TWO-STEP EM Algorithm

Assumes K spherical gaussians, separation $||\mu_k^{true} - \mu_{k'}^{true} \ge C\sqrt{d\sigma_k}$

- 1. Pick $K' = \mathcal{O}(K \ln K)$ centers μ_k^0 at random from the data
- 2. Set $\sigma_k^0 = \frac{d}{2} \min_{k \neq k'} ||\mu_k^0 \mu_{k'}^0||^2$, $\pi_k^0 = 1/K'$
- 3. Run one E step and one M step $\implies \{\pi_k^1, \mu_k^1, \sigma_k^1\}_{k=1:K'}$
- 4. Compute "distances" $d(\mu_k^1, \mu_{k'}^1) = \frac{||\mu_k^1 \mu_{k'}^1||}{\sigma_k^1 \sigma_{k'}^1}$
- 5. Prune all clusters with $\pi_k^1 \leq 1/4K'$
- 6. Run Fastest First Traversal with distances $d(\mu_k^1, \mu_{k'}^1)$ to select K of the remaining centers. Set $\pi_k^1 = 1/K$.
- 7. Run one E step and one M step $\implies \{\pi_k^2, \mu_k^2, \sigma_k^2\}_{k=1:K}$

Theorem For any $\delta, \epsilon > 0$ if d large, n large enough, separation $C \ge d^{1/4}$ the Two step EM algorithm obtains centers μ_k so that

$$||\mu_k - \mu_k^{true}|| \le ||\text{mean}(C_k^{true}) - \mu_k^{true}|| + \epsilon \sigma_k \sqrt{d}$$