

9/25/23

Lecture 6: SoS and Gaussian mixtures:

Let $x_1, \dots, x_n \in \mathbb{R}^d$ be samples from

$$q = \frac{1}{k} \sum_{j=1}^k N(\mu_j, \text{Id})$$

Define $\Delta \stackrel{\text{def}}{=} \min_{i \neq j} \|\mu_i - \mu_j\|_2$

$$N \stackrel{\text{def}}{=} \frac{n}{k} \quad (\approx \# \text{ pts in each component})$$

Let t (SoS degree) be power of 2. Suppose

$$\Delta \gg \sqrt{t} k^{1/t}.$$

SoS program

Variables: a_1, \dots, a_n (1-dimensional)
 μ (d-dimensional)

Constraints:

- 1) $a_i^2 = a_i$ (Boolean indicators)
- 2) $\sum_{i=1}^n a_i = N$ (selects out enough points for one component)
- 3) $\frac{1}{N} \sum_i a_i x_i = \mu$ (μ is empirical mean of points selected)
- 4) $\frac{1}{N} \sum_i a_i \langle u, x_i - \mu \rangle^t \leq 2^{t/2} \|u\|_2^t$ (selected points have Gaussian moment bounds)
for all vectors u

Let $S_j \subset [n]$ denote samples from $N(\mu_j, \text{Id})$

For convenience, we will pretend $|S_j| = N$ exactly $\forall j$.

* Seems like infinite constraints... see below for how to quantify over all $u \in \mathbb{R}^{d-1}$

For now, assume $d=1$

So constraint 4) becomes $\frac{1}{N} \sum_i a_i (x_i - \mu)^t \leq 2t^{t/2}$.

Warmup lemma: Let $S = S_j, \mu = \mu_j$ for any $j \in [k]$.

Overlap
b/w our
points and
 S^c

There is a $\text{deg-}O(t)$ proof that:

$$\left(\sum_{i \in S^c} a_i \right)^t (\mu - \mu^c)^t \leq 2^{O(t)} \left(\sum_{i \in S^c} a_i \right)^{t-1} \cdot N \cdot t^{t/2}$$

$$\begin{aligned} \text{PF: } & \left(\sum_{i \in S^c} a_i \right)^t \cdot (\mu - \mu^c)^t \\ &= \left(\sum_{i \in S^c} a_i (\mu - \mu^c) \right)^t \\ &= \left(\sum_{i \in S^c} a_i \left[(\mu - x_i) - (\mu^c - x_i) \right] \right)^t \end{aligned}$$

by degenet
Hölder's
inequality
(SOS-able)

$$\left(\sum_i b_i c_i \right)^t = \left(\sum_i \underbrace{b_i}_{L_p} \underbrace{c_i}_{L_q} \right)^t \leq \left(\sum_i b_i \right)^{t-1} \left(\sum_i c_i \right)^t$$

for $p = \frac{t}{t-1}, q = t$
so $\frac{1}{p} + \frac{1}{q} = 1$

$$\leq \left(\sum_{i \in S^c} a_i \right)^{t-1} \cdot \left(\sum_{i \in S^c} a_i \left[(\mu - x_i) - (\mu^c - x_i) \right]^t \right)$$

$$\begin{aligned} & (a-b)^t \leq 2^t (a+b)^t \\ (*) & \leq 2^t \left(\sum_{i \in S^c} a_i \right)^{t-1} \cdot \left(\sum_{i \in S^c} a_i \left[(\mu - x_i)^t + (\mu^c - x_i)^t \right] \right) \end{aligned}$$

$$\text{Note: } \sum_{i \in S^c} a_i (\mu - x_i)^t \leq \sum_{\text{all } i} a_i (\mu_i - x_i)^t \leq N \cdot 2t^{t/2}$$

moment bound, i.e. constraint 4

$$\sum_{i \in S^0} a_i (\mu^0 - x_i)^t \leq \sum_{i \in S^0} (\mu^0 - x_i)^t \leq N \cdot 2^{t+1/2}$$

Boundedness,
i.e. constraint 1

assuming t -th
empirical moment of actual
samples from component
concentrate

So

$$\left(\sum_{i \in S^0} a_i \right)^t (\mu - \mu^0)^t \leq 2^{(t+1)} \left(\sum_{i \in S^0} a_i \right)^{t-1} \cdot N \cdot 2^{t/2} \quad \square$$

Note, if we could "divide on both sides" and take t -th roots, we would get

$$|\mu - \mu^0| \leq \left(\frac{1}{N} \sum_{i \in S^0} a_i \right)^{-1/t} \cdot \sqrt{t} \quad (\dagger)$$

i.e. if overlap between our points (chosen by a_i) and true points in component S^0 is large, then our μ is close to the mean of S^0 .

Claim 1: If a_i 's were real indicators of a set S satisfying (6) for every center $\mu^0 = \mu_j$, then

Component S_{j^*} with largest overlap with S satisfies $|S \cap S_{j^*}| = (1 - \delta)N$ for $\delta \leq kt^{1/2} \cdot O(1/\Delta)^t$.
($\ll 1$)

Pf: Note $1 - \delta \geq \frac{1}{k}$ by averaging,

and $|S \cap S_j| \geq \frac{\delta}{k} \cdot N$ for some $j \neq j^*$. So:

$$|\mu_{j^*} - \mu| \leq (1 - \delta)^{-1/t} \sqrt{t} \leq k^{1/t} \sqrt{t} \ll \Delta/2,$$

so $|\mu_j - \mu| > \frac{\delta}{2}$. Thus, by $(*)$ applied to comp. j ,
Eq. $(*)$ applied to comp. j

$$\frac{\Delta}{2} < |\mu_j - \mu| \leq \left(\frac{\delta}{k}\right)^{-1/t} \sqrt{t},$$

$$\text{so } \delta^{1/t} \leq \frac{k^{1/t} \sqrt{t}}{\Delta} \quad (\ll 1),$$

and thus $\delta \leq kt^{1/2} \cdot O(1/\Delta)^t$ as claimed. \square

i.e. a_i 's must have $(\ll 1)$ overlap with some component!

Σ OU's:

- SoS version of claim 1? \curvearrowright closely related
- rounding SoS solution? \curvearrowright related
- $d > 1$?

Issue with Claim 1 is it breaks symmetry across clusters. Makes it unclear how to round.

Claim 2 (symmetric version of Claim 1 - still not SoS):

If a_i 's are indicator of S , then

$$\sum_{j=1}^k \left(\frac{|S_j \cap S|}{N} \right)^2 \geq 1 - k^{2+1/2} \cdot O(1/\Delta)^t$$

Note: This implies Claim 1. Define $c_j = \frac{|S_j \cap S|}{N}$

so $\sum_j c_j = 1$. Thus

$$1 - k^{2+1/2} \cdot O(1/\Delta)^t \leq \sum_j c_j^2 \leq (\max_j c_j) \cdot \sum_j c_j = \max_j c_j,$$

i.e. exists j s.t. $\frac{1}{N} |S_j \cap S| \geq 1 - k^{2+1/2} \cdot O(1/\Delta)^t$,

which recovers Claim 1 w/ extra (but unimportant) k factor.

Pf: Define $c_j = \frac{|S_j \cap S|}{N}$. Then

$$1 = \left(\sum_j c_j \right)^2 = \sum_j c_j^2 + 2 \sum_{i < j} c_i c_j$$

we'll show
these are small

$$c_i^{1/t} c_j^{1/t} \leq c_i^{1/t} c_j^{1/t} \frac{|m_i - m| + |m_j - m|}{\Delta}$$

≥ 1 by triangle inequality

$$\leq c_i^{1/t} \frac{|m_i - m|}{\Delta} + c_j^{1/t} \frac{|m_j - m|}{\Delta}$$

Recall by warmup lemma, in particular (b),

$$|m_i - m| \leq c_i^{-1/t} \cdot \sqrt{t}, \text{ so}$$

$$\leq \frac{\sqrt{t}}{\Delta},$$

thus $\sum_{i \neq j} c_i \cdot c_j \leq k^2 \cdot O\left(\frac{\sqrt{t}}{\Delta}\right)^t$ as desired. \square

Next. "SoS-ize" Claim 2

Claim 3 (SoS version of Claim 2):

For any deg- t pseudodistribution over $\{a_i\}$, μ ,

$$\mathbb{E} \left[\sum_{j=1}^k \left(\frac{1}{N} \sum_{i \in S_j} a_i \right)^2 \right] \geq 1 - k^{2+t/2} \cdot O(1/\Delta)^t$$

Pf: Define $c_j \stackrel{\Delta}{=} \frac{1}{N} \sum_{i \in S_j} a_i$ (now a deg-1 polynomial)

Recall the only thing we have proved about \mathbb{E}^h is that for all $j \in [k]$,

$$c_j^+ \cdot (\mu - \mu_j)^+ \leq O(t)^{+1/2} \cdot c_j^{t-1}. \quad (1)$$

Next, note that

$$\begin{aligned} \Delta^+ &\leq (\mu_i - \mu_j)^+ \\ &= [(\mu_i - \mu) - (\mu_j - \mu)]^+ \\ &\leq 2^+ [(\mu_i - \mu)^+ + (\mu_j - \mu)^+], \end{aligned}$$

$$\text{So } \frac{(\mu_i - \mu)^+ + (\mu_j - \mu)^+}{(\Delta/2)^+} \leq 1 \quad (2)$$

Combining (1) and (2) yields

$$c_i^+ c_j^+ \stackrel{(2)}{\leq} c_i^+ c_j^+ \frac{(\mu_i - \mu)^+ + (\mu_j - \mu)^+}{(\Delta/2)^+}$$

$$\leq (2/\Delta)^+ \cdot \left[c_j^+ \underbrace{c_i^+ (\mu_i - \mu)^+}_{\leq 1} + c_i^+ \underbrace{c_j^+ (\mu_j - \mu)^+}_{\leq 1} \right]$$

$$\begin{aligned}
 (1) \\
 &\leq \left(2\sqrt{t}/\Delta\right)^t \begin{pmatrix} t & t-1 \\ c_j^t c_i^{t-1} & + c_i^t c_j^{t-1} \end{pmatrix} \\
 &\stackrel{c_i, c_j \leq 1}{\leq} 2 \cdot \left(2\sqrt{t}/\Delta\right)^t c_i^{t-1} c_j^{t-1}
 \end{aligned}$$

So we have proved in $\text{deg-}O(t)$ SoS that

$$c_i^t c_j^t \leq O(\sqrt{t}/\Delta)^t c_i^{t-1} c_j^{t-1}.$$

To avoid working with odd powers, square both sides to get

$$c_i^{2t} c_j^{2t} \leq O(t/\Delta^2)^t c_i^{2t-2} c_j^{2t-2}$$

Thus,

$$\mathbb{E}^h \left[c_i^{2t} c_j^{2t} \right] \leq O(t/\Delta^2)^t \mathbb{E}^h \left[c_i^{2t-2} c_j^{2t-2} \right] \quad (\text{td})$$

Q: How do we simulate taking t -th roots in SoS?

A: "pseudo-expectation Cauchy-Schwarz / Hölder's inequalities"

Fact ("Cauchy-Schwarz"):

$$\mathbb{E}^h \left[p(x) q(x) \right] \leq \mathbb{E}^h \left[p(x) \right]^{1/2} \cdot \mathbb{E}^h \left[q(x) \right]^{1/2}$$

for any p, q of degree $\leq t/2$ and \mathbb{E}^h any $\text{deg-}t$ pseudo-expectation.

Fact ("Hölder's"):

$$\tilde{\mathbb{E}}[p(x)^{t-2}] \leq \tilde{\mathbb{E}}[p(x)^t]^{\frac{t-2}{t}}$$

For any deg- l sum of squares polynomial p and $\tilde{\mathbb{E}}$ any deg- t pseudo-expectation.

Pfs: Pset 2.

Applying pseudo-exp. Hölder's to (36), we get

$$\tilde{\mathbb{E}} \left[\begin{matrix} 2t & 2t \\ c_i & c_j \end{matrix} \right] \leq O(t/\Delta^2)^t \tilde{\mathbb{E}} \left[\begin{matrix} 2t & 2t \\ c_i & c_j \end{matrix} \right]^{\frac{t-1}{t}}$$

Now we
can divide
freely

$$\Rightarrow \tilde{\mathbb{E}} \left[\begin{matrix} 2t & 2t \\ c_i & c_j \end{matrix} \right] \leq O(t/\Delta^2)^{t^2}$$

Applying pseudo-exp Cauchy-Schwarz, we have

$$\begin{aligned} \tilde{\mathbb{E}}[c_i c_j] &= \tilde{\mathbb{E}}[\underbrace{c_i c_j}_{\text{red}} \cdot \underbrace{1}_{\text{green}}] \\ &\leq \tilde{\mathbb{E}}[(c_i c_j)^2]^{1/2} \cdot \cancel{\tilde{\mathbb{E}}[1]^{1/2}} \end{aligned}$$

and repeating this $\log_2 t$ times, get

$$\tilde{\mathbb{E}}[c_i c_j] \leq \tilde{\mathbb{E}}[\underbrace{c_i^{2t} c_j^{2t}}]^{1/2t}.$$

Thus, $\tilde{\mathbb{E}}[c_i c_j] \leq O(t/\Delta^2)^{t/2} \quad \forall i \neq j,$

and thus

$$\tilde{\mathbb{E}}\left[\sum_j c_j^2\right] = \tilde{\mathbb{E}}\left[\underbrace{\left(\sum_j c_j\right)^2}_i - \sum_{i \neq j} c_i c_j\right]$$

$$\geq 1 - k^2 t^{t/2} O(1/\Delta)^t$$

as desired. □

Rounding:

Can't just output $\tilde{\mathbb{E}}[M]$ b/c $\{a_i\}$'s don't
prefer any particular component...

How do we know $\{a_i\}$'s are indicating a fixed component,
or a dist over components?

Trick: entropy maximization

Want pseudo-dist over $\{a_i\}$'s to resemble uniform distribution over true indicators $\{a_i^{(j)}\}$'s, where

$$a_i^{(j)} \stackrel{\Delta}{=} \mathbb{1}[x_i \text{ from } N(\mu_j, \text{Id})]$$

This distribution has high "entropy" as quantified by

$$\left\| \bigoplus_j \left[a^{(j)} (a^{(j)})^T \right] \right\|_F^2 \quad (\text{ENT})$$

Note:

$$\begin{aligned} (\text{ENT}) &= \sum_{j, j'=1}^k \frac{1}{k^2} \langle a^{(j)} (a^{(j)})^T, a^{(j')} (a^{(j')})^T \rangle \\ &= \sum_{j, j'=1}^k \frac{1}{k^2} \underbrace{\langle a^{(j)}, a^{(j')} \rangle^2}_{\substack{= \\ 0 \text{ if } j \neq j' \\ \text{b/c components disjoint,}}} \\ &= \frac{1}{k^2} \sum_{j=1}^k \|a^{(j)}\|_2^4 = \frac{N^2}{k} \end{aligned}$$

We pick the pseudo-distribution solving

$$\max_{\tilde{\mathbb{P}}} \left\| \bigoplus \left[\underbrace{(a_1, \dots, a_n)}_{\triangleq a} (a_1, \dots, a_n)^T \right] \right\|_F$$

subject to $\tilde{\mathbb{P}}$ satisfying constraints of the program.

Lemma: This $\hat{\mathbb{E}}^n$ satisfies

$$\begin{aligned} & \left\| \hat{\mathbb{E}}^n [aa^T] - \mathbb{E}_j \left[a^{(j)} (a^{(j)})^T \right] \right\|_F^2 \quad (f) \\ & \leq \left\| \mathbb{E}_j \left[a^{(j)} (a^{(j)})^T \right] \right\|_F^2 \underbrace{\left(k^{2+t/2} \cdot O(1/\delta)^t \right)}_{\ll 1} \end{aligned}$$

Pf: Because unif distribution over $\{ \{a_i\}_i, \mu_j \}_j$ is a feasible solution, $\| \hat{\mathbb{E}}^n [aa^T] \|_F^2 \leq \frac{N}{k}$, so

$$\begin{aligned} (f) &= \frac{2N^2}{k} - \frac{2}{k} \sum_{j=1}^k \hat{\mathbb{E}}^n \left[\langle a, a^{(j)} \rangle^2 \right] \\ &= \frac{2N^2}{k} - \frac{2}{k} \sum_{j=1}^k \hat{\mathbb{E}}^n \left[\underbrace{\left(\sum_{i \in S_j} a_i \right)^2}_{= Nc_j} \right] \end{aligned}$$

$$= \frac{2N^2}{k} \left(1 - \sum_j c_j^2 \right)$$

$$\leq \frac{N^2}{k} \left(k^{2+t/2} \cdot O(1/\delta)^t \right).$$

□

Note,

$$\bigoplus_j \left[a^{(j)} (a^{(j)})^T \right] = \overset{n}{\left(\begin{array}{c} \frac{1}{\kappa} \\ \frac{1}{\kappa} \\ \frac{1}{\kappa} \\ \frac{1}{\kappa} \\ \dots \end{array} \right)}$$

(after row/col permutation),

so Lemma implies that we can read off clustering from $\hat{\bigoplus}^n [aa^T]$!

Final IOU: ... this was all for $d=1$!

Warmup lemma and main Claim 3 easy to generalize, e.g.

Before:

$$\left(\sum_{i \in S_j} a_i \right)^t (n - n_j)^t \leq 2^{O(t)} \left(\sum_{i \in S_j} a_i \right)^{t-1} N \cdot t^{t/2}$$

After

$$\left(\sum_{i \in S_j} a_i \right)^t \| \mu - \mu_j \|_2^t \leq 2^{O(t)} \left(\sum_{i \in S_j} a_i \right)^{t-1} N \cdot t^{t/2}$$

But $\|\mu - \mu_j\|_2^2 = \langle \mu - \mu_j, \mu - \mu_j \rangle$, so
 can just "project" data along $\mu - \mu_j$ direction
 and reduce to 1D proof.

(need to be careful b/c $\mu - \mu_j$ is not a real vector
 because μ is an SAS variable)

trickier: how to impose constraint

$$\frac{1}{N} \sum_{i=1}^n a_i \langle u, x_i - \mu \rangle^t \leq 2t^{t/2} \|u\|_2^t$$

for all $u \in \mathbb{R}^d$?

Because we will apply this to $u = \mu - \mu_j$, need
 this to make sense even when u is not a real vector...

Idea: constrain via

$\binom{***}{\text{u d + polynomial constraints}}$

$$\left\| \frac{1}{N} \sum_{i=1}^n a_i (X_i - \mu)^{\otimes t/2} \left[(X_i - \mu)^{\otimes t/2} \right]^T - \bigoplus_{g \sim N(0, Id)} \left(g^{\otimes t/2} \left(g^{\otimes t/2} \right)^T \right) \right\|_F^2 \leq 1$$

(satisfied by $a_i = \binom{[i]}{i}$ and $\mu = \mu_j$, if n large enough)

i.e. pick out subset s.t. empirical order- t moments are close to those of $N(0, Id)$.

Fact: For an SoS variable u ,

$$\mathbb{E}_{g \sim N(0, Id)} \langle g, u \rangle^t \leq t^{t/2} \cdot \|u\|_2^t$$

has a deg- t SoS proof in u .

Pf:

$$\mathbb{E}_g \langle g, u \rangle^t = \sum_{\substack{\text{deg-}t \text{ monomials } \alpha \\ \text{s.t. every} \\ \text{variable appears} \\ \text{even \# times}}} u_\alpha \mathbb{E}[g_\alpha]$$

$$\leq t^{t/2} \sum_{\alpha} u_\alpha^2$$

$$= t^{t/2} \|u\|_2^t.$$

i.e. $N(0, Id)$ is "certifiably t -hypercontractive" □

Note, if we take constraint $(\delta \delta \delta)$ and hit it on both sides with $\left[(\mu - \mu_j)^{\otimes t/2} \right]^T \cdots (\mu + \mu_j)^{\otimes t/2}$, we get:

$$\frac{1}{N} \sum_{i=1}^n a_i \langle \mu - \mu_j, x_i - \mu \rangle^+ - \mathbb{E}_g \langle \mu - \mu_j, g \rangle^+ \leq \|\mu - \mu_j\|_2^+$$

⇓ (using Fact above)

$$\begin{aligned} \frac{1}{N} \sum_{i=1}^n a_i \langle \mu - \mu_j, x_i - \mu \rangle^+ &\leq (1 + t^{+1/2}) \|\mu - \mu_j\|_2^{+1/2} \\ &\leq O(t)^{+1/2} \|\mu - \mu_j\|_2^{+1/2}, \end{aligned}$$

which is sufficient to prove high-dim generalization of warmup lemma and its consequences.