# Lecture 7: SoS and Gaussian mixtures:

Let $x_1, \ldots, x_n \in \mathbb{R}^d$ be samples from
$$q = \frac{1}{k} \sum_{j=1}^{k} N(\mu_j, \text{Id})$$

Define $\Delta \triangleq \min_{i \neq j} \|\mu_i - \mu_j\|_2$

$\quad N \triangleq \frac{n}{k}$ $\quad$ ($\approx$ # pts in each component)

Let $t$ (sos degree) be power of 2. Suppose
$$\Delta \gg \sqrt{t} \, k^{1/t}.$$

### SoS program

Variables: $\quad a_1, \ldots, a_n \quad$ (1-dimensional)

$\quad\quad\quad\quad\quad \mu \quad\quad\quad\quad$ (d-dimensional)

Constraints:

1). $a_i^2 = a_i$ $\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad$ (Boolean indicators)

2). $\sum_{i=1}^{n} a_i = N$ $\quad\quad\quad\quad\quad\quad\quad$ (selects out enough points for one component)

3). $\frac{1}{N} \sum_i a_i x_i = \mu$ $\quad\quad\quad\quad\quad$ ($\mu$ is empirical mean of points selected)

4). $\frac{1}{N} \sum_i a_i \langle u, x_i - \mu \rangle^t \leq 2t^{t/2} \|u\|_2^t$ $\quad$ (selected points have Gaussian moment bounds)
$\quad\quad$ for all vectors $u$ ${}^{\textcolor{red}{*}}$

---

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

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

---

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

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

Overlap b/t our points and $S^\bullet$

There is a deg-$O(t)$ proof that:

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

Pf: $\left( \sum_{i \in S^\bullet} a_i \right)^t \cdot (\mu - \mu^\bullet)^t$

$$= \left( \sum_{i \in S^\bullet} a_i (\mu - \mu^\bullet) \right)^t$$

$$= \left( \sum_{i \in S^\bullet} a_i \left[ (\mu - x_i) - (\mu^\bullet - x_i) \right] \right)^t$$

by degree-$t$ Hölder's inequality (SOS-able):

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

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

so $\frac{1}{p} + \frac{1}{q} = 1$

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

$(a-b)^t \leq 2^t (a^t + b^t)$

(☆)

$$\leq 2^t \left( \sum_{i \in S^\bullet} a_i \right)^{t-1} \cdot \left( \sum_{i \in S^\bullet} a_i \left[ (\mu - x_i)^t + (\mu^\bullet - x_i)^t \right] \right)$$

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

↑ moment bound, i.e. constraint 4

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

Bodeanity,
i.e. constraint 1

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

So

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

---

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

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

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

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

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

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

and $|S \cap S_j| \ge \frac{\delta}{k} \cdot N$ for some $j \ne j^\bullet$. So:

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

$\qquad\qquad$ Eq. (s)$^T$ applied to comp. $j^\bullet$

So $|\mu_j - \mu| > \frac{\Delta}{2}$. Thus, by $(\ast)$ applied to comp. $j$,

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

So $\delta^{1/t} \lesssim \frac{k^{1/t} \sqrt{t}}{\Delta} \; (\ll 1)$,

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

i.e. $a_i$'s must have $1-o(1)$ overlap with some

Component!

$\sum OU$'s:
- SoS version of claim 1 ?  $\Big)$ closely related
- rounding SoS solution ?
- $d > 1$ ?

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

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

---

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

so $\sum_j c_j = 1$. Thus

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

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

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

---

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

$$1 = \left( \sum_j c_j \right)^2 = \sum_j c_j^2 + 2 \underbrace{\sum_{i<j} c_i c_j}_{\text{we'll show these are small}}$$

$$c_i^{1/t} \, c_j^{1/t} \leq c_i^{1/t} \, c_j^{1/t} \, \underbrace{\frac{|\mu_i - \mu| + |\mu_j - \mu|}{\Delta}}_{\geq 1 \text{ by triangle ineq}}$$

$$\leq c_i^{1/t} \, \frac{|\mu_i - \mu|}{\Delta} + c_j^{1/t} \, \frac{|\mu_j - \mu|}{\Delta}$$

Recall by warmup lemma, in particular $(\textstyle \pm)$,

$$|\mu_i - \mu| \leq c_i^{-1/t} \cdot \sqrt{t}, \quad \text{so}$$

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

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

Next. "SoS-ize" Claim 2  (the following is pedantic and can be ignored upon first reading)

Claim 3 ( SoS version of Claim 2):

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

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

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

Recall the only thing we have proved about $\overset{\wedge}{\mathcal{E}}$ is that for all $j \in [k]$,

$$c_j^t \cdot \left(\mu - \mu_j\right)^t \leq O(t)^{t/2} \cdot c_j^{t-1}. \qquad (1)$$

Next, note that

$$\Delta^t \leq \left(\mu_i - \mu_j\right)^t$$
$$= \left[\left(\mu_i - \mu\right) - \left(\mu_j - \mu\right)\right]^t$$
$$\leq 2^t \left[\left(\mu_i - \mu\right)^t + \left(\mu_j - \mu\right)^t\right],$$

So $\dfrac{\left(\mu_i - \mu\right)^t + \left(\mu_j - \mu\right)^t}{\left(\Delta/2\right)^t} \leq 1 \qquad (2)$

Combining (1) and (2) yields

$$c_i^t \, c_j^t \overset{(2)}{\leq} c_i^t \, c_j^t \, \frac{\left(\mu_i - \mu\right)^t + \left(\mu_j - \mu\right)^t}{\left(\Delta/2\right)^t}$$

$$\leq (2/\Delta)^t \cdot \left[c_j^t \, c_i^t \left(\mu_i - \mu\right)^t + c_i^t \, c_j^t \left(\mu_j - \mu\right)^t\right]$$

(1)
$$\leq \left(2\sqrt{t}/\Delta\right)^t \left(c_j^t c_i^{t-1} + c_i^t c_j^{t-1}\right)$$

$c_i, c_j \leq 1$
$$\leq 2 \cdot \left(2\sqrt{t}/\Delta\right)^t c_i^{t-1} c_j^{t-1}$$

So we have proved in 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,
$$\tilde{\mathbb{E}}\left[c_i^{2t} c_j^{2t}\right] \leq O(t/\Delta^2)^t \tilde{\mathbb{E}}\left[c_i^{2t-2} c_j^{2t-2}\right] \quad (\star\star)$$

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

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

Fact ("Cauchy-Schwarz"):
$$\tilde{\mathbb{E}}\left[p(x)\,q(x)\right] \leq \tilde{\mathbb{E}}\left[p(x)^2\right]^{1/2} \cdot \tilde{\mathbb{E}}\left[q(x)^2\right]^{1/2}$$

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

# Fact ("Hölder's"):

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

for any deg-$\ell$ __sum of squares__ polynomial $p$ and $\tilde{\mathbb{E}}$ any deg-$t\ell$ pseudo-expectation.

__Pfs:__  Pset 2.  $\square$

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

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

Now we can divide freely

$$\Rightarrow \tilde{\mathbb{E}}\left(c_i^{2t} c_j^{2t}\right) \leq O(t/\Delta^2)^{t^2}$$

Applying pseudo-exp Cauchy-Schwarz, we have

$$\tilde{\mathbb{E}}\left(c_i c_j\right) = \tilde{\mathbb{E}}\left(c_i c_j \cdot \underline{1}\right)$$
$$\leq \tilde{\mathbb{E}}\left[(c_i c_j)^2\right]^{1/2} \cdot \cancel{\tilde{\mathbb{E}}\left[1^2\right]^{1/2}}$$

and repeating this $\log_2 2t$ times, get

$$\overset{n}{\underset{}{\mathbb{E}}}\left(c_i c_j\right) \leq \overset{n}{\underset{}{\mathbb{E}}}\left(\underline{c_i^{2t}\ c_j^{2t}}\right)^{1/2t}.$$

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

and thus

$$\overset{n}{\underset{}{\mathbb{E}}}\left[\sum_j c_j^2\right] = \overset{n}{\underset{}{\mathbb{E}}}\left[\underbrace{\left(\sum_j c_j\right)^2}_{\text{"}1\text{"}} - \sum_{i \neq j} c_i c_j\right]$$

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

as desired. ⬠

## Rounding:

Can't just output $\overset{n}{\underset{}{\mathbb{E}}}[\mu]$ b/c $\{a_i\}$'s don't preference any particular component...

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

Trick: <u>entropy maximization</u>

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

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

This distribution has high "entropy" as quantified by

$$\left\| \mathop{\mathbb{E}}_{j} \left[ a^{(j)}(a^{(j)})^T \right] \right\|_F^2 \qquad (\text{ENT})$$

Note:

$$(\text{ENT}) = \sum_{j,j'=1}^{k} \frac{1}{k^2} \left\langle a^{(j)}(a^{(j)})^T, \; a^{(j')}(a^{(j')})^T \right\rangle$$

$$= \sum_{j,j'=1}^{k} \frac{1}{k^2} \underbrace{\left\langle a^{(j)}, a^{(j')} \right\rangle^2}_{\substack{= \\ 0 \text{ if } j \neq j' \\ \text{b/c components disjoint,}}}$$

$$= \frac{1}{k^2} \sum_{j=1}^{k} \left\| a^{(j)} \right\|_2^4 = \frac{N^2}{k}$$

We pick the pseudo-distribution solving

$$\min_{\tilde{\mathbb{E}}} \left\| \tilde{\mathbb{E}} \left[ \underbrace{(a_1, \ldots, a_n)(a_1, \ldots, a_n)^T}_{\triangleq a} \right] \right\|_F^2$$

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

**Lemma** : The $\hat{\mathbb{E}}$ satisfies

$$\left\| \hat{\mathbb{E}}\left[aa^T\right] - \mathbb{E}_j\left[a^{(j)}(a^{(j)})^T\right] \right\|_F^2 \quad (\dagger)$$

$$\leq \left\| \mathbb{E}_j\left[a^{(j)}(a^{(j)})^T\right] \right\|_F^2 \underbrace{\left(k^2 \cdot t^{t/2} \cdot O(1/\delta)^t\right)}_{<<1}$$

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

$$(\dagger) = \frac{2N^2}{k} - \frac{2}{k}\sum_{j=1}^{k} \hat{\mathbb{E}}\left[\langle a, a^{(j)}\rangle^2\right]$$

$$= \frac{2N^2}{k} - \frac{2}{k}\sum_{j=1}^{k} \hat{\mathbb{E}}\left[\underbrace{\left(\sum_{i\in S_j} a_i\right)^2}_{= Nc_j}\right]$$

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

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

Note,

$$\mathop{\mathbb{E}}_{j}\left[a^{(j)}(a^{(j)})^{T}\right] = \hat{} \begin{pmatrix} \boxed{\frac{1}{k}} & & & & \\ & \boxed{\frac{1}{k}} & & & \\ & & \boxed{\frac{1}{k}} & & \\ & & & \boxed{\frac{1}{k}} & \\ & & & & \ddots \end{pmatrix}^{n}$$

(after row/col permutation),

so Lemma implies that we can read off clustering from $\hat{\mathbb{E}}\left[a a^{T}\right]$!

---

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

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

Before:

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

After

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

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.

<span style="color:red">(need to be careful b/c $\mu - \mu_j$ is not a real vector

because $\mu$ is an SoS variable)</span>

trickier: how to impose constraint

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

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

<span style="color:red">Because we will apply this to $u = \mu - \mu_j$, need

this to make sense even when $u$ is not a real vector...</span>

<u>Idea</u>: Constrain via

<span style="color:orange">in d t
polynomial
constraints</span>

$$(\bm{***}) \left\{ \left\| \frac{1}{N} \sum_{i=1}^{\hat{}} a_i (x_i - \mu)^{\otimes t/2} \left[ (x_i - \mu)^{\otimes t/2} \right]^T - \mathop{\mathbb{E}}_{g \sim N(0, Id)} \left( g^{\otimes t/2} \left( g^{\otimes t/2} \right)^T \right) \right\|_F^2 \le 1 \right.$$

<span style="color:green">(Satisfied by $a_i := a_i^{(j)}$ and $\mu = \mu_j$, if $n$ large enough)</span>

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

**Fact**: For an S.o.S variable $u$,

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

has a deg-$t$ SoS proof in $u$.

**Pf**:

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

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

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

i.e. $N(0, Id)$ is $\boxed{\text{"certifiably } t\text{-hypercontractive"}}$

---

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

$$\frac{1}{N}\sum_{i=1}^{\hat{N}} a_i \langle \mu - \mu_j , x_i - \mu \rangle^t - \mathop{\mathbb{E}}_{g}\langle \mu - \mu_j, g \rangle^t \leq \|\mu - \mu_j\|_2^t$$

$$\Downarrow \text{ (using Fact above)}$$

$$\frac{1}{N}\sum_{i=1}^{\hat{N}} a_i \langle \mu - \mu_j, x_i - \mu \rangle^t \leq \left(1 + t^{t/2}\right) \|\mu - \mu_j\|_2^t$$

$$\leq O(t)^{t/2} \|\mu - \mu_j\|_2^t \,,$$

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