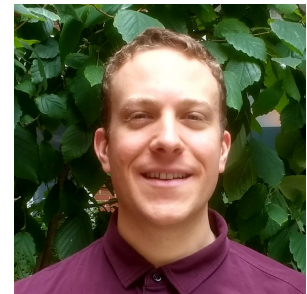


# Computational Barriers to Estimation from Low-Degree Polynomials

Tselil Schramm (Stanford)

with Alex Wein (NYU)



# High-dimensional statistics



Can **fast** algorithms use **noisy** data to give solid conclusions?

information - computation tradeoff

# Simple planted models

We (simply, but faithfully) model our data as **signal** + **noise**

e.g. the spiked Wigner matrix model

observe:  $\underbrace{M = \lambda \cdot \underbrace{uu^\top} + \underbrace{G}}_{\in \mathbb{R}^{d \times d}}$ , with i.i.d.  $G_{ij} = G_{ji} \sim N\left(0, \frac{1}{d}\right)$

goal(s):

- **detection**: Is the signal  $u$  present? (hypothesis testing with null hypothesis  $M \sim N\left(0, \frac{1}{d}\right)^{d \times d}$ )
- **estimation/recovery**: find  $\hat{u}$  with  $\|u - \hat{u}\|$  as small as possible

Several variations on a theme

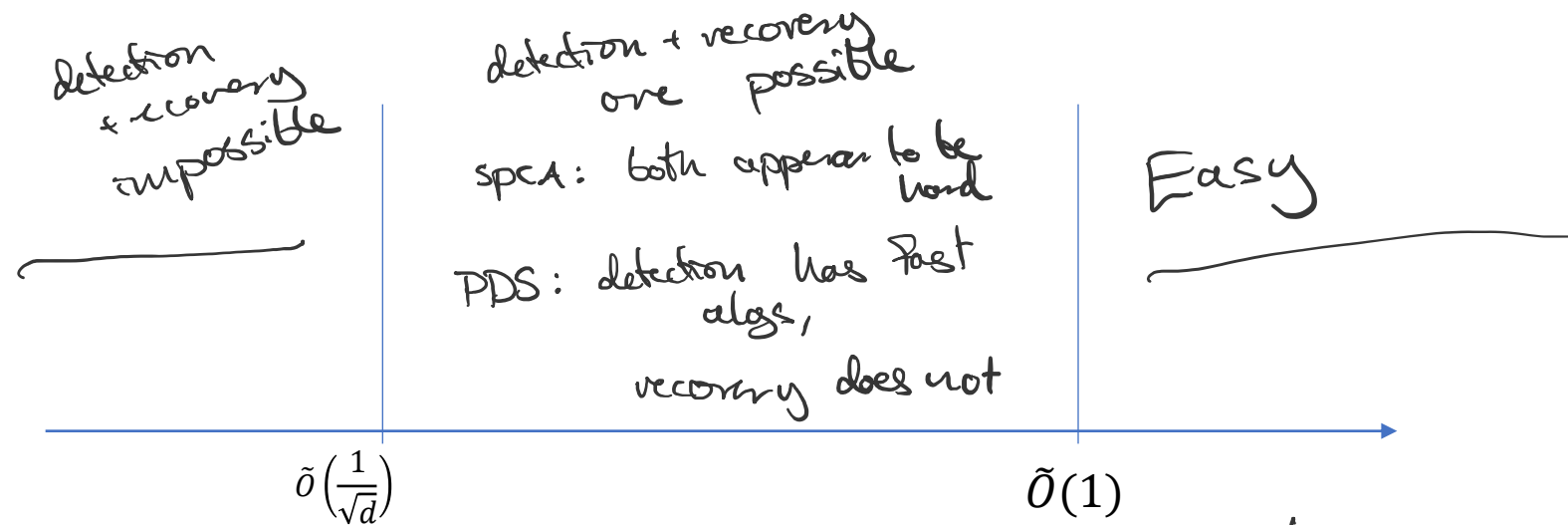
- Sparse PCA: i.i.d.  $u_i \sim \sqrt{\frac{1}{k}} \cdot \text{Ber}\left(\frac{k}{d}\right) \cdot \text{Rad}\left(\frac{1}{2}\right)$
- Planted Dense Submatrix: i.i.d.  $u_i \sim \sqrt{\frac{1}{k}} \cdot \text{Ber}\left(\frac{k}{d}\right)$

**These models are \*famous\*!**

[Johnstone-Lu'09, Baik-Ben-Arous-Peche'05, Brennan-Bresler'19, Chen-Xu'16, Ding-Kunisky-Wein-Bandeira'19, Deshpande-Montanari'14, Holtzman-Soffer-Vilenchik'20, Butucea-Ingster'13, Barbier-Macris-Rush'20, plus **dozens more**...]

# Information-computation gaps

When are detection and recovery possible (with fast algorithms)?



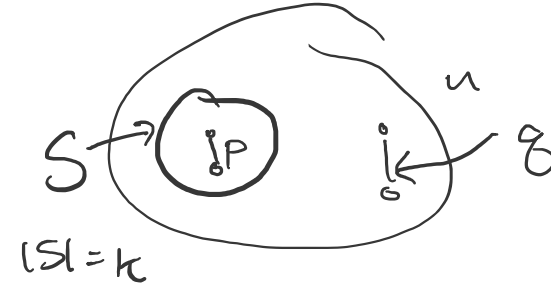
signal-to-noise parameter:  $\frac{\sqrt{d}}{k}$   
 (|signal entry| / |noise entry|)

$$M = \frac{1}{100} uu^T + G \in \mathbb{R}^{d \times d}, G_{ij} = G_{ji} \sim N\left(0, \frac{1}{d}\right)$$

Sparse PCA:  $u_i \sim \sqrt{\frac{1}{k}} \cdot \text{Ber}\left(\frac{k}{d}\right) \cdot \text{Rad}\left(\frac{1}{2}\right)$

Planted Dense Submatrix:  $u_i \sim \sqrt{\frac{1}{k}} \cdot \text{Ber}\left(\frac{k}{d}\right)$

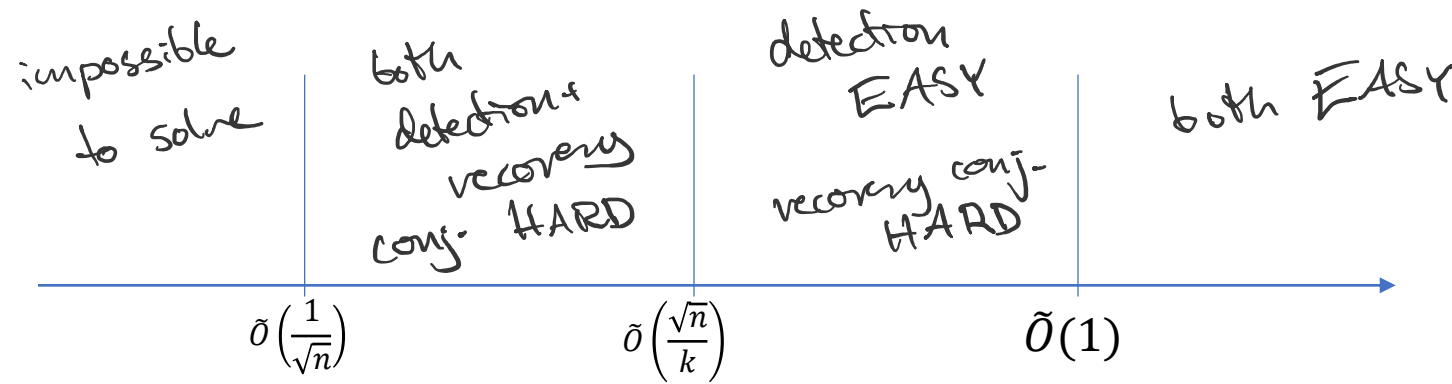
# Another simple model



## Planted dense subgraph:

Observe a graph  $G = (V, E)$  on  $n$  vertices with a hidden subgraph on  $S \subset V$ ,  $|S| = k$  where

$$\Pr[(i, j) \in E(G)] = \begin{cases} p & \text{if } i, j \in S \\ q & \text{otherwise} \end{cases}$$



signal-to-noise parameter:  $\frac{pk}{\sqrt{qn}}$

(avg degree in  $S$  / std. dev. of degree in  $G$ )

# Explaining intractability

Planted Dense  
Subgraph

$$p=1$$
$$q=\frac{1}{2}$$



$$p=1$$
$$q=\frac{1}{4}$$

Can we give rigorous evidence for computational barriers?

## 1. Reductions

HARD  
PROBLEM



OUR  
PROBLEM

[Ma-Wu'15], [Chen-Xu'16], [Brennan-Bresler'19], [Brennan-Bresler-Huleihel'19], etc...



## 2. Lower bounds against restricted models of computation

statistical query lower bounds, convex program (sum-of-squares) lower bounds, approximate message passing/belief propagation lower bounds, “energy barriers”, low-degree polynomials

# Low-degree polynomials

$f$  : Graph  
adj  
matrix  $\rightarrow$

YES/1 if graph  
has clique  
of size  $> k$   
NO/0 o/w

Classical complexity theory:  $f: \{0,1\}^n \rightarrow \{0,1\}$ .

What is the degree of  $f$  as a polynomial (over  $\mathbb{R}, \mathbb{F}_2$ )?

How well can a degree- $d$  polynomial approximate  $f$ ?

degree- $d$  function  
in time  $n^d$

Complexity of statistics:

How well can a degree- $d$  polynomial detect/estimate?

$f$ : data  $\rightarrow$  conclusion

# Low-degree polynomials in high-dimensional statistics

Why low-degree polynomials?

$$\begin{aligned} \mu &, \text{ spectral gap} & \mu &= \lambda u u^T + G \\ & & \lambda &> (1+\varepsilon) \lambda_{\max}(G) \\ u &\propto \mu^{\log u / \varepsilon} & g &\leftarrow \text{random vector} \end{aligned}$$

- Many algorithms are (approximately) low-degree  
e.g. many spectral algorithms, message passing, [folklore]  $O(\log u)$   
“reasonable” statistical query algorithms [Brennan-Bresler-Hopkins-Li-S’20],  
sum-of-squares semidefinite programs ? [Barak-Hopkins-Kelner-Kothari-Moitra-Potechin’16]

**Conclusion:** if we rule out low-degree polynomials, it is unlikely that other go-to algorithms work

- Accurately predict current computational thresholds for **detection!**  
degree  $\omega(\log n)$  required above known computational threshold



# Computational barriers to *detection* from low-degree polynomials

$$\begin{aligned}\text{planted:} & \quad \lambda uu^T + G \\ \text{null:} & \quad G \leftarrow\end{aligned}$$

Predictions consistent with the detection threshold for many problems:

Planted Clique [Barak-Hopkins-Kelner-Kothari-Moitra-Potechin'16], sparse PCA [Ding-Bandeira-Kunisky-Wein'19], tensor PCA [Bandeira-Kunisky-Wein'19], community detection in block models [Hopkins-Steurer'17], random CSPs...

Convenient closed form solution for detection when testing against null measure with product structure (e.g. sparse PCA, null =  $N(0, I_d)$ )

**But** for some problems, we observe a **detection-recovery gap**

planted dense submatrix, planted dense subgraph, overcomplete tensor decomposition, graph matching in mildly correlated random graphs,...

# Our results: computational barriers to **estimation** from low-degree polynomials $f(\mathbf{A}) = \hat{\mathbf{u}}$

$$\mu = \lambda \frac{\mathbf{u}\mathbf{u}^\top}{\tau} + \mathbf{G} \quad f: \mathbb{R}^{d \times d} \rightarrow \mathbb{R}^d$$

**Theorem (Planted Submatrix):** ↙

If  $\sqrt{d}/k \leq d^{-\delta}$ , no degree- $O(d^{\varepsilon(\delta)})$  polynomial outperforms the trivial estimator  
(which guesses every coordinate is equally likely to be in the support)

**Theorem (Planted Dense Subgraph):** ↙

resolve open problems from

[Kolar-Balakrishnan-Rinaldo-Singh'11], [Ma-Wu'15], [Chen-Xu'16]

↗ If  $pk/\sqrt{qn} \leq n^{-\delta}$ , no degree- $O(n^{\varepsilon(\delta)})$  polynomial outperforms the trivial estimator  
(which guesses every coordinate is equally likely to be in the subgraph)

Both derive from **more general result** characterizing the minimum mean squared error of best degree- $\ell$  estimator in **additive Gaussian** or **binary observation** model.

# Outline

- Background: degree lower bounds for detection
- Degree lower bounds for estimation
  - Framework
  - Results for additive gaussian and binary observation models

# Degree lower bounds for detection (background)

Setting: “null” model  $Q$  (e.g.  $N(0, I_n)$ ) and “planted” model  $P$  (e.g.  $\int N(\mu, I_n) d\mu$ )

Goal: find  $f \in \mathbb{R}[x]^{\leq \ell}$  with  $|\mathbb{E}_P f - \mathbb{E}_Q f| \gg \sqrt{\max(\mathbb{V}_Q f, \mathbb{V}_P f)}$

Convex program:

Strategy: upper bound Lemma: the optimizer is  $f^* \propto \left(\frac{dP}{dQ} - 1\right)^{\leq \ell}$ , with objective value  $\left\|\left(\frac{dP}{dQ} - 1\right)^{\leq \ell}\right\|_Q$ .

$$\max_{f \in \mathbb{R}[x]^{\leq \ell}} \mathbb{E}_P f$$

$$\text{s.t. } \mathbb{E}_Q f = 0,$$

$$\mathbb{V}_Q f \leq 1,$$

~~$$\mathbb{V}_P f \leq 1.$$~~

Conclusion: if  $\left\|\left(\frac{dP}{dQ} - 1\right)^{\leq \ell}\right\|_Q = o(1)$ , then there is no degree- $\ell$  polynomial which meets our criteria  
(no tests with difference in mean larger than the standard deviation)

# Degree lower bounds for detection

Strategy: compute  $\left\| \left( \frac{dP}{dQ} - 1 \right)^{\leq \ell} \right\|_Q$ , rule out degree- $\ell$  polynomials if  $o(1)$ .

This framework was used to give consistent detection lower bounds for many problems.

Planted clique [Barak-Hopkins-Kelner-Kothari-Moitra-Potechin'16], sparse PCA [Ding-Bandeira-Kunisky-Wein'19], tensor PCA [Bandeira-Kunisky-Wein'19], community detection in block models [Hopkins-Steurer'17], random CSPs...

Successful when  $\mathbb{R}[x]^{\leq \ell}$  has a nice orthogonal (w.r.t  $\langle \cdot, \cdot \rangle_Q$ ) basis (e.g. product measures).

Also, used to give new algorithms! (Evaluate  $\left( \frac{dP}{dQ} - 1 \right)^{\leq \ell}$  and threshold)

Overlapping block models [Hopkins-Steurer'17], graph matching [Barak-Chou-Lei-S-Sheng'19]

# Outline

- Background: degree lower bounds for detection ✓
- Degree lower bounds for estimation
  - Framework
  - Results for additive gaussian and binary observation models

# Degree lower bounds for estimation

$$g(\overset{x}{\underset{\mu}{i}}) = g(x) \in \mathbb{R}^n$$

Setting: “planted” model  $P$  (e.g.  $\int N(\underline{\mu}, I_n) d\mu$ )

Goal: find  $g \in (\mathbb{R}[x]^{\leq \ell})^{\otimes n}$  with  $\mathbb{E}_{(\mu, x) \sim P} \|g(x) - \underline{\mu}\|^2 = o(\mathbb{E}_{\mu \sim P} \|\mu\|^2)$

degree- $\ell$  MMSE:

$$\min_{g \in (\mathbb{R}[x]^{\leq \ell})^{\otimes n}} \mathbb{E}_P \|g(x) - \mu\|^2$$

scaling  
↓

$$\begin{aligned} & \max_{h \in (\mathbb{R}[x]^{\leq \ell})^{\otimes n}} \mathbb{E}_P \langle h(x), \mu \rangle \\ & \text{s.t. } \mathbb{E}_P \|h(x)\|^2 \leq 1 \end{aligned}$$

symmetry  
→

$$\begin{aligned} \text{OPT}_i &= \max_{h_i} \langle h_i, \mu_i^{\leq \ell} \rangle_P \\ & \text{s.t. } \|h_i\|_P \leq 1 \end{aligned}$$

$$\text{Lemma: } \text{OPT}_i = \|(\mu_i)^{\leq \ell}\|_P.$$

Conclusion: if for all  $i$ ,  $\text{OPT}_i \leq \sqrt{\alpha \cdot \mathbb{E}_P \mu_i^2}$ , error  $\geq (1 - \alpha) \cdot \mathbb{E}_P \|\mu\|^2$ .

# A familiar story?

So... why not just compute  $\|(\mu_i)^{\leq \ell}\|_P$  for each  $i$ ? (like for detection)

For planted distributions of interest,  $\mathbb{R}[x]^{\leq \ell}$  has a nasty orthogonal basis w.r.t  $\langle \cdot, \cdot \rangle_P$ !

Compare to hypothesis testing of  $Q = N(0, I)$  vs.  $P = \mathbb{E}_\mu N(\mu, I)$ .

$$\begin{array}{ccccc} \text{OPT}_i = \max_{h_i} \langle h_i, \mu_i^{\leq \ell} \rangle_P & \xrightarrow{\text{change-of-basis}} & \text{OPT}_i = \max_{f_i} \left\langle f_i, \left( \frac{dP}{dQ} \mu_i \right)^{\leq \ell} \right\rangle_Q & \xrightarrow{\text{closed form}} & \text{OPT}_i = \left\| \underbrace{(B^{-1})^\top}_{\text{change-of-basis}} \left( \frac{dP}{dQ} \mu_i \right)^{\leq \ell} \right\|_Q \\ \text{s.t. } \|h_i\|_P \leq 1 & & \text{s.t. } \|B f_i\|_Q \leq 1 & & \end{array}$$



# A solution for additive Gaussian models

Suppose  $P = \mathbb{E}_\mu N(\mu, I)$ .

For  $x = \underbrace{\mu + G}_{\sim P}$ , think of first sampling the **signal**  $\mu$ , then the noise  $G \sim Q = N(0, I)$ .

$$\begin{aligned} \forall i \in [n], \text{OPT}_i &= \max_{h_i \in \mathbb{R}[x]^{\leq \ell}} \mathbb{E}_{x \sim P} h_i(x) \mu_i \\ \text{s. t. } \underbrace{\mathbb{E}_{x \sim P} h_i(x)^2}_{\leq 1} &\quad \xrightarrow{\text{relaxation, by Jensen}} \quad \text{s. t. } \underbrace{\mathbb{E}_G}_{\leq 1} \left[ \underbrace{\mathbb{E}_\mu}_{\leq 1} h_i(\mu + G) \right]^2 \leq 1 \end{aligned}$$

intuitively: not too lossy when recovery is impossible

Let  $f_i(G) = \mathbb{E}_\mu h_i(\mu + G)$ ; in additive Gaussian models, we can write  $f_i = Ah_i$  for  $A$  **upper triangular**.

$\forall i \in [n], \text{rOPT}_i = \|(A^{-1})^T c_i\|_Q$  is tractable to compute, and we get a closed form.

# Estimation lower bounds in additive Gaussian + Binary observation models

Closed form for  $\text{OPT}_i$  for **additive Gaussian** planted models  
in which we observe  $x \sim P = \mathbb{E}_\mu N(\mu, I_n)$

Also for **Binary observation** models

in which  $\mu \sim D([0,1]^n)$ , we observe  $x_i \sim \text{Ber}(\mu_i) \forall i$  (e.g. planted dense subgraph)

Exact expression for degree- $\ell$  MMSE.

Special case: lower bounds for estimation in **planted submatrix** and **planted dense subgraph**.

# Conclusion

tl;dr : we extend methods for **lower bounding the polynomial degree** of hypothesis tests to lower bounding the polynomial degree of **estimators**.

We give the first(ish) **rigorous evidence** for hardness (against a restricted class of algorithms) of **planted submatrix** and **planted dense subgraph** below known algorithmic thresholds.

## Open:

*ish*

- Optimal degree-vs.-estimation tradeoff?
- Estimation lower bounds for other models with detection-recovery gaps?  
(favorite problem: overcomplete tensor decomposition)
- Extending consequences to other models  
(degree lower bounds imply SoS lower bounds?)

Thank you!