# Learning to Cluster Using SOS

Simple (but still interesting) version
of problem:

K = number of clusters
$d$ = dimension
$n$ = # of samples

We have fixed spherical (but unknown)
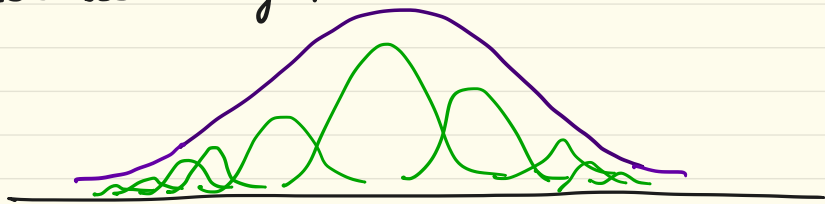gaussians $\mathcal{D}_1, ..., \mathcal{D}_K$; with means
$\mu_1, ..., \mu_K$ , and variance 1.

<u>Input</u> is : $n$ samples $X_1 ... X_n$ drawn
from mixture

<u>Output</u> : means $\mu_1 ... \mu_K$
     or clusters (say which distrib each
     sample is from)

Hard Info Theretically if means are
too close together:



These 2 different mixtures look too similar
- K gaussians have var distance $\sim 2^{-K}$
        from single gaussian. Need $\exp(K)$ samples

Regime of parameters we want: $K \sim d$

$n$ (# samples) $\sim poly(d)$

runtime $poly(d)$

$\Delta$ as small as possible

Separation Assumption $\forall i \neq j \; |\mu_i - \mu_j| \geq \Delta$

Info LB: For $\Delta \lesssim \sqrt{\log d}$ clustering impossible
    using $poly(d)$ samples

## History

d dimensions, k gaussians          $k \sim d$

(1.)     radius of clusters $\sim \sqrt{d}$

$\Delta > 4\sqrt{d}$     greedy clustering
                            <u>easy</u>

(2)   Dasgupta (spectral)

polytime alg for $\Delta = \varepsilon \sqrt{K}$

(3)   Dasgupta-Schulman (EM)

$\Delta \sim K^{1/4}$     polytime

(4)   Regev-Vijayaraghavan (MLE)

$\Delta = O(\sqrt{\log d})$   poly (d) samples
                        but runtime exp (d)

$\Delta = \| \mu_i - \mu_j \|$     ($L_2$ distance in $R^d$)

## New    3 papers    STOC '18

[Hopkins-Li], [Kothari-Steinhardt],
       [Diakonikolas-Kane-Stewart]


**Theorem**  $\exists$ constant $c$

For  $\Delta = c\sqrt{\log d}$,  there is a

quasipoly $(d)$ -time  alg,  error $\frac{1}{\text{poly}(d)}$

For  $\Delta = d^{\varepsilon}$,  poly$(d)$ -time alg, error $\frac{1}{\text{poly}(d)}$


* Also  works  for  <u>robust</u> versions
  of the problem, and for
  more general distributions

## Actual Theorem Statement

Fix $d \approx k$, $\exists t$ s.t. for $n = d^{O(t)} \cdot k^{O(1)}$, there is an algorithm running in time $\text{poly}(n)$ that takes as input random samples $X_1, \ldots, X_n \in \mathbb{R}^d$

$$\left[ \begin{array}{l} S_1, \ldots, S_k \text{ is the true partition of } [n] \text{ s.t.} \\ \{X_i \mid i \in S_i\} \text{ is from } \mathcal{D}_i \text{, and } |S_i| = \frac{n}{k} = N \end{array} \right]$$

and outputs a partition $T_1, \ldots, T_k$ of $[n]$ s.t. $|T_i| = N$, and whp $\forall i$

$$\frac{|S_i \cap T_i|}{N} \geq 1 - k^{10} \cdot \left( \frac{2\sqrt{t}}{\Delta} \right)^t$$

For $\Delta = O\sqrt{\log d}$, choose $t \sim O(\log k)$

$\Delta = k^{\varepsilon}$, choose $t \sim 1000$

# Overview ($d = 1$)

- In order to get an algorithm running in quaspoly time need to show quasipoly many samples suffice That is it is <u>necessary</u> to prove quasipoly sample complexity bounds

- We'll give low degree SOS sample complexity bounds, which will automatically [by SOS automatizability] give us an efficient learning algorithm

This isn't exactly how the algorithm goes...
but correct at a high level

# Important Property of the Samples $\left(\begin{array}{l} d=1, \\ \text{generalizes} \\ \text{easily} \end{array}\right)$

Let $\mathcal{D} = N(\mu, 1)$ be gaussian

and $Y_1 \dots Y_N$ be samples from $\mathcal{D}$.

Then for $N = N(t)$ large enough, whp

$$\mathbb{E}_{j \sim [N]} |Y_j - \bar{\mu}|^t \leq 1.1 \; t^{t/2}$$

↑

sample mean

$$\begin{bmatrix} \text{concentration of} \\ t^{th} \text{ empirical moments} \end{bmatrix}$$

✱ Our algorithm will assume

that the samples $X_1 \dots X_n$ come from

a true Partition $S_1 \dots S_K$, $S_i \subseteq [n]$, $|S_i| = \frac{n}{K} = N$

s.t. $\forall i \in [K]$, $\mathbb{E}_{j \sim S_i} |X_j - \mu_i|^t \leq 2 \cdot t^{t/2}$

↑

empirical avg of

$\{X_j \mid j \in S_i\}$

Can easy calculation show this is true whp.

## Algorithm in the $d=1$ case $\begin{bmatrix} \text{on input} \\ X_1 \dots X_n \end{bmatrix}$

Let $A$ be the following equations:

$$w_i^2 = w_i \qquad i \in [n]$$

$$\sum_{i \in [n]} w_i = N \qquad (N = \frac{n}{k})$$

$$\frac{1}{N} \sum_{i \in [n]} w_i (X_i - \mu)^t \leq 2 \cdot t^{k/2},$$

$$\text{where} \quad \mu = \frac{1}{N} \sum_{i \in [n]} w_i X_i$$

[variables are vectors $w_1 \dots w_n$]

Solve SDP of degree $O(t)$ SOS lift of $A$
to minimize $\lVert w w^T \rVert^2$

$n \times n$.
matrix

Frobenius norm
$L_2^2$ as a vector

<u>Claim</u> If $vv^T$ is solution to degree $O(t)$ SOS lift, then $\| vv^T - aa^T \|^2 \leq 2(\frac{n^2}{k^3} - \langle vv^T, aa^T \rangle)$

<u>Proof</u> Assume WLOG

$$\underbrace{X_1 \ldots X_N}_{S_1} \underbrace{X_{N+1} \ldots X_{2N}}_{S_2} \cdots \cdots \underbrace{X_{(k-1)N+1} \ldots X_{kN}}_{S_K}$$

$\Bigg($ so first $N = \frac{n}{k}$ samples are drawn from $\mathcal{D}_1$, next $N$ samples " " " $\mathcal{D}_2$, and so on. $\Bigg)$

Then these values for $w$ satisfy $A$:

$$q_1 = \begin{bmatrix} \boxed{1} \\ \\ \\ \\ \\ \end{bmatrix} \quad q_2 = \begin{bmatrix} \\ \boxed{1} \\ \\ \\ \end{bmatrix} \quad \cdots \quad q_k = \begin{bmatrix} \\ \\ \\ \boxed{1} \end{bmatrix} \updownarrow \tfrac{N}{2}$$

The corresponding solutions to degree 2 variables $ww^T$ ($w_i w_j$) are:

$$q_1 q_1^T = \begin{bmatrix} \boxed{1} & & \\ & \cdot & \\ & & \cdot \end{bmatrix} \quad q_2 q_2^T = \begin{bmatrix} & & \\ & \boxed{1} & \\ & & \cdot \end{bmatrix} \quad \cdots \quad q_k q_k^T = \begin{bmatrix} \cdot & & \\ & \cdot & \\ & & \boxed{1} \end{bmatrix}$$

($\leftarrow N \rightarrow$ marked on first matrix)

So the avg $aa^T$ also satisfies $A$:

$$\begin{bmatrix} \boxed{\tfrac{1}{k}} & & & \\ & \boxed{\tfrac{1}{k}} & & \\ & & \boxed{\tfrac{1}{k}} & \\ & & & \ddots \\ & & & & \boxed{\tfrac{1}{k}} \end{bmatrix}$$

$$\therefore \quad \| aa^T \|^2 \quad \text{(Frobenius norm = $L_2$ norm as a vector)}$$

$$= \frac{1}{K^2} \cdot K N^2 = \frac{1}{K} \frac{n^2}{K^2} = \frac{n^2}{K^3}$$

Say degree $O(t)$ SDP finds a solution $vv^T$
satisfying degree $O(t)$ SOS of $A$, and
minimizing $\| ww^T \|$. Then $\| vv^T \| \le \| aa^T \| \le \frac{n^2}{K^3}$

Then
$$\| vv^T - aa^T \|^2 = \| vv^T \|^2 + \| aa^T \|^2 - 2 \langle vv^T, aa^T \rangle$$

$$\le 2 \frac{n^2}{K^3} - 2 \langle vv^T, aa^T \rangle$$

$$= 2 \left( \frac{n^2}{K^3} - \langle vv^T, aa^T \rangle \right)$$

■ end of claim

<u>MAIN LEMMA</u>   Let $X_1 .. X_n \in \mathbb{R}$, $S_1 .. S_k$

a partition of $[n]$, $|S_i| = \frac{n}{k}$, s.t. $\forall i \in [k]$:

$$\underset{j \sim S_i}{\mathbb{E}} |X_j - \mu_i|^t = 2 \cdot t^{k/2}$$

Let $w$ be a solution to degree $O(t)$ SDP for $A$.
            (So $v$ is a degree $O(t)$ pseudodistri$t$)

Then $w$ satisfies

$$\sum_{i \in [k]} \left( \frac{|T \cap S_i|}{N} \right)^2 \geq 1 - \frac{2^{O(t)} t^{t/2} k^2}{\Delta^t}$$

where $|T \cap S_i| = \sum_{j \in S_i} w_j$    $\boxed{\begin{array}{c} \text{Proof} \\ \text{after} \end{array}}$

$\overline{\langle ww^T, aa^T \rangle} =$



$$= \sum_{i \in [k]} \underbrace{\left( \sum_{j \in S_i} w_j \right)^2}_{(T \cap S_i)^2} \geq K \cdot N^2 \left( \underline{1 - \frac{2^{O(t)} t^{t/2} k^2}{\Delta^t}} \right)$$

**Corollary** (of claim + MAIN Lemma)

$$\| vv^T - aa^T \| \leq \| aa^T \| \cdot \underbrace{\left( \frac{2^{o(t)} t^{t/2} K^2}{\Delta^t} \right)^{1/2}}_{\varepsilon}$$

**Pf**

By claim

$$\| vv^T - aa^T \|^2 \leq 2\left( \frac{n^2}{K^3} - \langle vv^T, aa^T \rangle \right)$$

By Main Lemma

$$\langle vv^T, aa^T \rangle \geq \frac{1}{K} N^2 \left( 1 - \frac{2^{o(t)} t^{t/2} K^2}{\Delta^t} \right)$$

So

$$\| vv^T - aa^T \| \leq \left[ 2\left( \underbrace{\frac{n^2}{K^3} - \frac{N^2}{K}}_{0} - \underbrace{\frac{N^2}{K}}_{\substack{= \frac{n^2}{K^3} = \| aa^T \|^2}} \frac{2^{o(t)} t^{t/2} K^2}{\Delta^t} \right) \right]^{1/2}$$

$$= \| aa^T \| \cdot \left( \underbrace{\frac{2^{o(t)} t^{t/2} K^2}{\Delta^t}} \right)^{1/2}$$

∎

<u>Recap</u>  $X_1 \ldots X_n$ ,  $S_1 \ldots S_k$

we solve degree $o(t)$ SOS SDP to get

pseudodistribution $ww^T$ s.t.

$$\| ww^T - aa^T \| \leq \| aa^T \| \cdot \left( \text{small fraction} \right)$$

so  $ww^T$ is very close to the
"good" solution  $aa^T$

<u>Note</u> from $aa^T$ we know $S_1 \ldots S_k$
and $ww^T$ is very close entry-wise to $aa^T$

<u>Rounding Alg</u>

(1) Let $I = [n]$ be active indices
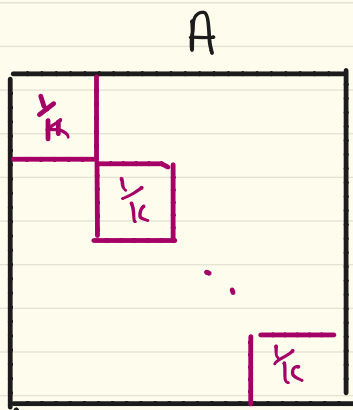
(2) Pick $i \sim I$ uniformly & Let
$S \subseteq I$ be the indices $j$ s.t. $\| M_i - M_j \| \leq \delta \sqrt{\frac{n}{k}}$
( so the rows that are almost )
( the same as row $i$ )

add $S$ to list of clusters & Let $J = I \setminus S$

(3) If $|I| \geq \frac{n}{2k}$ go to (2)

(4) Assign remaining indices to clusters
til all have size $\frac{n}{k}$

A

M



$aa^T$

bad rows

$ww^T$ (close to $aa^T$)

A row $i$ is $\underline{good}$ if $\| M_i - A_i \| < \frac{1}{100} \sqrt{\frac{n}{k}} = \frac{1}{100} \| A_i \|$

There are $\ll \varepsilon^2 n$ bad rows (by avging)

If rounding alg never picks a bad row, then alg will cluster all good rows correctly

Prob. round alg never picks a bad row
is $\le k^2 \varepsilon^2$

since $\sum_{i \in n} \| M_i - A_i \|^2$

$= \| M - A \|^2 \le \varepsilon^2 \| A \|^2$