

Moreover, the algorithm works in higher dimensions too (still poly-time)

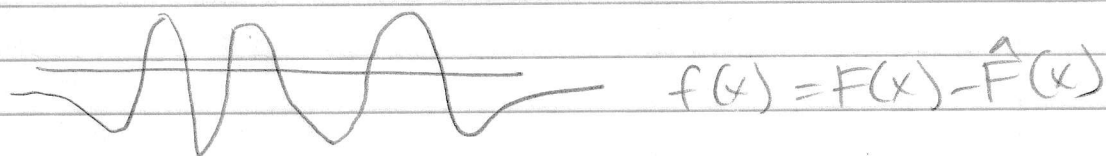
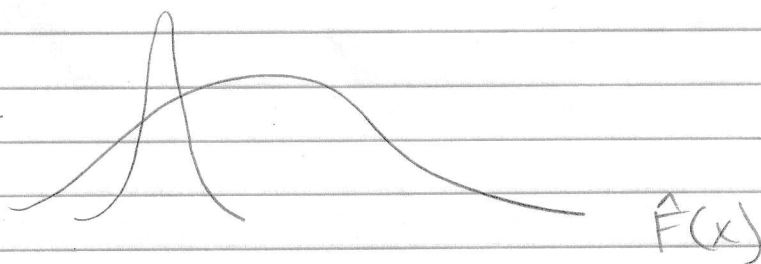
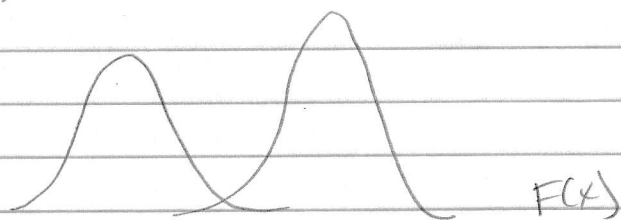
### Six Moments Suffice

What if we are given the first six moments exactly? Does this determine the parameters?

Alternatively:

Do any two different mixtures  $F$  and  $\hat{F}$  necessarily differ on at least one of their first six moments?

Consider the difference between their pdfs



Lemma: If  $f(x)$  has at most six zero crossings, then  $F(x)$  and  $\hat{F}(x)$  differ on one of their first six moments

Proof: we can find a degree  $\leq 6$  polynomial  $p(x)$  that agrees in sign with  $f(x)$ .

$$\text{Then } 0 < \left| \int p(x) f(x) dx \right| = \left| \int \sum_{r=1}^6 p_r x^r f(x) dx \right|$$

$$\leq \sum_{r=1}^6 |p_r| \left| \int x^r f(x) dx \right|$$

$$= \sum_{r=1}^6 |p_r| |M_r(F) - M_r(\hat{F})|$$

So  $\exists r \in \{1, \dots, 6\}$  s.t.  $|M_r(F) - M_r(\hat{F})| > 0$   $\square$

What remains is to show:

Proposition If  $f(x) = \sum_{i=1}^k \alpha_i \mathcal{N}(\mu_i, \sigma_i^2, x)$  is not  
positive/negative

identically zero, then  $f(x)$  has at most  $2k-2$  zero crossings

Can check with  $k=4 \Rightarrow$  six moments suffice

we'll use the heat equation:

Question: If the initial heat distribution on a 1-d infinite rod ( $x$ ) is  $f(x) = f(x, 0)$ , what is the heat distribution at time  $t$ ?

Probabilistic interpretation ( $\sigma^2 = 2kt$ )

$$f(x, t) = \mathbb{E} [f(x+z, 0)]$$

$z \sim \mathcal{N}(0, \sigma^2)$

Alternatively this is a convolution

$$f(x, t) = \int_{-\infty}^{\infty} f(x+z) \mathcal{N}(0, \sigma^2, z) dz$$

$$\stackrel{\Delta}{=} f(x) * \mathcal{N}(0, \sigma^2, x)$$

Theorem [Hummel, Gidas] Suppose  $f(x): \mathbb{R} \rightarrow \mathbb{R}$  is analytic and has  $N$  zeros. Then

$$f(x) * \mathcal{N}(0, \sigma^2, x)$$

has at most  $N$  zeros (for any  $\sigma^2 > 0$ )

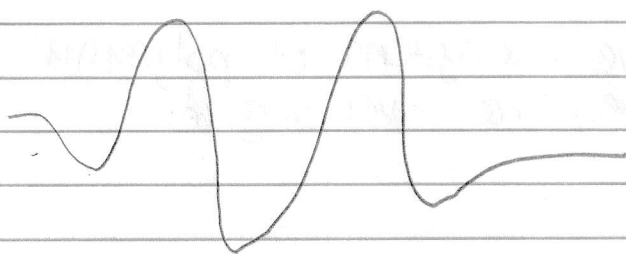
i.e. convolving by a Gaussian / running the heat equation does not increase # of zeros

Last ingredient:

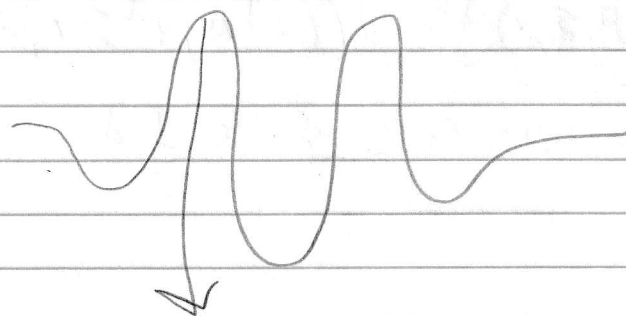
Fact:  $\mathcal{N}(\mu_1, \sigma_1^2, x) * \mathcal{N}(\mu_2, \sigma_2^2, x) =$

$$\mathcal{N}(\mu_1 + \mu_2, \sigma_1^2 + \sigma_2^2, x)$$

has at most 4 zc's. Hence



when we add back in the delta function,  
we add at most two new zc's,



and convolving by  $\mathcal{N}(0, \sigma_{\min}^2, x)$  gets us  
back to the original linear combination,  
but does not increase the number of zc's.  $\square$

### Sampling Noise

What if we only have estimates of the  
first six moments?

definition Let  $\Theta$  be the set of valid parameters,  
i.e.  $w_i \in [0, 1]$  and  $\sigma_i^2 \geq 0$

what we just showed is:

$$\{ \hat{\theta} \in \Theta \mid M_r(\hat{\theta}) = M_r(\theta) \text{ for } r=1 \text{ to } 6 \}$$

the only solutions are  $\theta = (w_1, \mu_1, \sigma_1^2, w_2, \mu_2, \sigma_2^2)$   
and the relabeling  $\theta' = (1-w_1, \mu_2, \sigma_2^2, w_1, \mu_1, \sigma_1^2)$

Are these equations stable, when given noisy estimates?

Again using deconvolution, can show

Proposition  $\exists$  constants  $c, C$  s.t. if  $\epsilon < c$ , the means and variances  $\leq \frac{1}{\epsilon}$ , and mixing weights are in  $[\epsilon, 1-\epsilon]$  and

$$|M_r(\theta) - M_r(\hat{\theta})| \leq \epsilon^C \quad r=1 \text{ to } 6$$

then there is a permutation  $\pi$  s.t.

$$\sum_{i=1}^2 |w_i - \hat{w}_{\pi(i)}| + |\mu_i - \hat{\mu}_{\pi(i)}| + |\sigma_i^2 - \hat{\sigma}_{\pi(i)}^2| \leq \epsilon$$

Hence close enough estimates for the first six moments guarantee close parameters too

This is called polynomial identifiability

reduction from high d to 1-d

A View from Algebraic Geometry

Following Belkin and Sinha, we'll give a more general framework

definition: A class of distributions  $F(\theta)$  is a polynomial family if  $\forall r$

$$M_r(\theta) \stackrel{\Delta}{=} \mathbb{E}[x^r]$$

$x \sim F(\theta)$

is a polynomial in  $\theta = (\theta_1, \dots, \theta_k)$

e.g. GMMs, mixtures of uniform, exponential, Poisson, gamma

definition: the moment generating function (mgf) is defined as

$$f(t) = \sum_{n=0}^{\infty} \mathbb{E}[x^n] \frac{t^n}{n!}$$

Fact: If the mgf converges in a neighborhood of zero then

$$\{M_r(\theta) = M_r(\hat{\theta}) \forall r\} \Rightarrow F(\theta) = F(\hat{\theta})$$

i.e. the infinite sequence of moments determines the density function

Now we'll need some notions/tools from algebraic geometry

definition Given a ring  $R$ , an ideal  $I$  generated by  $g_1, \dots, g_n \in R$  is

$$\langle g_1, \dots, g_n \rangle \stackrel{\Delta}{=} I = \left\{ \sum r_i g_i \mid \forall r_i \in R \right\}$$

Moreover:

definition: A Noetherian ring is a ring s.t.  
for any sequence of ideals

$$I_1 \subseteq I_2 \subseteq \dots$$

$$\exists N \text{ st. } I_N = I_{N+1} = I_{N+2} = \dots$$

And our main tool:

Theorem [Hilbert's Basis theorem] If  $R$  is  
a Noetherian ring, then  $R[x]$  is too.

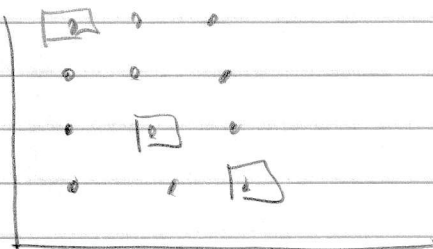
↑  
polynomials w/ coeffs in  $R$

Let's pause for some intuition

Lemma [Dickson] For any subset  $K \subseteq \mathbb{N} \times \mathbb{N}$   
there are a finite number of minimal elements

i.e.  $(a,b) \in K$  s.t.  $\nexists (a',b') \in K$   
with  $a' \leq a, b' \leq b$  and one  
inequality is strict

we can visualize this



containing only  
If  $I$  is an ideal generated by polynomials of the form  
 $x^a y^b$

then we know  $I$  is Noetherian, and hence has a finite list of generators. This is what Dickson's lemma tells us too

Warning: An ideal  $I$  over the ring  $\mathbb{R}[x, y]$  is Noetherian, but there is no effective bound of how many generators we need

Theorem [Belkin, Sinha] Let  $F(\theta)$  be a polynomial family. If the mgf converges in a neighborhood of zero then  $\exists N$  s.t.

$$F(\theta) = F(\hat{\theta}) \iff M_r(\theta) = M_r(\hat{\theta}) \text{ for } r=1 \text{ to } N$$

Proof: Let  $Q_r(\theta, \hat{\theta}) = M_r(\theta) - M_r(\hat{\theta})$

Now let

$$I_1 = \langle Q_1 \rangle, I_2 = \langle Q_1, Q_2 \rangle, \text{ etc.}$$

By Hilbert's basis theorem,  $\exists N$  s.t.

$$I_N = I_{N+j} \quad \forall j \geq 0$$

thus we have

$$Q_{N+j}(\theta, \hat{\theta}) = \sum_{i=1}^N P_{ij}(\theta, \hat{\theta}) Q_i(\theta, \hat{\theta})$$



for some polynomials  $p_{ij} \in \mathbb{R}[\theta, \hat{\theta}]$ . Thus if

$$M_r(\theta) = M_r(\hat{\theta}) \text{ for } r=1 \text{ to } N$$

we have  $\varphi_r(\theta, \hat{\theta}) = 0$  for  $r=1$  to  $N$  which implies

$$\varphi_{N+1}(\theta, \hat{\theta}) = 0 \text{ too!}$$

Now from the fact, we get that  $F(\theta) = F(\hat{\theta})$   $\square$

So we know finitely many moments suffice,  
but don't have an effective bound.

Belkin and Sinha also gave a stability analysis  
through quantifier elimination.

def: A set  $S$  is semialgebraic if  $\exists$  polynomials s.t.

$$S = \{x_1, \dots, x_d \mid \bigvee_{i=1}^N p_i(x_1, \dots, x_d) = 0\}$$

or if  $S$  is the finite union or intersection of such sets

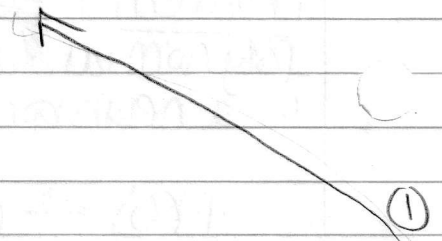
we can define the projection i.e.

$$T = \{x_1, \dots, x_{d-1} \mid \exists x_d \text{ s.t. } (x_1, \dots, x_d) \in S\}$$

Theorem: [Tarski] the projection of a semialgebraic set is also semi-algebraic.

Notice, if you only use polynomial equations,  
you might need polynomial inequalities for the proj.

Fact: the complement  $\bar{S}$  of a semialgebraic set is also semialgebraic.



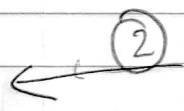
Corollary: The set  $U = \{x_1, \dots, x_{d+1} \mid \exists x_d (x_1, \dots, x_d) \in S\}$  is also semialgebraic.

Proof: We can write

$$\bar{U} = \{x_1, \dots, x_{d+1} \mid \exists x_d (x_1, \dots, x_d) \notin S\}$$

Thus we have

$$S \xrightarrow{\text{Fact}} \bar{S} \xrightarrow{\text{Tarski}} \bar{U} \xrightarrow{\text{Fact}} U$$



are all semialgebraic  $\square$

Now let's prove stability via quantifier elimination

Consider the set  $H(\epsilon, \delta)$  defined to be the set of all  $\epsilon$  and  $\delta$  s.t.

$$\left. \begin{aligned} \forall \theta, \hat{\theta} \quad |M_r(\theta) - M_r(\hat{\theta})| \leq \delta \text{ for } r=1 \text{ to } 6 \end{aligned} \right\} (P)$$

$$\Downarrow$$

$$d_p(\theta, \hat{\theta}) \leq \epsilon$$

Claim:  $H(\epsilon, \delta)$  is a semialgebraic set

Proof: Consider the tuple  $(\epsilon, \delta, \theta, \hat{\theta})$  that satisfy the predicate  $(P)$ . Equivalently

$$\{ |M_r(\theta) - M_r(\hat{\theta})| \leq \delta \text{ for } r=1 \text{ to } 6 \} \cup \{ d_p(\theta, \hat{\theta}) \leq \epsilon \}$$

This is semialgebraic

Now applying quantifier elimination to get rid of the  $\forall$  completes the proof.  $\square$

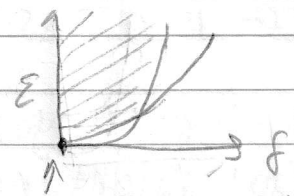
Theorem [Belkin, Sinha]  $\exists$  constants  $C_1, C_2$  s.t. <sup>and's</sup>  
 if  $\delta \leq 1/C_1$ , the  $\epsilon(\delta) \leq C_2 \delta^{1/5}$   
 (sketch)

Proof: Since  $H(\epsilon, \delta)$  is semialgebraic,

$$\epsilon^*(\delta) = \text{smallest } \epsilon \text{ s.t. } (\epsilon, \delta) \in H(\epsilon, \delta)$$

$$\epsilon^*(\delta) \geq \text{poly}(\delta)$$

for sufficiently small  $\delta$ .  $\square$



six moments suffice!

need to show  $\epsilon^*(\delta) > 0$   $\forall \delta > 0$

Proof by wiggling

# Robust Statistics

Basic estimation problem, but we'll go in a new direction:

Given samples from a 1-d Gaussian  $\mathcal{N}(\mu, \sigma^2)$ , can we estimate its parameters?

Of course! Use:

$$\hat{\mu} = \frac{1}{N} \sum_{i=1}^N x_i, \quad \hat{\sigma}^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \hat{\mu})^2$$

These are examples of maximum likelihood paradigm (Ronald Fisher 1912-1922)

(1) consistent: converges to true parameters as  $N \rightarrow \infty$  under tame conditions

(2) asymptotically normal: has smallest variance among all unbiased estimators

Main Question: But what if the samples are only approximately Gaussian?

definition: In the strong contamination model:

(1)  $m$  samples are drawn iid from  $P \in \mathcal{D}$

known class of distributions