

Beyond Clustering

What are the minimal assumptions we need to learn the parameters?

Assumption #1: $\epsilon \leq w \leq 1 - \epsilon$

If not, we might never get a sample from one of the components, even w/ $\frac{1}{\epsilon}$ samples

Assumption #2: $d_{TV}(F_1, F_2) \geq \epsilon$

If not, we might only get samples from the overlap region and we can't learn the mixing wt

Theorem [Kalai, Moitra, Valiant] There is a polynomial time/sample complexity algorithm for learning* a mixture of two Gaussians under Assumption 1 & 2 up to ϵ - i.e.

* output $\hat{F} = \hat{w}_1 \hat{F}_1 + \hat{w}_2 \hat{F}_2$ s.t. \exists a permutation $\pi: \{1, 2\} \rightarrow \{1, 2\}$ and for all $i \in \{1, 2\}$ we have

$$d_{TV}(F_i, \hat{F}_{\pi(i)}) \leq \epsilon \text{ and } |w_i - \hat{w}_{\pi(i)}| \leq \epsilon$$

Outline

(1) reduce the d -dimensional problem to a series of 2 -dimensional problems

(2) Give provable guarantees for Pearson's sixth moment test

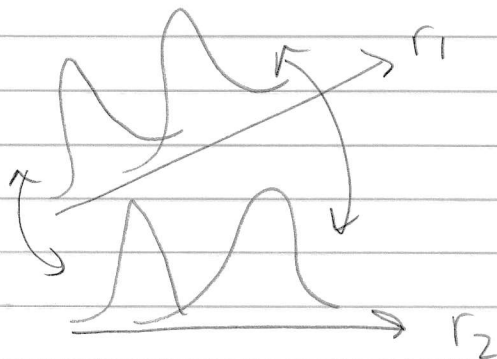
For (1) the key idea is

"learning the parameters of the projection gives linear constraints on the high-dim. parameters"

Fact: $\text{Proj}_r[F(x)] = w_1 \mathcal{N}(r^T \mu_1, r^T \Sigma_1 r)$
 $+ w_2 \mathcal{N}(r^T \mu_2, r^T \Sigma_2 r)$

i.e. each sample $x \sim F$ is mapped to $r^T x$

There are $\Theta(d^2)$ parameters so we need at least $\Theta(d^2)$ projections, but consider



Main Issue: How do you pair components up along different projections?

def. we say $F(x)$ is in isotropic position if

$$\textcircled{1} \mathbb{E}_{x \sim F} [x] = 0$$

$$\textcircled{2} \mathbb{E}_{x \sim F} [xx^T] = I$$

Lemma: If $\mathbb{E}_{x \sim F} [xx^T]$ is full rank, there

is an affine transformation that puts F in isotropic position

Proof: Let $\mu = \mathbb{E}_{x \sim F} [x]$ and consider

$$\mathbb{E}_{x \sim F} [(x-\mu)(x-\mu)^T] = M$$

Then by full rankness $M \succ 0$ and so by the Cholesky decomposition

$$M = BB^T$$

for invertible B . then set

$$y = B^{-1}(x-\mu) \quad \square$$

Now let's define a parameter distance in 1-d

$$d_p(N(\mu_1, \sigma_1^2), N(\mu_2, \sigma_2^2)) = |\mu_1 - \mu_2| + |\sigma_1^2 - \sigma_2^2|$$

Isotropic Projection Lemma: If F is in isotropic position and satisfies Assumption 1 & 2 then w.h.p for random r we have

$$d_p(\text{Proj}_r[F_1], \text{Proj}_r[F_2]) \leq \text{poly}\left(\frac{1}{d}, \epsilon\right) \triangleq \epsilon_3$$

Now if the directions are ϵ_2 -close, pairing is easy because the projected means/variances change by much less than ϵ_3 .

Main Issue: How do the errors in our estimates of the parameters of 1-d projections propagate to the high-dim estimates?

Condition Number Lemma: The condition number of the linear system

$$\left\{ \begin{array}{c} \left[\begin{array}{cc} r_i^T & 0 \\ 0 & \text{vec}(r_i r_i^T) \end{array} \right] \left[\begin{array}{c} \mu \\ \text{vec}(\Sigma) \end{array} \right] = \left[\begin{array}{c} \hat{\mu}_{r_i} \\ \hat{\Sigma}_{r_i} \end{array} \right] \end{array} \right\}_i$$

↑
unknowns

is at most $\text{poly}\left(n, \frac{1}{\epsilon_2}\right)$

Now we can estimate the parameters in 1-d up to accuracy $\epsilon_1 \ll \epsilon_2$ to get ϵ -error in high-dim