

Partial Specification and Agnostic Clustering.

Last time:

- Partial specification for linear regression
 - Non-robust method assuming Gaussianity
 - Robust method based on robust standard errors

This time.

- Clustering

Setting.

Distributions p_1, \dots, p_k

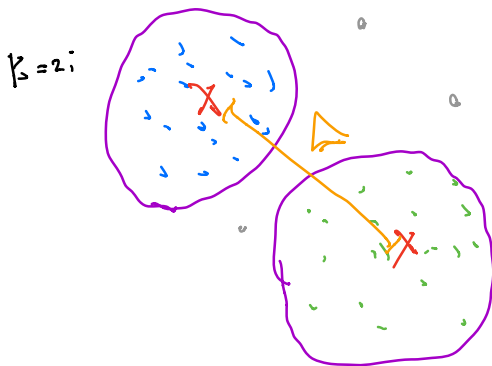
Points x_1, \dots, x_n sampled from $p = \alpha_1 p_1 + \alpha_2 p_2 + \dots + \alpha_k p_k$.

A priori bound: $\alpha_{\min} = \min_j \alpha_j$.

Two goals.

- Parameter recovery: estimate parameter e.g. $\mu_j = \mathbb{E}_{p_j}[x]$ of each p_j

- Cluster recovery: figure out which cluster p_j each point x_i came from.



Sometimes have separation condition Δ

Settings.

- Parametric: each p_j has known parametric form e.g. $N(\mu_j, \Sigma_j)$
 - Agnostic: no parametric form, only assume e.g. sub-Gaussianity, bounded moments, etc.
 - Robust: also allow ϵ -fraction of outliers, i.e. $\alpha_1 + \dots + \alpha_K = 1 - \epsilon$
-

Warm-up. Parametric setting

$$p_j = N(\mu_j, \Sigma_j)$$

$$p = \sum_j \alpha_j p_j = \sum_j \alpha_j N(\mu_j, \Sigma_j)$$

Key questions.

- Identifiability: are (μ_j, Σ_j) unique?
↳ If not, can't recover (μ_j, Σ_j) even when $n = \infty$.
- Statistical rate: assuming identifiable, how many samples do we need? $\rightarrow \frac{1}{\sqrt{n}}$ rate w/ MLE
complicated but "standard"

Only up to permutation

$$\alpha_1 \leftrightarrow \alpha_2$$

$$\mu_1 \leftrightarrow \mu_2$$

$$\Sigma_1 \leftrightarrow \Sigma_2$$

Proposition. As long as (μ_j, Σ_j) are all distinct,
 then $(\alpha_j, \mu_j, \Sigma_j)_{j=1}^k$ identifiable up to permutation.

Pf. Suppose not, then

$$\sum_j \alpha_j N(\mu_j, \Sigma_j) + \sum_j (-\alpha_j') N(\mu_j', \Sigma_j') = 0$$

\Rightarrow linear dependence among PDFs

ID. $\sum_{j=1}^m c_j \exp\left(-\frac{(x-\mu_j)^2}{\sigma_j^2}\right) / \sqrt{2\pi\sigma_j^2} = 0 \quad \forall x$

\Rightarrow Take MGF

$$\sum_{j=1}^m c_j \exp\left(\frac{1}{2} \sigma_j^2 \lambda^2 + \mu_j \lambda\right) = 0 \quad \forall \lambda$$

As $\lambda \rightarrow \infty$, only largest σ_j ,
 call it σ_{\max} , matters.
 \rightarrow Divide through by σ_{\max}

$$\Rightarrow \sum_{j=1}^m c_j \exp(\mu_j \lambda) \rightarrow 0 \quad \text{as } \lambda \rightarrow \infty$$

$$\text{as } \lambda \rightarrow -\infty$$

$$\mu_j > 0$$

and $\mu_j < 0 \Rightarrow$ contradiction.

\Rightarrow No ID linear relations.

n-D. Take random projection of n-D Gaussians
 Linear relation in n-D \Rightarrow relation in 1-D
 $\Rightarrow \Leftarrow$

$$\sum_{j=1}^m c_j N(\mu_j, \Sigma_j) = 0$$

Random projection onto v

$$\Rightarrow \sum_{j=1}^m c_j N(v^T \mu_j, v^T \Sigma_j v) = 0$$

Impossible unless parameters are
not actually distinct.

probability zero

$$v^T \mu_{j'} = v^T \mu_j$$

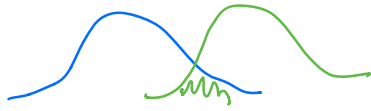
$$v^T \Sigma_{j'} v = v^T \Sigma_j v \quad \text{for } j \neq j'$$

Cluster recovery: need separation



hard to cluster in overlap

- Overlap btw $N(\mu, \Sigma)$ and $N(\mu', \Sigma')$?
- Simplify: $\Sigma = \Sigma' = \sigma^2 \cdot I$



okay if $k \cdot \Phi\left(\frac{\|\mu - \mu'\|_2}{\sigma}\right) \leq \epsilon$

$$\Delta = \min_j \|\mu_j - \mu_{j'}\|_2$$

$$\Rightarrow \text{Want } k \cdot \Phi\left(\frac{\Delta}{\sigma}\right) \leq \epsilon$$

$$\Delta \gg \sigma \sqrt{\log(k/\epsilon)}$$

Partial specification

- Just assume sub-Gaussian, resilience

- Take $k=1$, $\epsilon = 1-\alpha$ 10% good points
 $\alpha = 0.1$ 90% outliers

- List-decoding. Output $m = O(1/\alpha)$ "candidates"
 $\hat{\mu}_l$ s.t. $\|\mu - \hat{\mu}_l\|_2$ is small for some l .



Can actually handle $k > 1$
clusters automatically via
list-decoding.



Assumption, Set S of an "good" points
w/ mean μ , which is $(\rho, \frac{\alpha}{4})$ -resilient

Proposition, Given assumption, can output

$m \leq \frac{2}{\alpha}$ candidates $\hat{\mu}_1, \dots, \hat{\mu}_m$ s.t.

$$\|\hat{\mu}_j - \mu\|_2 \leq \frac{8\rho}{\alpha} \text{ for some } j.$$

Interpretation, sub-Gaussian $\Rightarrow (\sigma \alpha \sqrt{\log(1/\alpha)}, \alpha)$ -resilient

$$\rho = \mathcal{O}(\sigma \alpha \sqrt{\log(1/\alpha)})$$

$$\frac{8\rho}{\alpha} = \mathcal{O}(\sigma \sqrt{\log(1/\alpha)})$$

k clusters
each w/ $\alpha = \frac{1}{k}$
 $\Rightarrow \mathcal{O}(\sigma \sqrt{\log(k)})$

\Rightarrow good cluster recovery
if $\Delta \gg \sqrt{\log(k)}$.

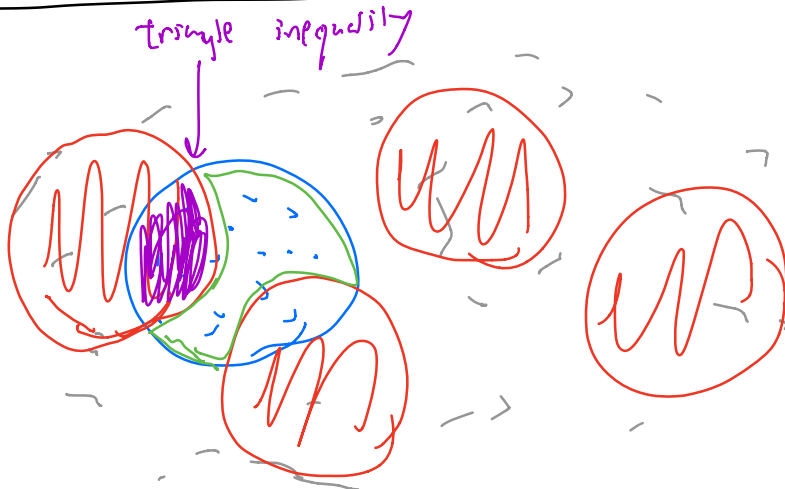
bounded covariance

$\Rightarrow (\sigma\sqrt{\alpha}, \alpha)$ -resilient

$$\frac{\delta f}{\alpha} = O(\sigma/\sqrt{\alpha}). \Rightarrow$$

$$\sigma\sqrt{k}$$

\uparrow Threshold
at which MLE
can be done
"easily" efficiently



Idea.
• Cover set of all points w/ resilient subsets

How to pick red sets: maximal collection of sets S_j
s.t.:

- $|S_j| \geq \frac{\alpha}{2} n$
 - S_j is $(\frac{4}{\alpha} p, 1 - \frac{\alpha}{2})$ -resilient
 - $S_j \cap S_{j'} = \emptyset \quad \forall j \neq j'$
- $\rightarrow \left(\frac{\frac{\alpha}{2} \cdot \frac{4}{\alpha} p}{1 - \frac{\alpha}{2}}, \frac{\alpha}{2} \right)$
 $\approx (2p, \frac{\alpha}{2})$ -resilient
- Lemma p is (p, ϵ) -resilient
 if and only if it is

$\hat{\mu}_j = \text{mean of } S_j$

Output $\hat{\mu}_1, \dots, \hat{\mu}_m$

$(\frac{1-\epsilon}{\epsilon} p, 1-\epsilon)$ -resilient.
 $\tau' \epsilon + (1-\epsilon)\tau = 0 \Rightarrow \tau' = -\frac{(1-\epsilon)\tau}{\epsilon}$

τ $1-\epsilon$

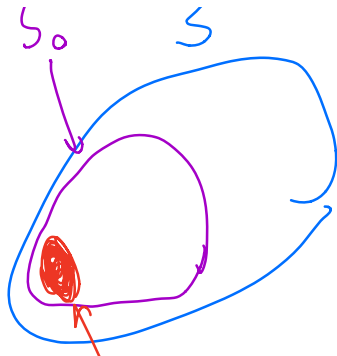
Claim Some S_j intersects $S \Rightarrow \hat{\mu}_j \approx \mu$

$$S_0 = S \setminus (S_1 \cup \dots \cup S_m)$$

① $|S_0| \geq \frac{\alpha}{2} n$. Then S_0 can be added to S_j .

S is $(p, \frac{\alpha}{4})$ -resilient

$\Rightarrow (\frac{4}{\alpha} p, 1 - \frac{\alpha}{4})$ -resilient (by Lemma)



\Rightarrow Any subset of $\frac{\alpha}{4}|S|$ points in S has mean at most $\frac{4}{\alpha}p$ away,

\Rightarrow Any subset of $\frac{\alpha}{2}|S_0|$ points in S_0 is a subset of $\frac{\alpha}{4}|S|$ points in $S_0 \subseteq S$.

$\frac{\alpha}{2}$ -fraction of S_0
 \Rightarrow $\frac{\alpha}{4}$ -fraction of S
 Can't be too far away.

(because $|S_0| \geq \frac{1}{2}|S|$).

$\Rightarrow S_0$ is $(\frac{4}{\alpha}p, 1 - \frac{\alpha}{2})$ -resilient.

$\Rightarrow S_0$ could have been added, $\rightarrow \mathbb{F}$.

(2) $|S_0| < \frac{\alpha}{2}n$. $|S| = \alpha n$

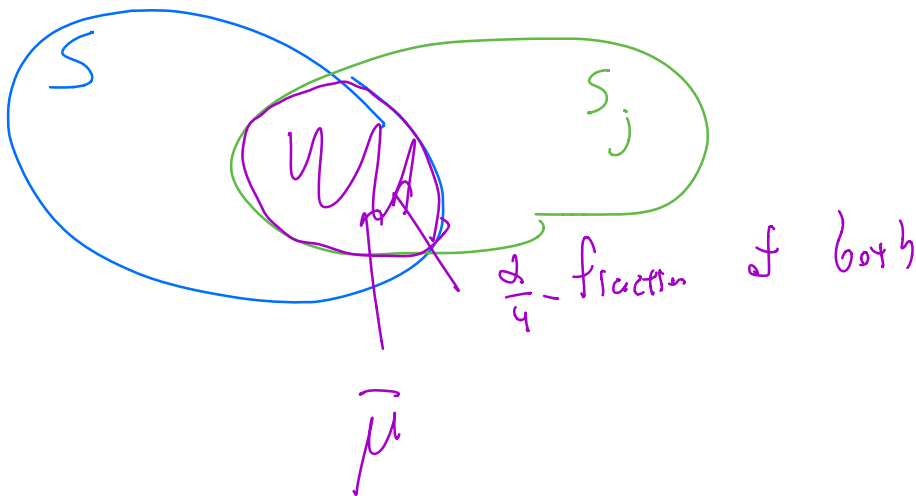
$$S: \underbrace{S \cap S_1, S \cap S_2, \dots, S \cap S_m, S_0}_{> \frac{\alpha}{2}n} \quad S_0 < \frac{\alpha}{2}n$$

$$\sum_j |S \cap S_j| \geq \frac{\alpha}{2}n$$

$$\sum_j |S_j| \leq n \quad (\text{by disjointness})$$

$$\Rightarrow |S \cap S_j| \geq \frac{\alpha}{2} |S_j| \quad \text{for some } j.$$

$$\geq \frac{\alpha}{4} |S| \quad \leftarrow \text{since } |S_j| \geq \frac{\alpha}{2} n \geq \frac{1}{2} |S|.$$



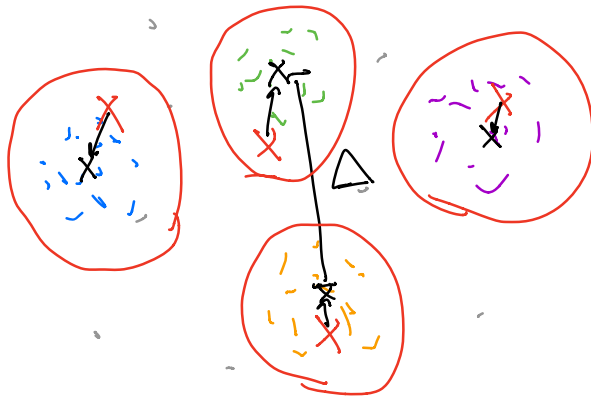
$$\| \bar{\mu} - \mu(S) \|_2 \leq \frac{4\rho}{\alpha} \quad (\text{by resilience w/ } \epsilon = 1 - \frac{\alpha}{4})$$

$$\| \bar{\mu} - \mu(S_j) \|_2 \leq \frac{4\rho}{\alpha}$$

$$\Rightarrow \| \mu(S) - \mu(S_j) \|_2 \leq \frac{8\rho}{\alpha}.$$

Some $\hat{\mu}_j$ within $\left(\frac{8\rho}{\alpha}\right)$ of μ , as claimed. \square

Better rescale assuming well-separated clusters.



Proposition. Suppose x_1, \dots, x_n can be partitioned into sets C_1, \dots, C_k of size $\alpha_1 n, \dots, \alpha_k n$, plus ϵn outliers (so $\alpha_1 + \dots + \alpha_k = 1 - \epsilon$).

Further suppose:

- $2\epsilon \leq \alpha_{\min} = \min_{j=1}^k \alpha_j$.

- The means are well-separated: $\Delta > \frac{cD}{\epsilon}$,

where $\Delta = \min_{j \neq j'} \|\mu_j - \mu_{j'}\|_2$.

- Each cluster is $(\rho, 2\epsilon/\alpha)$ -resilient.
 (ρ, ϵ) -resilient

∧

Then we can output clusters \hat{C}_j such that:

$$\textcircled{1} |C_j \Delta \hat{C}_j| \leq O\left(\frac{\epsilon}{\alpha}\right) |C_j| \quad (\text{cluster})$$

$$\textcircled{2} \|\mu(\hat{C}_j) - \mu(C_j)\|_2 \leq 2\rho. \quad (\text{parameter})$$

$$\rho = \sigma \sqrt{\log(K)}$$