

Learning mixtures of spherical Gaussians.

$$X = \sum_{i=1}^K w_i \mathcal{N}(\mu_i, I) \subset \mathbb{R}^n.$$

Given samples from  $X$  and you want to learn  $X$ .

(either: learn  $\hat{\lambda}$  s.t.  $d_{TV}(X, \hat{\lambda}) < \epsilon$ ,

learn  $\hat{\mu}_i$  s.t. (up to a perm.)  $|\mu_i - \hat{\mu}_i| < \epsilon$ )

Note: Learning Non-spherical Gaussians has  
 $n^{\Omega(k)}$  lower bounds in the  
SQ model.

Natural approach is method of moments.



[as long as  $\mu$  is not too big]

easy approximate 1st  
d moments of  $X$  for  
any constant  $d$ .

o

o

Question

What are these moments?

$$X = D * G$$

discrete distr.  
over means.

$\uparrow$   
 $N(0, I)$

$$E[X^{(d)}]$$

$$= \sum_i (i! \text{sgn}(d) \binom{d}{i} \mu_i)$$

Not hard to deconvolve.

approximate  $E[D^{(d)}]$

Q Can we learn  $D$  from its low order moments?

A No. Consider 1-d. problem.

first  ~~$d$~~   $d$  moments,  $d+1$  parameters

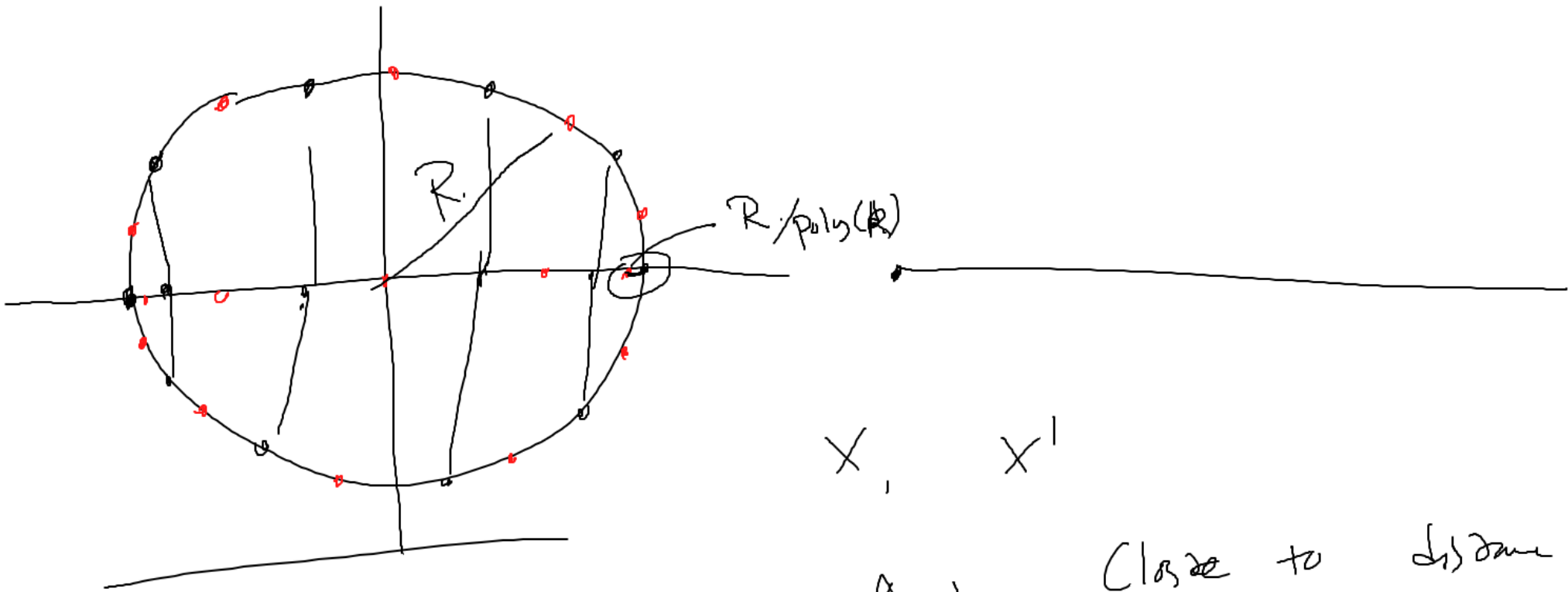
low  $X$  has  $2k$  parameters.

$\Rightarrow$  Unless  $d \geq 2k$ , moments don't determine  $D$ .

$D, D'$  match moments, (even up to small distance)

$N \cdot D, N \cdot D'$  also match moments.

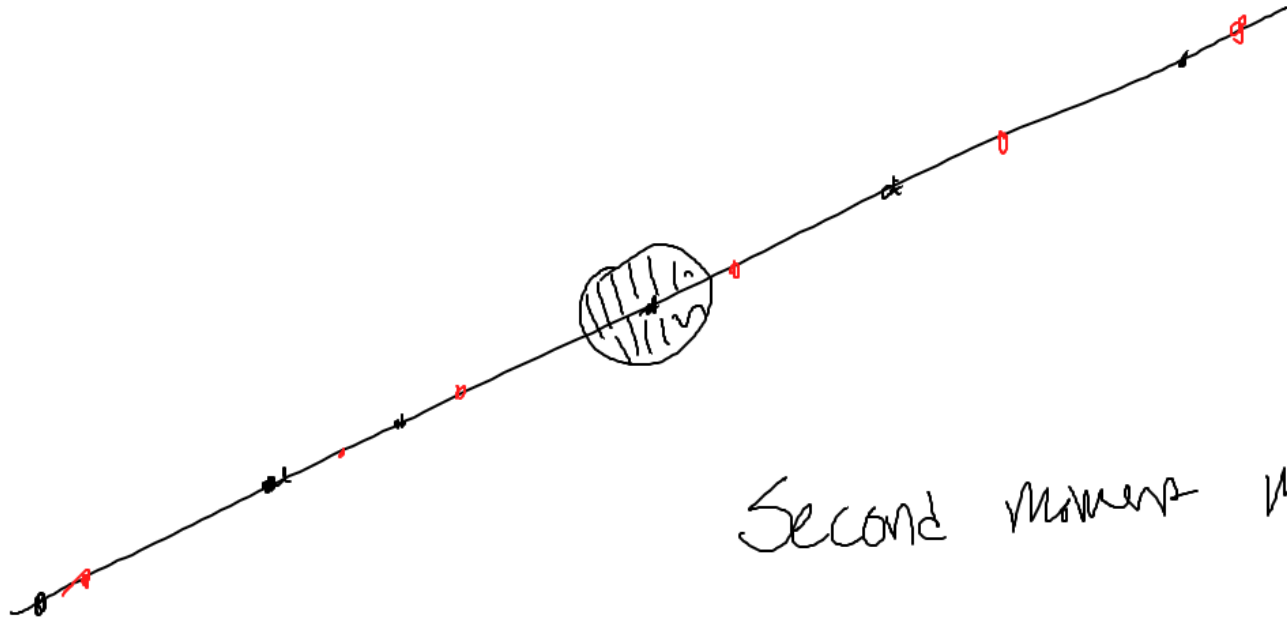




Any 1-d  $X$  approximation  
 $X$  by empirical  
 distribution  $*N(0, \Sigma)$

$X, X'$

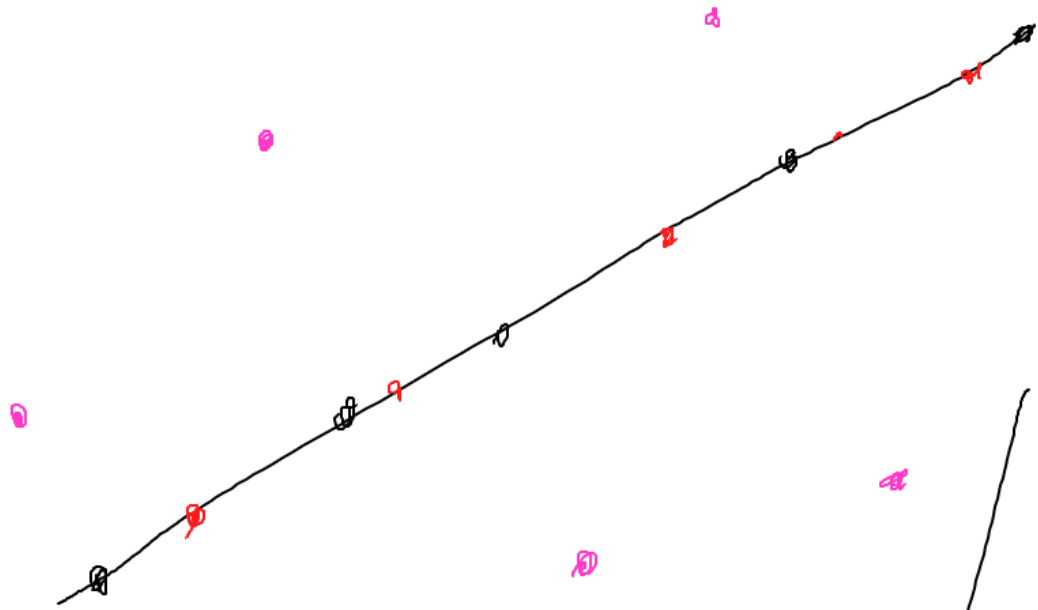
~~Attr~~ Close to distance 1  
 in dtv.  
 but first ~~two~~ moments match.



Second moment matrix of  $(X)$

$$= I + \text{Second moments of } (D)$$

~~Span~~ rank  $k$  ← Suppose on a ~~Span~~  $k$  Subspace of  $\text{dim.}$   
 $\text{Span} \subset \text{span of } \mu_i \text{'s.}$



reduce to  
k-dim  
problem.

to find lower bounds  
(say in SA)

Want is means that  
are spatially symmetric.

Impossible

$$X = D * G.$$

low deg moments of  $X \rightarrow$  low deg moments of  $D$

$$\rightarrow \sum_{i=1}^R w_i \mu_i \quad \text{for any low degree poly. } q$$

$$q(x) = p^2(x) \quad \text{for some } q$$

$$= \sum w_i p^2(\mu_i)$$

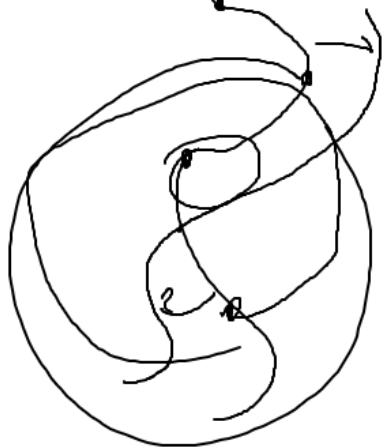
$$= 0 \quad \text{iff } \boxed{p(\mu_i) = 0} \text{ for all } i.$$



Are the low degree polys  $p$  s.t.

$P(M_i) = 0$  for all  $i$ ?

Yes



Linear transformation  $\begin{matrix} \left\{ \begin{matrix} \deg \leq d \\ \text{polys} \end{matrix} \right\} \end{matrix} \rightarrow \begin{matrix} \binom{n+d}{d} \\ \mathbb{R}^k \end{matrix}$

$(P(M_1), P(M_2), \dots, P(M_n))$

$\dim k$

Kernel of  $J$ .

Kernel has  $\dim \leq k$ .

$\boxed{\binom{n+d}{d}} \mathbb{R}^{n'}$

$\binom{n'+d}{d} > k$  if  $n' > d + k$ .

Choose  $n$  rows of  $X \rightarrow$  (can  $D$ )  $\rightarrow$  Choose space of low degree  
 Vanishing Polynomials  $\mathcal{P}$   
 (exactly).

$V =$  Variety defined by these polynomials

$=$  Set of pts for which these polys vanish.

$M_i \in V$ .

intuition  $V$  small

Proof  $\dim(\text{deg-}d \text{ polys on } V) \leq \text{Codim of space of defining polys.}$   
 $\leq R$ .

$\left( \begin{matrix} \dim(V) + d \\ d \end{matrix} \right)$

$\Rightarrow \underline{\dim(V) \leq d \cdot R/d.}$

$\square$

Idea "project" onto  $V$ .

do something exhaustive within  $\exp(dn(U))$

$$d = \lg(K)$$

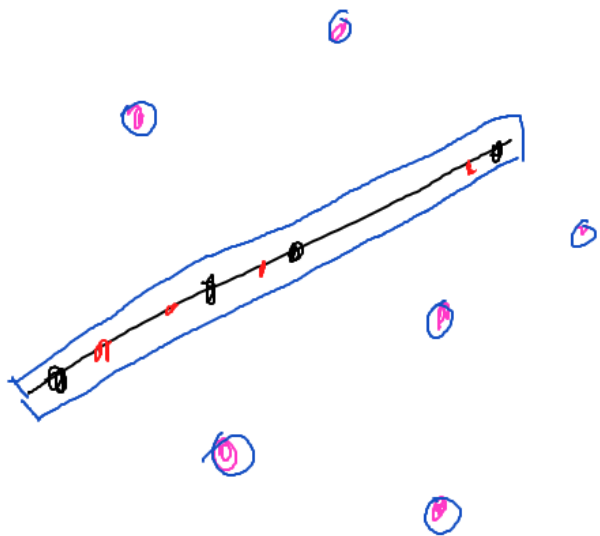
Complexity of answer data

$(nd)$  time.

reduce to  $d^{1/d} \approx d$ . dimensionality varies

exhaustive time taken  $\approx 2^d$  time.

quasi-poly alg



## Technical problems

- ①  $V$  is low dimension but is it "simple"?
- ② How do we project onto  $V$ ?
- ③ How do we compute on  $V$ ?
- ④ How do I deal with approx. defining polynomials?

Idea ~~is~~ replace  $V$  w/ a point cloud.

~~Set~~ Vector space  $U$  of degree  $d$  polys  $\mathbb{R}$

$$\text{Codim}(U) = k.$$

$$S = \{ x \in \mathbb{R}^n : |x| < R, |p(x)| < \underbrace{\delta}_{\text{we'll get back to this}} |p| \text{ for all } p \in U \}$$

Want for  $\epsilon > 0$  to find a small  $\epsilon$ -cover of  $S$ .

~~Thm~~ Thm If  $\epsilon > \frac{1}{\delta^{2d}} \text{poly}(R \text{ and } k)$  then  
there exists an  $\epsilon$ -cover of size at most  
 $(2(R/\epsilon) \cdot dk) O(d^2 k^4)$

expect

lower side should be  $\geq$

$O(R/\epsilon)$   $\dim(V) \leftarrow$

$d_k' \epsilon$

we get

$d^2 k' \epsilon$

$$\boxed{\epsilon > \delta^{1/2} (\dots)}$$

$\epsilon < \delta$

# Technicalities

①  $U \subset$  non. degree-d polynomials.

②  $\delta < (\epsilon/R)^{1/d}$

$$P(x) = A \cdot x^{\otimes d}$$

$$|P| := |A|_2$$

$|P|$  is spherically symmetric.

---

Define  $f(\underline{\epsilon}, \underline{R}, \underline{d}, \underline{k}, \underline{n}, \underline{\delta}) =$  ~~the~~ least case case size.

Prove a recursive upper bound on  $f$ .

$$\mathbb{R}^n = \mathbb{R}^{n'} \times \mathbb{R}^{n-n'}$$

$\downarrow$                        $\downarrow$   
 $(x, y)$

$n'$  as smallish

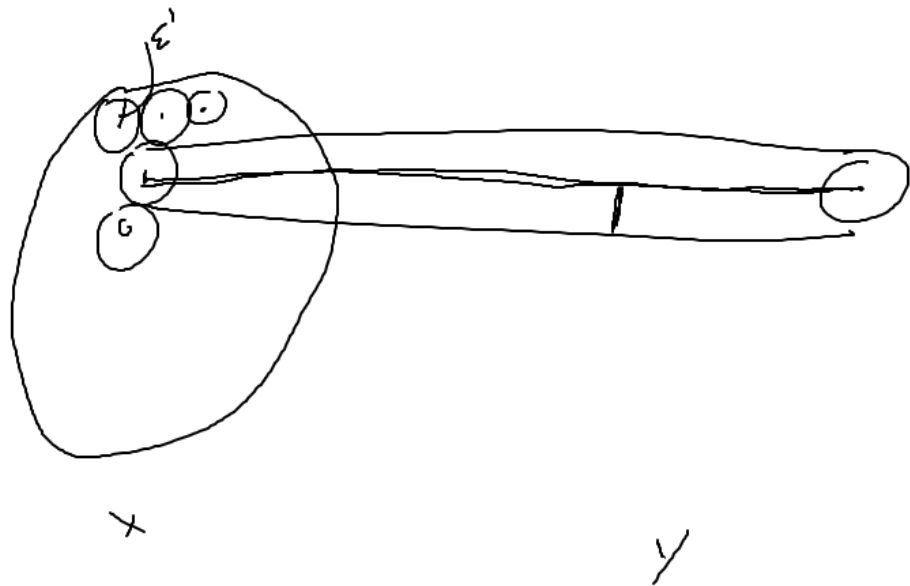
$W = \{ \text{polys in } U \text{ that are degree } \leq d \text{ in } x \text{ \& degree } \leq 1 \text{ in } y \}$

$$\text{Codim}(W) \leq k.$$

if  $U$  is  $k$ -dim in  $V$ .

then  $U \cap W$  is  $k$ -dim in  $U \cap W$ .





$\epsilon'$  very small

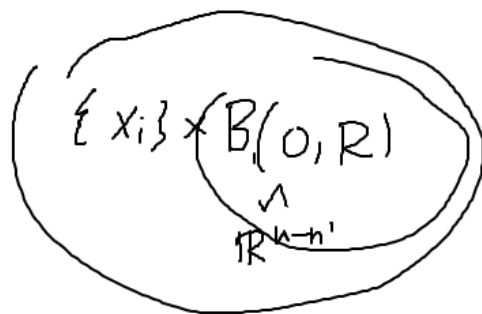
Small enough

If I change  $\delta$  by a bit...

It is enough to cover

this center line.  $\{x_i\} \times \mathbb{R}^{n-1}$

$\{x_i\}$



$A_{x_i} : W \rightarrow \{ \text{new deg } d-1 \text{ polys, in } y \}$

$P(x, y) \rightarrow P(x_i, y)$

$d$   
deg<sup>1</sup>, deg<sup>d-1</sup>  
polys, in  $U$ .

$S_{x_i} = \{ \text{new vanishing pts w/ } x = x_i \}$

$\eta$  not too small

$n' + (-1)$

If  $q = A_{x_i} P$

~~if~~  $|q| > \eta |p|$

then  $S_{x_i}$  must nearly vanish on  $q$ .

$(x_i, y)$   
nearly vanishes

$U_{x_i} = \text{Span of the singular vectors of } A_{x_i} \text{ w/ singular value } \geq \eta$

$\delta/|p| = (\delta/\eta) |q|$

$S_{x_i} \subset \{ \text{pts that nearly vanish on } U_{x_i} \}$

recursive version of Division Problem -

Say  $x_i$  is good if  $\text{Codim}(U_{x_i}) \leq R'$

If  $x_i$  is good then the cylinder can be covered by a set of size.

$$f(d-1, R', n-n', R, \tilde{\epsilon}, \tilde{\delta}/\eta)$$

What about the rest of good pts

---

Then I can cover all of the cylinders using all of the good

$$\delta \approx \epsilon^d \dots$$

$$\eta \approx \text{poly}(\epsilon)$$

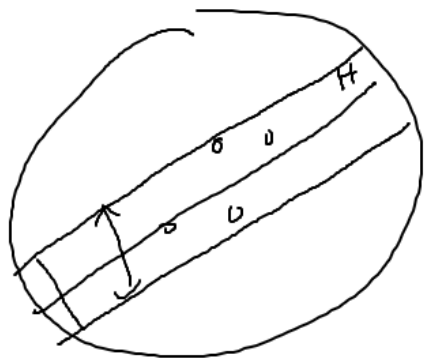
$$f(\dots) \leq (\dots)^{d^2} R^{kd}$$

$$R' \approx R^{\frac{d-1}{d}}$$

points w/ an appropriate # of covers.

Prop There exists a subspace  $H \subset \mathbb{R}^{n'}$   
 of dimension  $\leq 2k/k'$

s.t. all of the bad pts are  
close to  $H$



Cover  $H \times \mathbb{R}^{k-n'}$  as many as  
 $n' \gg k/k' \gg k^{1/2}$ .  
 Subspace of dim  $\mathbb{R}^{k-n'+2k/k'}$ . Then I  
 should be fine.

+ f(d, R,  $\tilde{\epsilon}$ ,  $\delta$ ,  $\underline{n-n'+2k/k'}$ )

How do I prove this prop.?

Contraction

If not.

Sequence of  
 $t \geq \mathbb{R}/\mathbb{R}'$

bad pts  
s.t. each

$x_1, x_2, x_3, \dots, x_t$

$x_i$  is far from

span  $x_1, \dots, x_{i-1}$

---

for each  $x_i$  I have

$p_{i1}, p_{i2}, \dots, p_{ik}$  ← sum of singular vectors of  $A_{x_i}$

$$p_{ij} \cdot (A_{x_i} \begin{matrix} \uparrow \\ w \end{matrix}) \leq \eta \|w\|_2$$

$$\|p_{ij}\| = 1$$

$$B_{i,j} = q_i(x) \cdot P_{i,j}(y)$$

$$q_i(x) = x \cdot x_i$$

Gram matrix determinant of the  $B_{i,j}$ 's

$$[\langle B_{i,j}, B_{i',j'} \rangle]$$

reasonably big

list  $B_{i,j}$ 's in order

$B_{1,1} \ B_{1,2} \ \dots \ B_{1,k} \ B_{2,1} \ \dots$

then each

$B_{i,j}$  has a reasonably large component orthogonal to previous ones.

$B_{i,j}$  vs

$B_{i',j}$

$$P_{i,j} \perp P_{i',j'}$$

$X_i$  had a component of size  $\geq \gamma$  orthogonal to

$X_1, \dots, X_{i-1}$

$$\langle p e_j, p' q' \rangle \approx \langle p, p' \rangle \langle q, q' \rangle$$

$B_{ij}$  has a comp of size  $\approx \gamma$  orthogonal to any

$q_j'(x) p_j(y)$  for  $j' < i$

$B_{ij}$  has a  $\gamma$ -sized comp. orthogonal to previous  $B_{ij}$ 's.

$\Rightarrow$  if we apply G-S. to  $B_{ij}$ 's we get  
comps of size  $\geq \gamma$ .

$\Rightarrow \det(\text{Gram Matrix}) \geq \gamma \cdot \prod_{i,j} (2R_{ij}) \cdot \prod_{i,j} z_k$

# Upper Bound

Claim

~~B.T~~

If  $P \in W$  then.

$\langle P(x,y), B_{ij} \rangle$  small.

$$B_{ij} = q_i(x) r_j(y)$$

~~≠~~  $P$  is not on basis for  $x$ .

$$P(x,y) = \underbrace{\sum q_i(x)}_{\text{orthogonal!}} P_i(y) + \underbrace{\dots}_{\text{orthogonal!}}$$

$$\begin{aligned} \langle P(x,y), B_{ij} \rangle &\approx \langle P(x_i, y), \underbrace{P_{ij}(y)}_{\text{small}} \rangle \\ &\quad \parallel \\ &\quad A_{x_i} P. \\ &\leq \nu \cdot |P|. \end{aligned}$$

Sum of errors of  $A_{x_i}$



~~3~~

$$M = \begin{bmatrix} | \\ B_{ij} \\ | \end{bmatrix}$$

We note  $M$  is  $2k$ -dimensional.

$W$  has  $\text{Codim} \leq k$ .

$\Rightarrow \geq k$ -dim'l subspace  $a$ . s.t.

Not hard to show that  $M^T M$  has no superfluous evs.

$$M a \in W.$$

For such

$\alpha$ 's.

$$M^T M a = M^T (M a).$$

$$\det(M^T M)$$

$$|M^T M a| \leq 2k \eta |a|.$$

$$= \underbrace{[ \langle B_{ij}, M a \rangle ]}_{\leq \eta \cdot |M a|}$$

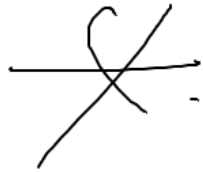
$$\leq (2k \eta)^k.$$

$\Rightarrow M^T M$  has at least  $k$  eigenvalues  $\leq 2k \eta$

$$\leq \eta \cdot |M a|$$
  
$$\uparrow$$
  
$$2k \cdot |a|$$

$$(\gamma)^{2k} < \det(M^T M) < (2k\eta)^k.$$

as long as  $\eta$  is small enough.



Note: This can be made computationally efficient.

---

Learning mixtures of Spherical Gaussians in TV

- ① reduce to  $k$  dimensions.  
distance.
- ① rough clustering - reduces to case where  
 $|M_i| = \text{Poly}(k)$
- ② Compute 1'st ~~2d~~ <sup>2d</sup> moments to error  $(\epsilon/k)^{O(d)}$   
parameter (moments of  $\mu_i$ ).

④ Define  $Q(p) \stackrel{\text{deg } d \text{ poly.}}{\approx} \sum w_i P^2(\mu_i)$

Let  $U$  be space of small singular vectors of  $Q$ .

for any  $i$  w/  $w_i$  not too small and any  $P \in U$   
 $|P(\mu_i)|$  small.

⑤ Use Thm compute a cone  $S$ .

every  $i$  with  $w_i \geq \epsilon/k$  has  $e_i$  within  $\epsilon$  of  
some elt of  $S$ .

$$|S| = (\epsilon/k)^{\alpha(d^2 k^{1/2})}.$$

$$X = \sum w_i N(\mu_i, \Sigma).$$

$$\approx \sum w_i N(\mu_i', \Sigma)$$

↑  
S.

Explicit set of distributions  $X_1, \dots, X_{|S|}$ .

$X$  is  $\epsilon$ -close to a mixture of the  $X_i$ 's

Claim Learn  $X$  given samples from  $X$  and  $X_i$ 's.

Point is compute MLE.  $\rightarrow$  EP. Convex program.

$$X + \varepsilon \left( \frac{x_1 + \dots + x_n}{n} \right)$$

Runtime

$$N = (k/\varepsilon) O(d^2 k^{1/d})$$

processes required

$$(k/\varepsilon) O(d)$$

Samples

$$d = \lg(k)$$

runtime

$$(k/\varepsilon)^{O(\lg^2(k))}$$

TV  
distance  
lower  
no  
separation.

Existing algorithms

run

If  $|u_i - u_j| \gg \sqrt{\lg(k)}$   
in quasi poly time samples learn ~~the~~ parameters to error  $\varepsilon$ .

If  $|M_i - M_j| \gg \frac{1}{\sqrt{\log(k)}}$

and if you have  $n$  samples (of size  $N$ ) means you can learn the mean for the set.

Then in  $\text{poly}(N/k)$  time you can learn the mean.  $\mathcal{O}(\text{poly}(n/\epsilon))$  samples.

~~$\mathcal{O}(k)$~~   $(k/\epsilon)^{\mathcal{O}(\log n)}$  samples  $(k/\epsilon)^{\mathcal{O}(\log^2 k)}$  time.

$$d = \log k$$

$|M_i - M_j| \gg k^{1/d}$   $\leftarrow$  precision we have a trade off.  $(k/\epsilon)^{\mathcal{O}(d)}$  samples  $(k/\epsilon)^{\mathcal{O}(d^2 k^{1/d})}$  time.

① Preprocessing  $R$  dim  $|M|$ 's not too big.

② Computed parameter moments

③ Compute a lower

④ use brute force.



# Mixture of linear regressions

$$X \sim N(0, I)$$

$$Y = \beta \cdot X + \text{noise.}$$

$\beta_1 \dots \beta_k$

weights

$w_1, \dots, w_k$

w/ prob  $w_i$

$$Y = \beta_i \cdot X + \text{noise.}$$

~~compute~~

penalty terms of  $\beta_i$ 's

$$\sum w_i \beta_i^{\otimes d}$$

# Non-negative Linear Combinations of ReLUs

$$\text{ReLU}(x) = \max(0, x)$$

$$\underline{F}(x) = \sum a_i \text{ReLU}(x \cdot \underline{w}_i) \quad (\underline{a_i \geq 0})$$

$(x, F(x))$   
+ noise

approximate  $F$ .

$$\sum a_i \underline{v}_i \otimes d.$$

$$\hookrightarrow \underline{\sum a_i \rho^2(v_i)}$$

$$\underline{\chi_i^2 = 1.}$$

$$\mathbb{R}^n = \mathbb{R}^{n'} \times \mathbb{R}^{n-n'}$$

$\downarrow$                        $\downarrow$   
 $(x, y)$

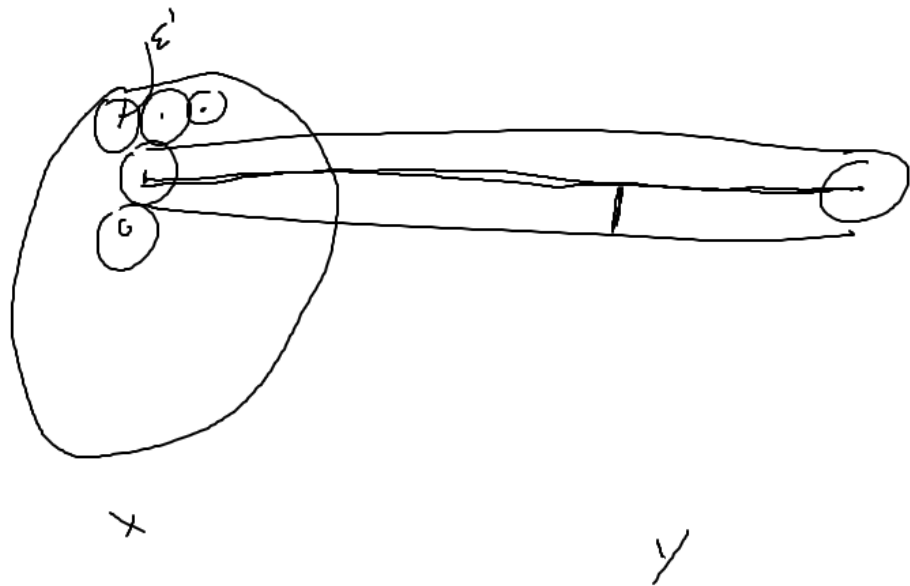
$n'$  as smallish

$W = \{ \text{polys in } U \text{ that are degree } \leq d \text{ in } x \text{ \& degree } \leq d-1 \text{ in } y \}$

$$\text{Codim}(W) \leq k.$$

if  $U$  is Codim  $k$  in  $V$ .

then  $U \cap W$  is Codim  $\leq k$  in  $V \cap W$ .



$\epsilon'$  very small

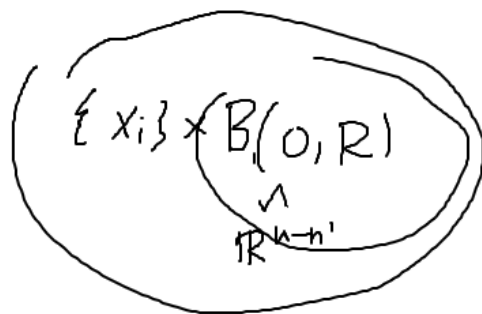
Small enough

If I change  $\delta$  by a bit...

It is enough to cover

this center line.  $\{x_i\} \times \mathbb{R}^{n-1}$

$\{x_i\}$



$A_{x_i} : W \rightarrow \{ \text{new deg } d-1 \text{ polys, in } y \}$

$P(x, y) \rightarrow P(x_i, y)$

$d$   
deg', deg  $d-1$   
polys, in  $U$ .

$S_{x_i} = \{ \text{new vanishing pts w/ } x = x_i \}$

$\eta$  not too small

$n' + \epsilon - 1$

If  $q = A_{x_i} P$

~~if~~  $|q| > \eta |p|$

then  $S_{x_i}$  must nearly vanish on  $q$ .

$(x_i, y)$   
nearly vanishes

$U_{x_i} = \text{Span of the singular vectors of } A_{x_i} \text{ w/ singular value } \geq \eta$

$\delta |p| = (\delta/\eta) |q|$

$S_{x_i} \subset \{ \text{pts that nearly vanish on } U_{x_i} \}$

recursive version of Divisor problem -

Say  $x_i$  is good if  $\text{Codim}(U_{x_i}) \leq R'$

If  $x_i$  is good then the cylinder can be covered by a set of size.

$$f(d-1, R', n-n', R, \tilde{\epsilon}, \tilde{\delta}/\eta)$$

What about the rest of good pts

---

Then I can cover all of the cylinders using all of the good

$$\delta \approx \epsilon^d \dots$$

$$\eta \approx \text{poly}(\epsilon)$$

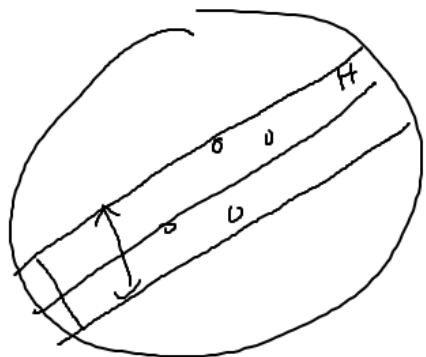
$$f(\dots) \leq (\dots)^{d^2} R^{kd}$$

$$R' \approx R^{\frac{d-1}{d}}$$

points w/ an appropriate # of covers.

Prop There exists a subspace  $H \subset \mathbb{R}^{n'}$   
 of dimension  $\leq 2k/k'$

s.t. all of the bad pts are  
close to  $H$



Cover  $H \times \mathbb{R}^{k-n'}$  as many as  
 $n' \gg k/k' \gg k^{1/2}$ .  
 Subspace of dim  $\mathbb{R}^{k-n' + 2k/k'}$ .  
 Then I should be fine.

+ f(d, R,  $\tilde{\epsilon}$ ,  $\delta$ ,  $\underline{n-n'+2k/k'}$ )