

The Low Degree Method

Main Question: Can we decide if a sample was drawn from P or Q ?

A test is any function

$$f: S \rightarrow \{p, a\}$$

sample
space

and its type I error is

$$\alpha(f) = \mathbb{P}[f(x) = p \mid x \sim Q]$$

and its type II error is

$$\beta(f) = \mathbb{P}[f(x) = a \mid x \sim P]$$

What is the optimal tradeoff between α and β ?

def: The likelihood ratio test with threshold κ is

$$L_{\kappa}(x) = \begin{cases} p & \text{if } \frac{P(x)}{Q(x)} \geq \kappa \\ a & \text{else} \end{cases}$$

Proposition [Neyman, Pearson] For any $\kappa \geq 0$, among all tests f with

$$\alpha(f) \leq \alpha(L_{\kappa})$$

L_{κ} maximizes the power $1 - \beta(f)$

Proof: Let $R_f = \{x \mid f(x) = P\}$. Then

$$\begin{aligned} 1 - \beta(f) &= P[f(x) = P \mid x \sim P] \\ &= \int_{R_f} dP = \int_{R_f} L(x) d\phi \end{aligned}$$

Now by assumption

$$\alpha(f) = \phi[R_f] \leq \phi[L(x) > \tau] = \alpha(L_\tau)$$

Thus we can write the problem of maximizing the power as

$$\max \int_{R_f} L(x) d\phi$$

$$\text{s.t. } \phi[R_f] \leq \phi[L(x) > \tau]$$

Now let's look at the difference of powers

$$(1 - \beta(L_\tau)) - (1 - \beta(f))$$

$$= \int_{R_{L_\tau}} L(x) d\phi - \int_{R_f} L(x) d\phi$$

$$= \int_{R_{L_\tau}} L(x) d\phi - \int_{R_f} L(x) d\phi$$

$$\begin{array}{ccc} R_{L_\tau} \setminus R_f & \uparrow & R_f \setminus R_{L_\tau} \\ & \geq \tau & \leq \tau \end{array}$$

Thus the difference in power

$$\begin{aligned} &\geq \mathbb{N}(\phi[R_{L_n} \setminus R_f] - \phi[R_f \setminus R_{L_n}]) \\ &= \mathbb{N}(\underbrace{\phi[R_{L_n}] - \phi[R_f]}_{\geq 0}) \end{aligned}$$

□

However computing the likelihood can be hard, b/c it involves summing over all partitions, etc

Main Idea [Hopkins]: Consider the projection onto low degree polynomials

Preliminary ideas appeared in [Barak, Hopkins, Kelner, Kothari, Maitra, Potechin] for SOS lower bounds for planted clique.

Consider the following inner-product structure

$$\langle f, g \rangle \stackrel{\Delta}{=} \mathbb{E}_{x \sim \mathcal{P}} [f(x)g(x)]$$

In particular

$$\|f\|^2 \stackrel{\Delta}{=} \langle f, f \rangle$$

Now consider the following optimization problem:

$$\max_{x \sim P} \mathbb{E}[f(x)]$$

$$\text{s.t. } \mathbb{E}[f(x)^2] = \langle f, f \rangle = 1$$

claim: the optimal solution is

$$f^* = \frac{L(x)}{\|L\|}, \text{ normalized LRT}$$

Proof: We can rewrite the objective as

$$\mathbb{E}_{x \sim P}[f(x)] = \mathbb{E}_{x \sim P}[L(x)f(x)] = \langle L, f \rangle$$

And by C-S we have

$$\langle L, f \rangle \leq \|L\| \|f\| = \|L\|$$

$\underbrace{\|f\|}_{=1 \text{ from constraints}}$

with equality iff $f = cL$ for some constant c . \square

Intuitively we want f to be large on P while being bounded on \mathcal{Q}

Now let's define a low degree variant f . Let

$$\mathcal{U}_n^{\leq D} = \text{polynomials on } S \text{ of degree } \leq D$$

e.g. for planted clique, the variables correspond to edges

$$x_{ij} = \begin{cases} 1 & \text{if } (i,j) \in E \\ -1 & \text{else} \end{cases}$$

and we can count subgraphs, e.g.

$$\left(\frac{a+1}{2}\right) \left(\frac{b+1}{2}\right) \left(\frac{c+1}{2}\right) = \begin{cases} 1 & \text{if } a=b=c=1 \\ 0 & \text{if } a=-1 \text{ or } \\ & b=-1 \text{ or } \\ & c=-1 \end{cases}$$

and so the polynomial

$$\sum_{i,j,k} \left(\frac{x_{ij}+1}{2}\right) \left(\frac{x_{ik}+1}{2}\right) \left(\frac{x_{jk}+1}{2}\right)$$

counts the # of triangles

Now consider

$$\max_{x \sim P_n} \mathbb{E} [f(x)]$$

$$\text{s.t. } \mathbb{E} [f(x)^2] = 1$$

$$f(x) \in \mathcal{D}_n^{\leq D}$$

The optimal solution is

$$f^* = \frac{L^{\leq D}}{\|L^{\leq D}\|} \quad \text{low-degree projection of LRT}$$

For "nice" distributions

Conjecture [Hopkins] (informal) If for some $\epsilon > 0$ and $D \geq (\log n)^{1+\epsilon}$ we have

$$\|L_n^{\leq D}\| = o(1)$$

then there is no polynomial time algorithm for distinguishing P and Q with prob $1 - o(1)$

More generally:

$$\begin{array}{l} \text{degree } D \\ \text{polynomials} \end{array} \approx \begin{array}{l} n^{\tilde{O}(D)} \\ \text{time} \\ \text{algorithms} \end{array}$$

Further Discussion

① sum-of-squares

[Parrilo]
[Lasserre]

Powerful semidefinite programming hierarchy

Goal: Find an operator that behaves like the expectation over a distribution on solutions

$$\tilde{\mathbb{E}}: \mathcal{V}_n^{\leq D} \rightarrow \mathbb{R}$$

Pseudoexpectation

Constraints

1. $\tilde{\mathbb{E}}$ is linear

2. $\tilde{\mathbb{E}}[1]$ is 1

3. $\tilde{\mathbb{E}}[p^2] \geq 0$

for all poly. p
of $\deg(p) \leq \frac{D}{2}$

Specialization to P.C.

4. $\tilde{\mathbb{E}}[x_i^2 p] = \tilde{\mathbb{E}}[x_i p]$

booleanity

5. $\tilde{\mathbb{E}}[\sum x_i] = k$

clique size

6. $\tilde{\mathbb{E}}[x_i x_j p] = 0$

$\forall (i,j) \notin E, \deg(p) \leq D-2$

Given a clique, can define a valid $\tilde{\mathbb{E}}$ as

$$\tilde{\mathbb{E}}[p(x_1, \dots, x_n)] = p(a_1, a_2, \dots, a_n)$$

where a_1, \dots, a_n is indicator vector of a k -clique

Claim: This $\tilde{\mathbb{E}}$ satisfies 1. - 6.

Proposition: There is an $n^{o(D)}$ time algorithm for finding a valid $\tilde{\mathbb{E}}$, if it exists

Main idea: Constraint 3. is a PSD-ness condition

Thm [Barak et al] For any $D = o(\log n)$ SOS fails to detect cliques of size $n^{1/2 - o(1)}$

Many other problems where we believe there is a computational vs statistical

tradeoff are open

Does SOS fail? Probably, but very difficult to show

(2) Overlap gap property

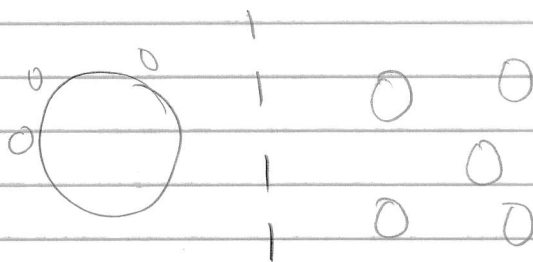
Given $G \sim G(n, 1/2)$, we know that w.h.p the largest clique has size $(2 \pm o(1)) \log n$

Best known efficient algorithm only finds a clique of size $\log n$.

Main Idea [Gamarnik, Sudan], many others
Does the landscape exhibit a gap?

* pair of near optimal solutions has distance $\leq \alpha$ or $\geq \beta$

Let's look at space of k cliques



giant component $k = \log n$ shatters

③ statistical query lower bounds

[Kreuzer] Don't directly see samples but can query functions $f: \mathbb{R}^n \rightarrow [-1, 1]$

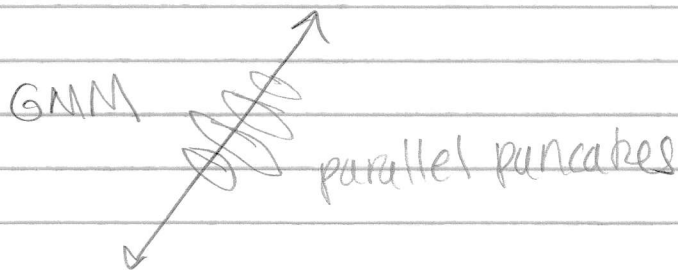
[Feldman et al]

VSTAT(ϵ): Get response

$$v \in [E_{x \sim p}[f(x)] \pm \epsilon]$$

$$\text{where } \epsilon = \max \left\{ \frac{1}{t}, \sqrt{\frac{\text{var}_{x \sim p}[f(x)]}{t}} \right\}$$

Example of computational vs statistical tradeoff in distribution learning



[Diakonikolas, Kane, Stewart] Any SQ algorithm for learning k -GMMs need $n^{\Omega(k)}$ queries or tolerance $n^{-o(k)}$

Intuition: It's hard to find the interesting direction, and o.w. it looks like a single Gaussian

