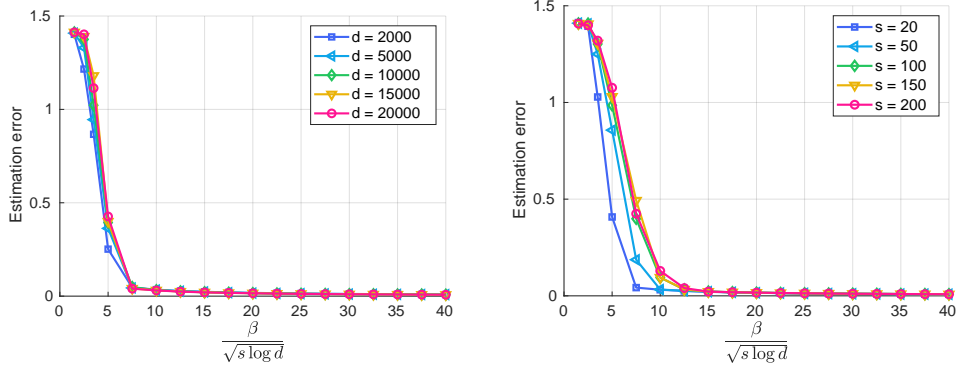


Figure 8: Estimation error versus scaled signal strength for TPM initialized by column thresholding (Algorithm 3), under varying dimensions (left) and sparsities (right). Experimental settings match those in Figure 2, except that the nonzero entries of the true spike \mathbf{u} have uniform amplitudes.

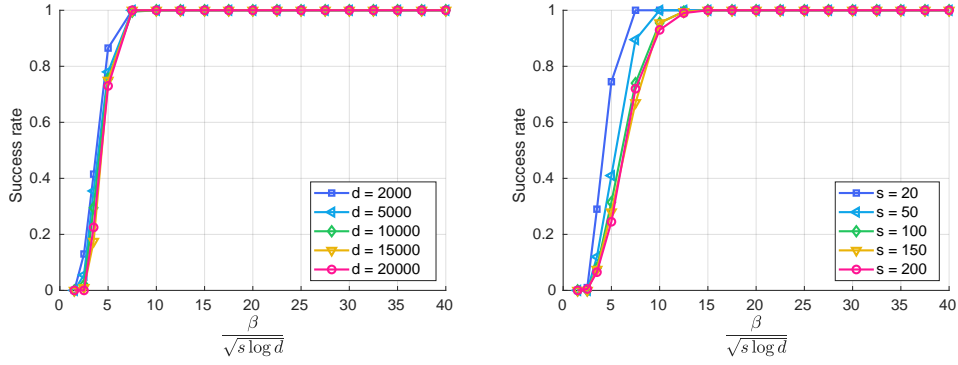


Figure 9: Success rate versus scaled signal strength for TPM initialized by column thresholding (Algorithm 3), under varying dimensions (left) and sparsities (right). Experimental settings match those in Figure 3, except that the nonzero entries of the true spike \mathbf{u} have uniform amplitudes.

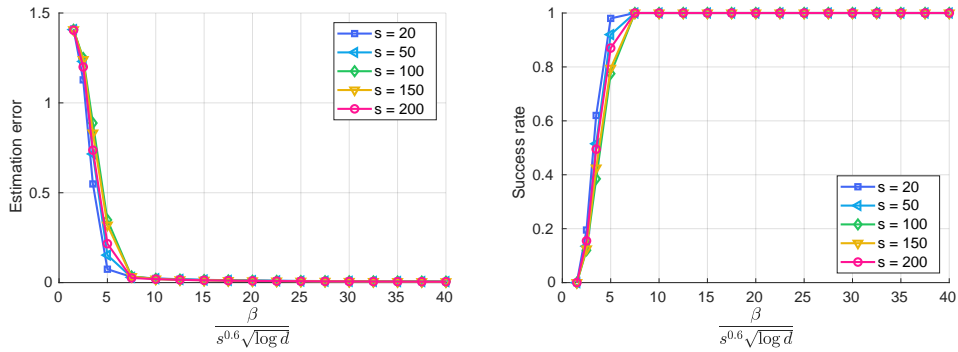


Figure 10: Estimation error (left) and success rate (right) versus scaled signal strength $\frac{\beta}{s^{0.6} \sqrt{\log d}}$ for the truncated power method initialized by column thresholding (Algorithm 3) under varying dimensions. Experimental settings match those in Figures 8 and 9.

In Figure 10, we instead scale the signal strength by $\beta/(s^{0.6} \sqrt{\log d})$. Under this scaling, the curves for different s align much more closely for both estimation error and support recovery. This collapse provides empirical evidence for an $s^{0.6} \sqrt{\log d}$ signal strength requirement for the column-thresholding-initialized truncated power method when the spike has uniform amplitudes.

Taken together, these experiments suggest that the ℓ_∞ condition in our analysis is not merely technical, but crucial for achieving the $\sqrt{s \log d}$ signal strength rate. In the uniform-amplitude setting, the algorithm does not empirically attain the $\sqrt{s \log d}$ scaling, in line with our theory. Nonetheless, the method remains practically attractive even when the ℓ_∞ condition is violated.

B.4 GROWING- s EXPERIMENT AND PHASE-TRANSITION BEHAVIOR

In this subsection we empirically investigate how the performance of the column-thresholding step evolves as the sparsity level s grows with the dimension d . We adopt the same phase-diagram parametrization as in Figure 1:

$$s = \tilde{\Theta}(d^\phi), \quad \frac{s}{\beta} = \tilde{\Theta}(d^\psi).$$

For any fixed ψ , increasing ϕ moves the problem from the “Impossible” region, through the “Hard” region, and eventually into the “Easy” region (see Figure 1).

In this experiment we fix $\psi = 0.2$, for which the phase diagram predicts a transition from “Impossible” to “Hard” at $\phi = 0.4$ and from “Hard” to “Easy” at $\phi = 0.7$. To probe the growing- s behavior along this slice, we consider three representative sparsity scalings

$$s = d^{0.4}, \quad s = d^{0.5}, \quad s = d^{0.6},$$

corresponding to $\phi = 0.4, 0.5, 0.6$, respectively, and set $\beta = 10sd^{-0.2}$ so that $s/\beta \asymp d^{0.2}$ in all cases. Figure 11 reports the empirical success probability over 500 Monte Carlo trials as a function of d for these three values of ϕ .

The results in Figure 11 are consistent with the theoretical phase diagram and clearly illustrate the growing- s transition. When $s = d^{0.4}$ (i.e., $\phi = 0.4$, at the boundary between the “Impossible” and “Hard” regions), the success rate remains close to zero across the range of dimensions considered, indicating that the algorithm almost never identifies the true support. When $s = d^{0.6}$ ($\phi = 0.6$, in the “Hard” region), the success probability rises to 1 (exact support recovery in every trial) over the range of d considered.

Overall, these experiments provide a quantitative illustration of the transition in algorithm performance as s grows with d , and show that the empirical behavior of the support-recovery step closely matches the theoretically predicted thresholds in (ϕ, ψ) -space.

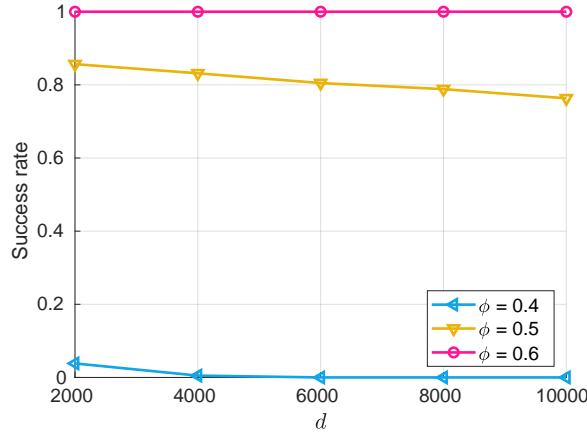


Figure 11: Empirical success rate of support recovery by column thresholding as a function of the dimension d for three sparsity scalings: $s = d^{0.4}$, $s = d^{0.5}$, and $s = d^{0.6}$, with $\psi = 0.2$ fixed. The success rate is computed over 500 Monte Carlo trials.