

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?

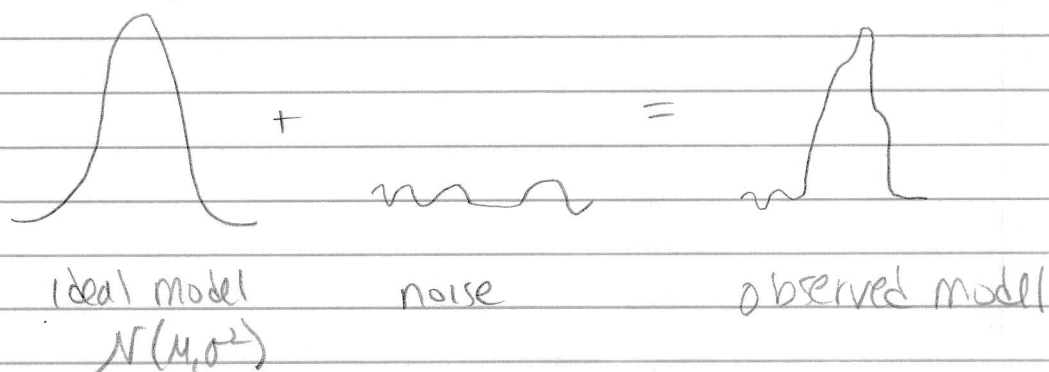
definition: In the strong contamination model:

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

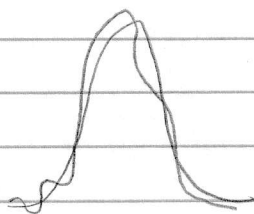
known class of distributions

(2) adversary is allowed to arbitrarily corrupt an ϵ -fraction of samples

Pictorially:



we can think of the area between the curves



as representing the samples the adversary has added/deleted

definition: the total variation distance between r.v.s. with pdfs $f(x)$ and $g(x)$ is

$$d_{TV}(f, g) = \frac{1}{2} \int |f(x) - g(x)| dx$$

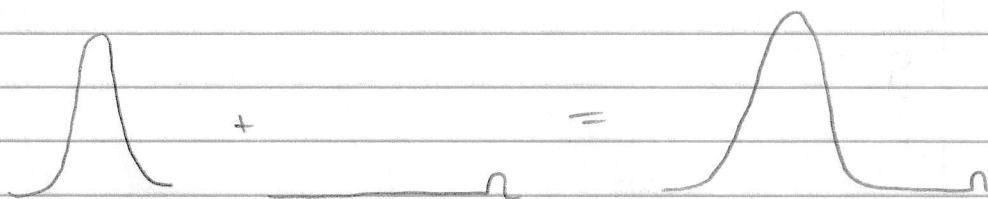
In our model, we have

$$d_{TV}(\underbrace{\text{smooth curve}}_{\text{ideal}}, \underbrace{\text{jagged curve}}_{\text{observed}}) \leq O(\epsilon)$$

Can we estimate the true Gaussian in $O(\epsilon)$ in TV?

(MLE)
Observation: The empirical mean/variance do not work!

e.g.



but as the bump $\rightarrow \infty$, $\hat{\mu}$ and $\hat{\sigma}$ diverge

So what should we do? Consider

$$\hat{\mu} = \text{median}(\{x_i | \xi_i\})$$

$$\text{MAD} = \text{median}(\{|x_i - \hat{\mu}| | \xi_i\})$$

$$\hat{\sigma} = \frac{\text{MAD}}{\Phi^{-1}(3/4)}$$

\nwarrow
c.d.f of a standard Gaussian

Proposition [folklore] Given ϵ -corrupted samples from a 1-D Gaussian $\mathcal{N}(\mu, \sigma^2)$ we have

$$d_{\text{TV}}(\mathcal{N}(\hat{\mu}, \hat{\sigma}^2), \mathcal{N}(\mu, \sigma^2)) \leq O(\epsilon)$$

provided $m \geq \frac{c \ln^{1/8}}{\epsilon^2} \nwarrow$ failure prob.

In the nomenclature of TCS

"properly agnostically learning a 1-d Gaussian"

↖
output something
from the class

↑
it's not actually
Gaussian, but want
to do well if it's close

Main Question: what about in high-dimensions?

Given ϵ -corrupted samples from a d -dimensional Gaussian $\mathcal{N}(\mu, \Sigma)$, can we efficiently estimate st.

$$d_{\text{TV}}(\mathcal{N}(\hat{\mu}, \hat{\Sigma}), \mathcal{N}(\mu, \Sigma)) \leq \tilde{O}(\epsilon)?$$

Special cases:

① unknown mean: $\mathcal{N}(\mu, \mathbf{I})$

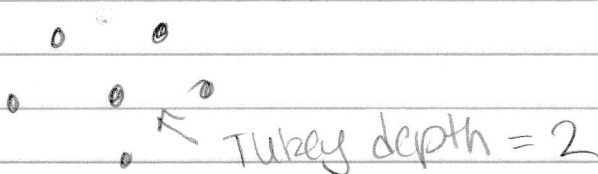
② unknown covariance: $\mathcal{N}(0, \Sigma)$

What's known in robust statistics?

def: The Tukey depth of a point x w.r.t. a dataset x_1, \dots, x_m is

$$\min_{1\text{-d proj.}} \min(\# \text{ points to left / right of } x)$$

e.g.

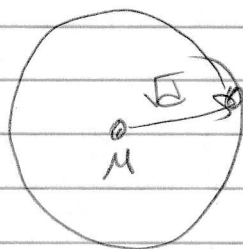


Fact: Given ϵ -corrupted samples from $\mathcal{N}(\mu, \Sigma)$ the Tukey median (Tukey deepest point over all space) satisfies $d_{TV}(\mathcal{N}(\hat{\mu}, \hat{\Sigma}), \mathcal{N}(\mu, \Sigma)) \leq O(\epsilon)$

Unfortunately:

Lemma: the Tukey Median is NP-hard to compute.

Alternatively we could take coordinatewise median, but that would only get $TV \leq \epsilon \sqrt{d}$



Because direction of corruption might not be axis aligned

Theorem: [Diakonikolas, Li, Kamath, Kane, Moitra, Stewart]
There is a polynomial time/sample complexity algorithm that finds $\hat{\mu}, \hat{\Sigma}$ satisfying

$$d_{TV}(\mathcal{N}(\hat{\mu}, \hat{\Sigma}), \mathcal{N}(\mu, \Sigma)) \leq O(\epsilon \log^{3/2} 1/\epsilon)$$

in the ϵ -strong contamination model

[Lai, Rao, Vempala] also gave an algorithm satisfying

$$TV \leq O(\sqrt{\epsilon} \log d)$$

when the covariance is bounded

General Recipe

- ① Find an appropriate parameter distance
- ② Detect when naive estimator has been compromised via method of moments
- ③ win-win: Find good parameters, or make progress by filtering out corruptions

Unknown Mean

Consider the special case when we get ϵ -corrupted samples from $\mathcal{N}(\mu, I)$

Definition: The KL-divergence is

$$d_{KL}(f \parallel g) = \int_{-\infty}^{\infty} f(x) \ln \frac{f(x)}{g(x)} dx$$

Fact: For two Gaussians, we have

$$d_{KL}(\mathcal{N}(\mu_1, \Sigma_1) \parallel \mathcal{N}(\mu_2, \Sigma_2)) =$$

$$\frac{1}{2} \left(\ln \frac{\det(\Sigma_2)}{\det(\Sigma_1)} + \text{Tr}(\Sigma_2^{-1} \Sigma_1) + (\mu_1 - \mu_2)^T \Sigma_2^{-1} (\mu_1 - \mu_2) - d \right)$$

when $\Sigma_1 = \Sigma_2 = I$ this simplifies to:

$$d_{KL}(\mathcal{N}(\hat{\mu}, I), \mathcal{N}(\mu, I)) = \frac{1}{2} \|\hat{\mu} - \mu\|_2^2$$

Fact [Pinsker's Inequality]

$$d_{\text{TV}}(f, g)^2 \leq \frac{1}{2} d_{\text{KL}}(f, g)$$

Putting it all together, we have

Lemma: If we can estimate

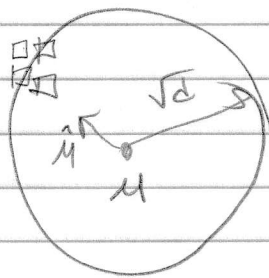
$$\|\hat{\mu} - \mu\|_2 \leq \tilde{O}(\epsilon)$$

then we'd get $d_{\text{TV}}(\mathcal{N}(\hat{\mu}, \mathbb{I}), \mathcal{N}(\mu, \mathbb{I})) \leq \delta(\epsilon)$

Thus we have our parameter distance!

Detecting Corruptions

How can the adversary move the empirical mean by $\epsilon \sqrt{d}$?



But in this case, the projected variance on the direction $\hat{\mu} - \mu \gg 1$

Takeaway: To mess up the first moment, an adversary would have to mess up the second moment too

Key Lemma: If X_1, \dots, X_m are ϵ -corrupted samples from $\mathcal{N}(\mu, \Sigma)$ and

$$(1) m \geq c \frac{d \ln 1/\delta}{\epsilon^2}$$

$$(2) \hat{\mu} = \frac{1}{m} \sum X_i, \hat{\Sigma} = \frac{1}{m} \sum (X_i - \hat{\mu})(X_i - \hat{\mu})^T$$

then we have w/ probability $\geq 1 - \delta$

$$\|\mu - \hat{\mu}\|_2 \geq C' \epsilon \sqrt{\log 1/\delta} \Rightarrow \|\hat{\Sigma} - \Sigma\|_2 \geq C'' \epsilon \log 1/\delta$$

Thus we can detect when the empirical mean has been corrupted

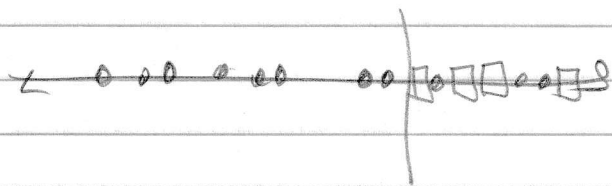
Spectral Filtering

If $\|\hat{\Sigma} - \Sigma\|_2 < C'' \epsilon \log 1/\delta$ then we can just output $\hat{\mu}$ and are guaranteed

$$\|\mu - \hat{\mu}\|_2 \leq C' \epsilon \sqrt{\log 1/\delta} \Rightarrow$$

$$d_{TV}(\mathcal{N}(\hat{\mu}, \Sigma), \mathcal{N}(\mu, \Sigma)) \leq O(\epsilon \sqrt{\log 1/\delta})$$

Otherwise consider $v =$ direction of largest variance



we can compute a threshold T s.t. throwing

out the samples above T results in throwing out more corrupted than uncorrupted points

Unknown Covariance

Using the formula for KL-divergence and Pinsker's inequality it can be shown

Fact: For two Gaussians, $\mathcal{N}(0, \Sigma)$ and $\mathcal{N}(0, \hat{\Sigma})$

$$d_{TV}(\mathcal{N}(0, \hat{\Sigma}), \mathcal{N}(0, \Sigma)) \leq O\left(\|I - \hat{\Sigma}^{-1/2} \Sigma \hat{\Sigma}^{-1/2}\|_F\right)$$

↑

Mahalanobis distance

Can we use the fourth moment to detect corruptions in the second moment?

Lemma: Let $X \sim \mathcal{N}(0, \Sigma)$. Then consider

$$M = \mathbb{E}[(X \otimes X)(X \otimes X)^T]$$

restricted to flattenings of symmetric $d \times d$ matrices we have

$$M = 2\Sigma^{\otimes 2} + (\Sigma^b)(\Sigma^b)^T$$

Now imagine we get ε -corrupted samples X_1, \dots, X_m from $\mathcal{N}(0, \Sigma)$. Then define

$$\hat{\Sigma} = \frac{1}{m} \sum X_i X_i^T \text{ and}$$

$$Y_i = (\hat{\Sigma})^{-1/2} X_i$$

If $\Sigma = \hat{\Sigma}$ then $Y_i \sim \mathcal{N}(0, I)$ in which case
restricted to subspace of symmetric matrices

$$F = \frac{1}{m} \sum (Y_i \otimes Y_i) (Y_i \otimes Y_i)^T \sim \underbrace{2I}_{d \times d} + (I^b)(I^b)^T$$

If we consider

$$\max_{z \text{ s.t.}} z^T F z \sim 2$$

$z =$ flattening of
 $d \times d$ trace zero,
symmetric matrix

We show that if $\|I - \hat{\Sigma}^{-1/2} \Sigma \hat{\Sigma}^{-1/2}\|_F \gg \epsilon \log 1/\epsilon$
then we must have

$$\max_z z^T (F - 2I) z \gg \epsilon \log 1/\epsilon$$

$z =$ flattening...

And similarly if the generalized eigenvalues
of F are small, we can output $\hat{\Sigma}$

Otherwise we can find a direction (degree two polynomial)
to filter on

Assembling the Algorithm

Given ϵ -corrupted samples from $\mathcal{N}(\mu, \Sigma)$

① Doubling trick

$$x_i - x_i' \sim_{\delta(\epsilon)} \mathcal{N}(0, 2\Sigma)$$

Use algorithm for unknown covariance

② Agnostic isotropic position

$$\hat{\Sigma}^{-1/2} x_i \sim_{\delta(\epsilon)} \mathcal{N}(\hat{\Sigma}^{-1/2} \mu, \mathbf{I})$$

Use algorithm for unknown mean

More Robust Statistics

Many subsequent directions

① Handling more errors ($\epsilon > 1/2$) with list decoding

② Giving evidence of lower bounds, e.g. Statistical query algorithms can't get $O(\epsilon)$ error

③ weakening the distributional assumptions, just need bounds on moments

④ exploiting sparsity e.g. $\|u\|_0 \leq k$

⑤ more complex generative models

Let's return to GMMs:

Earlier we gave a polynomial time/sample complexity non-robust learner. Implicitly we should:

definition we say a family of distributions \mathcal{D} is polynomially identifiable if

$\forall P_1, P_2 \in \mathcal{D}$ that have ϵ -different parameters

$$\Rightarrow d_{TV}(P_1, P_2) \geq \text{poly}(\epsilon, \frac{1}{d})$$

[Kalai, Mottra, Valiant]

Corr: Mixtures of two Gaussians in d -dimensions are polynomially identifiable

Otherwise the algorithm wouldn't work

But to get robust algorithms we need a much stronger notion

definition: we say a family of distributions \mathcal{D} is robustly identifiable if

$\forall P_1, P_2 \in \mathcal{D}$ that have parameter discrepancy ϵ

$$\Leftrightarrow d_{TV}(P_1, P_2) \geq \text{poly}(\epsilon)$$

Notice that for unknown mean/covariance

we could take ℓ_2 / Mahalanobis distance in parameters

Theorem [Liu, Maity] Mixtures of two Gaussians in d -dimensions are robustly identified by a constant number of Hermite moments

The proof will be via generating functions and differential operators

Generating Functions

Consider a Gaussian $G = \mathcal{N}(\mu, I + \Sigma)$ and let

$$M(x) = X^T \mu \quad \text{and} \quad \Sigma(x) = X^T \Sigma X$$

vector of formal variables

Key Lemma

$$e^{M(x)y + \frac{1}{2}\Sigma(x)y^2} = \sum_{m=0}^{\infty} \frac{1}{m!} \mathbb{E}_{Z \sim G} [H_m(Z, X)] y^m$$

where $H_m(Z, X) =$ Hermite moment tensor, i.e.

$$H_m(Z, X=v) = \begin{matrix} m^{\text{th}} \text{ Hermite moment} \\ \text{of } 1\text{-d Gaussian r.v.} \\ v^T Z \end{matrix}$$

Let's make it even simpler to see why

Easier lemma. Let $G = \mathcal{N}(\mu, I + \sigma^2)$. Then

$$e^{\mu y + \frac{\sigma^2}{2} y^2} = \sum_{m=0}^{\infty} \frac{1}{m!} \mathbb{E}[h_m(z)] y^m$$

Proof: Expanding the LHS we get

$$1 + (\mu y + \frac{\sigma^2}{2} y^2) + \frac{(\mu y + \frac{\sigma^2}{2} y^2)^2}{2} + \frac{(\mu y + \frac{\sigma^2}{2} y^2)^3}{6} + \dots$$

Collecting terms

$$1 + \mu y + \left(\frac{\mu^2 + \sigma^2}{2}\right) y^2 + \left(\frac{\mu^3 + 3\mu\sigma^2}{6}\right) y^3 + \dots$$

$$\mathbb{E}[h_2(z)] \stackrel{\Delta}{=} \mathbb{E}[z^2 - 1] = \mu^2 + 1 + \sigma^2 - 1$$

$$\mathbb{E}[h_3(z)] \stackrel{\Delta}{=} \mathbb{E}[z^3 - 3z] = \mu^3 + 3\mu(1 + \sigma^2) - 3\mu$$

etc. ~~...~~

Now the Key Lemma extends immediately to mixtures

$$\mathcal{M} = w_1 \mathcal{N}(\mu_1, I + \Sigma_1) + \dots + w_k \mathcal{N}(\mu_k, I + \Sigma_k)$$

then we have:

Key Lemma:

$$\sum_{j=1}^k w_j e^{\mu_j(x)y + \frac{1}{2} \Sigma_j(x)y^2} = \sum_{m=0}^{\infty} \frac{1}{m!} \frac{\mathbb{E}[H_m(Z, x)]}{z^m} y^m$$

Back to Identifiability

We WTS that

$$\textcircled{1} \sum_{j=1}^k w_j e^{\mu_j(x)y + \frac{1}{2} \Sigma_j(x)y^2} = \sum_{j=1}^k \hat{w}_j e^{\hat{\mu}_j(x)y + \frac{1}{2} \hat{\Sigma}_j(x)y^2}$$



② the parameters match, i.e. the mixtures are the same on a component-by-component basis

Now consider the differential operator

$$\mathcal{D} = d_y - (\mu + \sigma^2 y)$$

applied to the generating function $e^{\mu y + \frac{1}{2} \sigma^2 y^2}$

Fact: $\mathcal{D}(e^{\mu y + \frac{1}{2} \sigma^2 y^2}) = 0$

This holds as a formal identity

Observation \mathcal{D} is a polynomial rearrangement on the series expansion

This works in high dimensions too. Let

$$\mathcal{D} \triangleq d_y - (\mu(x) - \Sigma(x)y)$$

then:

Fact: $\mathcal{D}(e^{\mu(x)y + \frac{1}{2}\Sigma(x)y^2}) = 0$

which is a complicated, but explicit multivariate polynomial rearrangement

Main Question: what happens when you apply \mathcal{D} to another component?

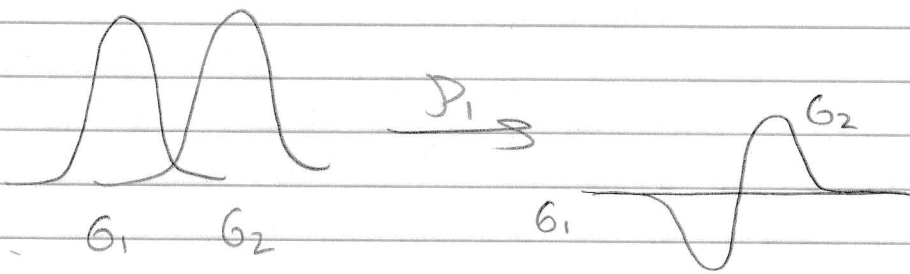
Fact: $\mathcal{D}(e^{\hat{\mu}(x)y + \frac{1}{2}\hat{\Sigma}(x)y^2}) = P e^{\hat{\mu}(x)y + \frac{1}{2}\hat{\Sigma}(x)y^2}$

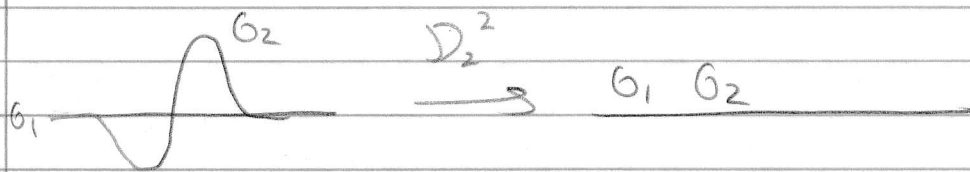
where $P = (\hat{\mu}(x) - \mu(x) + y(\hat{\Sigma}(x) - \Sigma(x)))$

And finally

Fact: $\mathcal{D}(P e^{\mu(x)y + \frac{1}{2}\Sigma(x)y^2}) = \frac{dP}{dy} e^{\mu(x)y + \frac{1}{2}\Sigma(x)y^2}$

Now we can use differential operators to isolate a component





Thus consider

$$\begin{array}{ccc}
 M & \begin{array}{c} 2k-2 \\ D_{k+1}^2 \end{array} & \text{vs.} \\
 \downarrow & D_{k+1}^{2k} & \hat{M} \\
 P(x, y) e^{M_1(x)y + \frac{1}{2} \sum_j (x) y^2} & & \hat{D}_1 \\
 & & \downarrow \\
 & & O
 \end{array}$$

This implies one of their $f(x)$ degree moments is different!

Can show this argument gets robust identifiability!

Theorem [Liu, Mottra] There is a polynomial time algorithm for robustly learning a GMM with accuracy that depends polynomially on the corruption rate

[Bakshi et al] get related results, but for weaker notion of proper density estimation

[Liu, Mottra] get $\tilde{O}(\epsilon)$ accuracy for semi-proper density estimation