

Average-case integrality gap for non-negative principal component analysis

Afonso Bandeira, Dmitriy Kunisky, Alexander Wein

MSML 2021

Principal Component Analysis (PCA)

Principal Component Analysis (PCA)

Classical computational problem, especially in statistics:

$$\begin{aligned} &\text{maximize} && \mathbf{x}^\top \mathbf{W} \mathbf{x} \\ &\text{subject to} && \mathbf{x} \in \mathbb{S}^{n-1} \end{aligned}$$

Optimizer is the **leading eigenvector** or “**principal component**” of \mathbf{W} (e.g., \mathbf{W} a sample covariance matrix).

Principal Component Analysis (PCA)

Classical computational problem, especially in statistics:

$$\begin{aligned} & \text{maximize} && \mathbf{x}^\top \mathbf{W} \mathbf{x} \\ & \text{subject to} && \mathbf{x} \in \mathbb{S}^{n-1} \end{aligned}$$

Optimizer is the **leading eigenvector** or “**principal component**” of \mathbf{W} (e.g., \mathbf{W} a sample covariance matrix).

Elaborated version, allowing **prior information**:

$$\begin{aligned} & \text{maximize} && \mathbf{x}^\top \mathbf{W} \mathbf{x} \\ & \text{subject to} && \mathbf{x} \in \mathcal{X} \subset \mathbb{S}^{n-1} \end{aligned}$$

This paper: **non-negative** PCA, with $\mathcal{X} = \mathbb{R}_+^n$.

Spiked Matrix Model

Well-studied toy model of PCA:

- Under \mathbb{Q} , observe $W \sim \text{GOE}(n)$ (i.i.d. Gaussian entries, symmetrized).
- Under \mathbb{P} , observe $W = W^{(0)} + \lambda \mathbf{x} \mathbf{x}^\top$ for $W^{(0)} \sim \text{GOE}(n)$, \mathbf{x} drawn from some signal distribution over $\mathcal{X} \subset \mathbb{S}^{n-1}$.

Spiked Matrix Model

Well-studied toy model of PCA:

- Under \mathbb{Q} , observe $W \sim \text{GOE}(n)$ (i.i.d. Gaussian entries, symmetrized).
- Under \mathbb{P} , observe $W = W^{(0)} + \lambda \mathbf{x} \mathbf{x}^\top$ for $W^{(0)} \sim \text{GOE}(n)$, \mathbf{x} drawn from some signal distribution over $\mathcal{X} \subset \mathbb{S}^{n-1}$.

Interesting questions:

- Can we **distinguish** \mathbb{Q} from \mathbb{P} ?
- Can we **estimate** the signal \mathbf{x} under \mathbb{P} ?
- Can we **optimize** the likelihood $\mathbf{x}^\top W \mathbf{x}$ over $\mathbf{x} \in \mathcal{X}$?

Spiked Matrix Model

Well-studied toy model of PCA:

- Under \mathbb{Q} , observe $W \sim \text{GOE}(n)$ (i.i.d. Gaussian entries, symmetrized).
- Under \mathbb{P} , observe $W = W^{(0)} + \lambda \mathbf{x} \mathbf{x}^\top$ for $W^{(0)} \sim \text{GOE}(n)$, \mathbf{x} drawn from some signal distribution over $\mathcal{X} \subset \mathbb{S}^{n-1}$.

Interesting questions:

- Can we **distinguish** \mathbb{Q} from \mathbb{P} ?
- Can we **estimate** the signal \mathbf{x} under \mathbb{P} ?
- Can we **optimize** the likelihood $\mathbf{x}^\top W \mathbf{x}$ over $\mathbf{x} \in \mathcal{X}$?

This paper: optimization for $\mathbb{Q} = \text{GOE}(n)$, $\mathcal{X} = \mathbb{R}_+^n \cap \mathbb{S}^{n-1}$.

Algorithmic Approaches [Montanari, Richard 2015]

Algorithmic Approaches [Montanari, Richard 2015]

Scale to have $\lambda_{\max}(\mathbf{W}) \approx 2$. Then, using **approximate message-passing** algorithm have

$$\lambda^+(\mathbf{W}) := \left\{ \begin{array}{l} \text{maximize } \mathbf{x}^\top \mathbf{W} \mathbf{x} \\ \text{subject to } \|\mathbf{x}\| = 1 \\ \mathbf{x}_i \geq 0 \end{array} \right\} \approx \sqrt{2}$$

Algorithmic Approaches [Montanari, Richard 2015]

Scale to have $\lambda_{\max}(W) \approx 2$. Then, using **approximate message-passing** algorithm have

$$\lambda^+(W) := \left\{ \begin{array}{l} \text{maximize } \mathbf{x}^\top W \mathbf{x} \\ \text{subject to } \|\mathbf{x}\| = 1 \\ \phantom{\text{subject to }} x_i \geq 0 \end{array} \right\} \approx \sqrt{2}$$

Alternative: **semidefinite program** with $X = \mathbf{x}\mathbf{x}^\top$ relaxed to

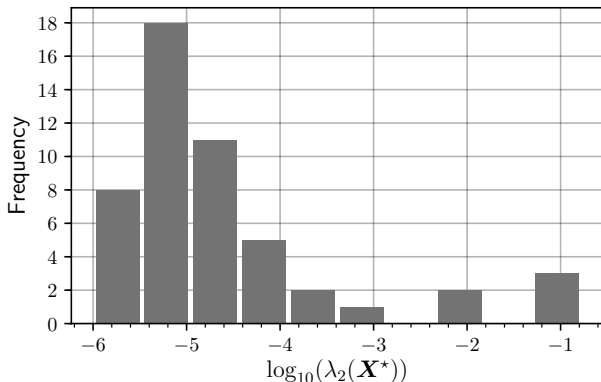
$$\text{SDP}(W) := \left\{ \begin{array}{l} \text{maximize } \langle X, W \rangle \\ \text{subject to } X \succeq \mathbf{0} \\ \phantom{\text{subject to }} \text{Tr}(X) = 1 \\ \phantom{\text{subject to }} X_{ij} \geq 0 \end{array} \right\} \geq \lambda^+(W)$$

Seems to recover under $\mathbb{P} \rightsquigarrow$ expect $\text{SDP}(W) \approx \sqrt{2}$ under \mathbb{Q} .

Small SDP Experiments ($n = 150$)

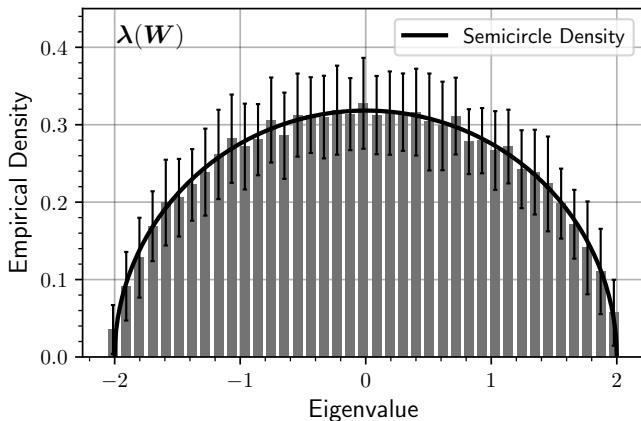
$$\text{SDP}(W) = \left\{ \begin{array}{l} \text{maximize} \quad \langle X, W \rangle \\ \text{subject to} \quad X \succeq \mathbf{0}, \text{Tr}(X) = 1, X_{ij} \geq 0 \end{array} \right\}$$

Optimizer X^* has numerical rank ≈ 1 :



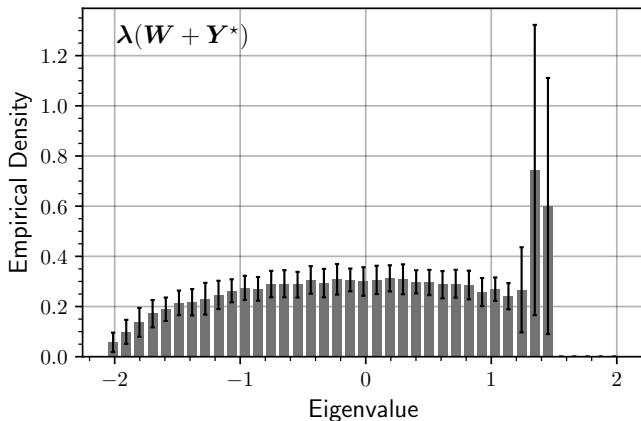
Small SDP Experiments ($n = 150$)

$$\text{SDP}(\mathbf{W}) := \left\{ \begin{array}{l} \text{minimize} \quad \lambda_{\max}(\mathbf{W} + \mathbf{Y}) \\ \text{subject to} \quad Y_{ij} \geq 0 \end{array} \right\}$$



Small SDP Experiments ($n = 150$)

$$\text{SDP}(W) := \left\{ \begin{array}{ll} \text{minimize} & \lambda_{\max}(W + Y) \\ \text{subject to} & Y_{ij} \geq 0 \end{array} \right\} \stackrel{?}{\approx} \sqrt{2}$$



An Unpleasant Surprise

Theorem 1: With high probability $\text{SDP}(\mathbf{W}) \approx 2 \approx \lambda_{\max}(\mathbf{W})$.

An Unpleasant Surprise

Theorem 1: With high probability $\text{SDP}(\mathbf{W}) \approx 2 \approx \lambda_{\max}(\mathbf{W})$.

Proof Sketch: Need \mathbf{X} feasible with $\langle \mathbf{X}, \mathbf{W} \rangle \approx \lambda_{\max}(\mathbf{W})$.

An Unpleasant Surprise

Theorem 1: With high probability $\text{SDP}(\mathbf{W}) \approx 2 \approx \lambda_{\max}(\mathbf{W})$.

Proof Sketch: Need \mathbf{X} feasible with $\langle \mathbf{X}, \mathbf{W} \rangle \approx \lambda_{\max}(\mathbf{W})$.

Let \mathbf{P} be the projector to top δn eigenspaces of \mathbf{W} .

$$\mathbf{X}^{(0)} := \frac{1}{\delta n} \mathbf{P}$$

almost feasible; only $X_{ij}^{(0)} < 0$ sometimes.

An Unpleasant Surprise

Theorem 1: With high probability $\text{SDP}(W) \approx 2 \approx \lambda_{\max}(W)$.

Proof Sketch: Need X feasible with $\langle X, W \rangle \approx \lambda_{\max}(W)$.

Let P be the projector to top δn eigenspaces of W .

$$X^{(0)} := \frac{1}{\delta n} P$$

almost feasible; only $X_{ij}^{(0)} < 0$ sometimes.

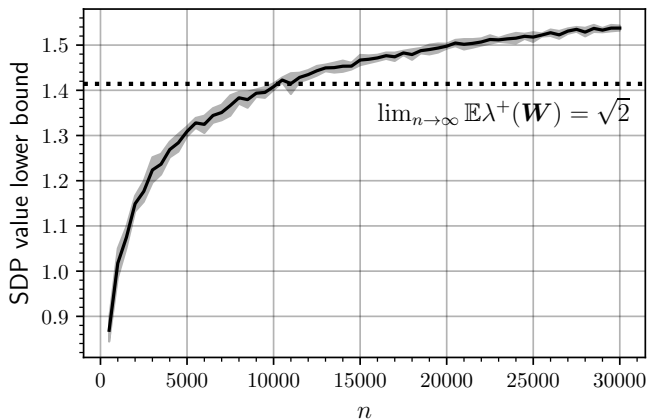
So, **nudge** towards feasible set:

$$X := \alpha \frac{1}{n} \mathbf{1}\mathbf{1}^\top + (1 - \alpha) X^{(0)}.$$

Analysis of off-diagonal entries \rightsquigarrow feasible for any $\alpha > 0$.

Larger Experiments ($\delta = 1/25$, α as small as possible)

Expect to need $n \gtrsim 10^4$ to see $\text{SDP}(\mathbf{W}) > \sqrt{2}$. Much bigger than scale of quickly solving SDP!



General Certification Algorithms

What's going on—is the SDP just not fancy enough and there are better similar methods?

General Certification Algorithms

What's going on—is the SDP just not fancy enough and there are better similar methods?

Natural follow-up question: does there exist a better efficient **certification algorithm**? That is, $c(W)$ with

$$\lambda^+(W) \leq c(W)$$

such that

- The bound holds **for all** $W \in \mathbb{R}_{\text{sym}}^{n \times n}$.
- $c(W) \leq 2 - \epsilon$ with high probability when $W \sim \text{GOE}(n)$, for some $\epsilon > 0$.

Low-Degree Lower Bound

A two-part argument that such an algorithm (running in subexponential time) **cannot exist**.

Low-Degree Lower Bound

A two-part argument that such an algorithm (running in subexponential time) **cannot exist**.

Part 1: Reduction. If such existed, it could be used to distinguish two models of **bottom** eigenspaces of W :

- $\mathcal{Y}_1, \dots, \mathcal{Y}_{(1-\delta)n} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_n)$.
- $\mathcal{Y}_1, \dots, \mathcal{Y}_{(1-\delta)n} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_n - \beta \mathbf{x} \mathbf{x}^\top)$ for \mathbf{x} close to positive orthant.

Low-Degree Lower Bound

A two-part argument that such an algorithm (running in subexponential time) **cannot exist**.

Part 1: Reduction. If such existed, it could be used to distinguish two models of **bottom** eigenspaces of W :

- $\mathcal{Y}_1, \dots, \mathcal{Y}_{(1-\delta)n} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_n)$.
- $\mathcal{Y}_1, \dots, \mathcal{Y}_{(1-\delta)n} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_n - \beta \mathbf{x} \mathbf{x}^\top)$ for \mathbf{x} close to positive orthant.

Part 2: Hardness. Low-degree polynomials of the observations **cannot** distinguish these two models.

\rightsquigarrow Assuming optimality of low-degree polynomial tests, no efficient certifier can beat $\text{SDP}(W) \approx \lambda_{\max}(W)$.

Takeaway Messages

Takeaway Messages

1. Small experiments can be misleading about the asymptotic behavior of semidefinite programs.

There are strong finite-size effects up to current “laptop scale” with off-the-shelf solvers.

Takeaway Messages

1. Small experiments can be misleading about the asymptotic behavior of semidefinite programs.

There are strong finite-size effects up to current “laptop scale” with off-the-shelf solvers.

2. Among certification algorithms (in particular convex relaxations), it is hard to beat the simple spectral bound $\lambda^+(\mathbf{W}) \leq \lambda_{\max}(\mathbf{W})$.

In this regard, non-negative PCA behaves like many other constrained PCA problems studied in previous literature.

Thank You!