

Convergence rate of EM Algorithm on Gaussian Mixture Models

settings : under-specified , exactly-specified , over-specified. ✓.

over-specified : the Gaussian Mixture model we want to fit has more mixture components than the true mixture model generating the data.

It's most common to use over-specified mixture models.. So we want to understand how EM Algo. behaves in over-specified mixture model,

especially on Gaussian mixture. (Others, like mixture of regressions are discussed in other places.)

Problem set-up :

denote $\phi(\cdot, \mu, \Sigma)$ multivariate Gaussian distribution.

To make things easier, we focus on Gaussian Mixture with 2 components.

Say, true model is :

$$f(x; \theta^*, \sigma, \pi) := \pi \cdot \phi(x; \theta^*, \sigma^2 \text{Id}) + (1-\pi) \cdot \phi(x; -\theta^*, \sigma^2 \text{Id}) \quad (1).$$

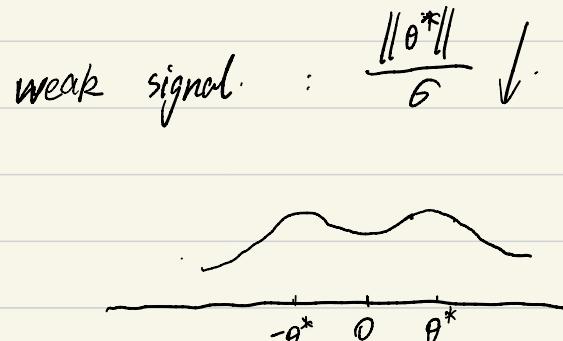
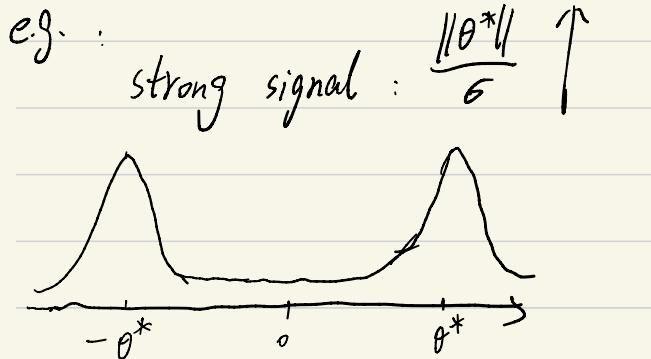
fixed, known = π, σ, Id .

Now we use EM to fit model

$$f(x; \theta, \sigma, \pi) := \pi \cdot \phi(x; \theta, \sigma^2 \text{Id}) + (1-\pi) \cdot \phi(x; -\theta, \sigma^2 \text{Id}). \quad (2)$$

and use solution $\hat{\theta}_n$ to estimate θ^*

"signal strength" : $\frac{\|\theta^*\|}{\sigma}$, i.e., the separation between the means of mixture components relative to the spread in the components.



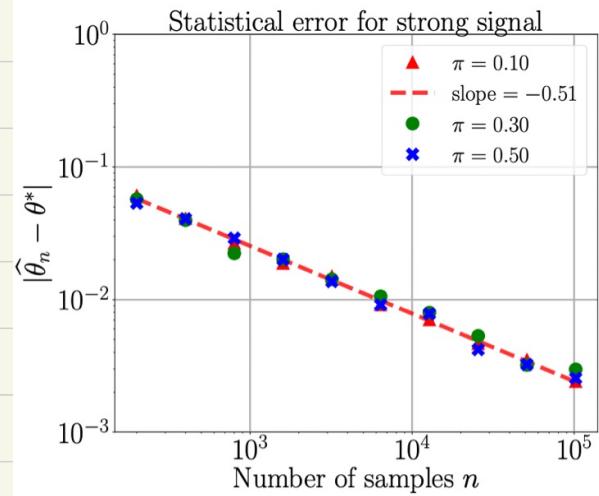
Particularly interested in the limit of weak signal case when there is no separation,

$$\text{i.e., } \|\theta^*\|_2 = 0.$$

Now the true model becomes a single Gaussian, $f(x) = \phi(x, 0, \sigma^2 \text{Id})$.

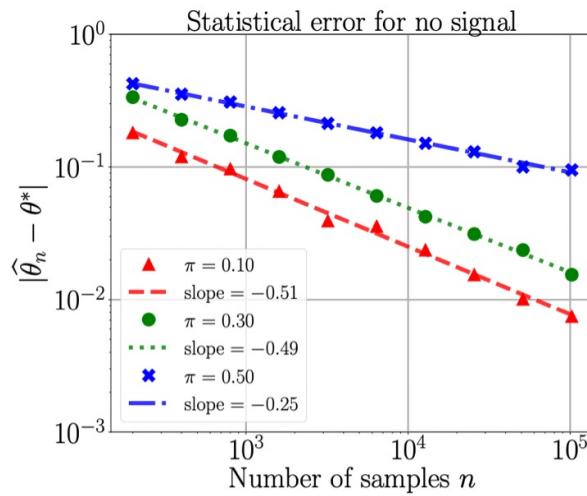
we are then in the setting of over-specified.

Why we interest in such limit weak signal, $\theta^* = 0$? See some simulations.



(a) $\theta^* = 5$

Figure 1.

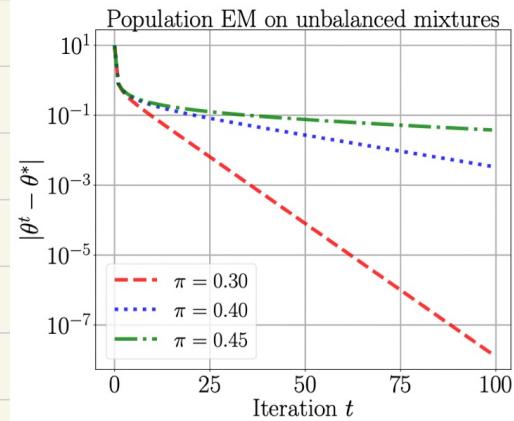


(b) $\theta^* = 0$

When true is $f(x) = \phi(x, 0, \sigma^2 \text{Id})$, but we fit $\pi \phi(x, \theta, \sigma^2 \text{Id}) + (1-\pi) \phi(x, -\theta, \sigma^2 \text{Id})$, different weights, error decays in different rate.

And with $\pi \uparrow 0.5$, i.e., more balanced, error decays more slowly.

In the over-specified setting with $\theta^* = 0$,
 weights affect EM convergence rate greatly..



(a)

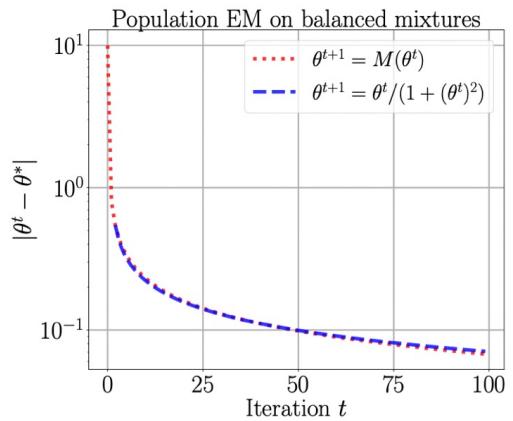


Figure 2.

(b)

(a) : unbalanced $\pi \neq \frac{1}{2}$: EM converges geometrically quickly , convergence rate $\sqrt{0}$ as $\pi \uparrow \frac{1}{2}$.

(b) balanced $\pi = \frac{1}{2}$: EM convergence rate is sub-geometric , much slower.

Schema:

While true data-generating model is a single Gaussian,

We analyze EM convergence rate of Gaussian Mixture with 2 components.

Unbalanced weights } { population-based : Thm 1(a)
sample-based : Thm 1(b).

balanced weights } { population-based : Thm 2
sample-based : Thm 3
+
Thm 4.

Population EM Algo. is the idealized version of EM based on an infinite sample size , the usual sample-based EM Algo is used in practice.

For theoretical purposes , we first analyze the convergence of the population EM updates , and then leverage these findings to understand the behavior of sample EM.

Population EM updates (§ 2.3).

true model: $\mathcal{N}(0, \sigma^2 I_d)$, $\theta^* = 0$.

we try to fit $\pi \phi(x, \theta, \sigma^2 I_d) + (1-\pi) \cdot \phi(x, -\theta, \sigma^2 I_d)$. (3).

latent variable Z : $P(Z=1) = \pi$, $P(Z=0) = 1-\pi$, $Z \sim \text{Bernoulli}(\pi)$.

joint (X, Z) : $X|Z=1 \sim \phi(x, \theta, \sigma^2 I_d)$, $X|Z=0 \sim \phi(x, -\theta, \sigma^2 I_d)$.

complete data: $(x_1, z_1), (x_2, z_2), \dots, (x_n, z_n)$, but we only have x_1, \dots, x_n .

population complete log-likelihood:

$$L(\theta'|X, Z) = E_X [\log P(X, Z|\theta')], \quad \text{where the expectation is taken over the true distribution } X \sim \mathcal{N}(0, \sigma^2 I_d).$$

$$= E_X \left\{ \log \left[(\pi \phi(\theta', \sigma^2 I_d))^Z ((1-\pi) \phi(-\theta', \sigma^2 I_d))^{1-Z} \right] \right\}$$

$$= E_X \left\{ Z \cdot \log(\pi \cdot \phi(\theta', \sigma^2 \text{Id})) + (1-Z) \cdot \log((1-\pi) \cdot \phi(-\theta', \sigma^2 \text{Id})) \right\}$$

Given θ ,

Take expectation of $\mathcal{L}(\theta' | X, Z)$ w.r.t. the posterior distribution of Z , i.e. $Z|X, \theta$, we get $Q(\theta', \theta)$, the lower bound of $\mathcal{L}(\theta' | X, Z)$.

$$\begin{aligned} Q(\theta', \theta) &= E_{Z|X, \theta} \mathcal{L}(\theta' | X, Z) \\ &= E_{Z|X, \theta} \left\{ E_X [\log p(X, Z | \theta')] \right\} \\ &= E_X \left\{ E_{Z|X, \theta} [\log p(X, Z | \theta')] \right\} \\ &= E_X \left\{ E_{Z|X, \theta} [Z \log(\pi \cdot \phi(\theta', \sigma^2 \text{Id})) + (1-Z) \cdot \log((1-\pi) \cdot \phi(-\theta', \sigma^2 \text{Id}))] \right\}, \end{aligned}$$

$$\text{denote } w_\theta(x) = E_{z|x,\theta}(z) = P(z=1|x,\theta) = \frac{\pi \cdot \exp\left(\frac{-\|x-\theta\|_2^2}{2\sigma^2}\right)}{\pi \cdot \exp\left(\frac{-\|x-\theta\|_2^2}{2\sigma^2}\right) + (1-\pi) \cdot \exp\left(\frac{-\|\theta+x\|_2^2}{2\sigma^2}\right)}$$

$$Q(\theta', \theta) = E_x \left[w_\theta(x) \left(\log \pi + \log \phi(\theta', \sigma^2 \text{Id}) \right) + (1-w_\theta(x)) \log (1-\pi) + (1-w_\theta(x)) \cdot \log \phi(-\theta', \sigma^2 \text{Id}) \right]$$

remove terms
unrelated to θ'

$$= -\frac{1}{2} E_x \left[w_\theta(x) \|x-\theta'\|_2^2 + (1-w_\theta(x)) \|x+\theta'\|_2^2 \right]. \quad (7)$$

$$\begin{aligned} \nabla_{\theta'} Q(\theta', \theta) &= -\frac{1}{2} E_x \left[w_\theta(x) (-X-X+2\theta') + (1-w_\theta(x)) (X+X+2\theta') \right] \\ &= E_x [(2w_\theta(x)-1)x] - \theta' = 0. \end{aligned}$$

so population EM operator:

$$M(\theta) = \arg \max_{\theta' \in \mathbb{R}^d} Q(\theta', \theta) = E_x [(2w_\theta(x)-1)x]. \quad (8)$$

population EM Algo : $\theta^{t+1} = M(\theta^t) = E_X \left[(2W_{\theta^t}(X) - 1) X \right], \quad t = 0, 1, 2, \dots$.

Analogously,

sample-based EM Algo : $\theta^{t+1} = M_n(\theta^t) = \frac{1}{n} \sum_{i=1}^n \left[(2W_{\theta^t}(X_i) - 1) X_i \right], \quad t = 0, 1, 2, \dots \quad (9).$

given the observed data $\{X_i\}_{i=1}^n$

Main results.

§ 3.1. Behavior of EM for unbalance mixtures.

$$\text{mixture model } \pi \cdot \phi(x, \theta, \sigma^2 \text{Id}) + (1-\pi) \cdot \phi(x, -\theta, \sigma^2 \text{Id}) \quad (3).$$

Thm 1: We fit an unbalanced ($\pi \neq \frac{1}{2}$) of mixture model (3) to $N(0, \sigma^2 \text{Id})$. $\theta^* = 0$.
Then

(a). The population EM operator (8), $M(\theta) = \mathbb{E}_x[(2w_\theta(x)-1)x]$, is globally strictly contractive, meaning that

population EM. (10a) $\|M(\theta)\|_2 \leq \left(1 - \frac{\rho^2}{2}\right) \|\theta\|_2 \text{ for all } \theta \in \mathbb{R}^d$.

$$\text{with } \rho := |1-2\pi| \in (0,1).$$

$$1 - \frac{\rho^2}{2} \in \left(\frac{1}{2}, 1\right)$$

(b) \exists constants c, c' s.t. for $\forall \delta \in (0, 1)$, and sample size $n \geq c \cdot \frac{\sigma^2}{\rho^4} \left(d + \log \frac{1}{\delta} \right)$,
 the sample EM updates $\theta^{t+1} = M_n(\theta^t)$ satisfies the upper bound.

sample EM.

$$(\text{Job}) \quad \|\theta^t\|_2 \leq \|\theta^0\|_2 \left(1 - \frac{\rho^2}{2}\right)^t + \frac{c' \left(\|\theta^0\|_2 \sigma^2 + \rho \sigma \right)}{\rho^2} \sqrt{\frac{d + \log \frac{1}{\delta}}{n}}$$

with probability at least $1 - \delta$.

By Thm 1 (a)

$$\|\theta^1\|_2 = \|M(\theta^0)\|_2 \leq \left(1 - \frac{\rho^2}{2}\right) \|\theta^0\|_2$$

$$\|\theta^2\|_2 \leq \left(1 - \frac{\rho^2}{2}\right) \|\theta^1\|_2 \leq \left(1 - \frac{\rho^2}{2}\right)^2 \|\theta^0\|_2$$

$$\vdots$$

$$\|\theta^T\|_2 \leq \left(1 - \frac{\rho^2}{2}\right)^T \|\theta^0\|_2$$

We want $\|\theta^T\|_2 \leq \left(1 - \frac{\rho^2}{2}\right)^T \|\theta^0\|_2 \leq \varepsilon$ $\forall \varepsilon > 0$. since $\theta^* = 0$.

$$\Rightarrow \left(1 - \frac{\rho^2}{2}\right)^T \leq \frac{\varepsilon}{\|\theta^0\|_2}$$

$$\Rightarrow T \geq \frac{1}{\log \frac{1}{\left(1 - \frac{\rho^2}{2}\right)}} \cdot \log \left(\frac{\|\theta^0\|_2}{\varepsilon} \right).$$

so $\|\theta^T\|_2 \leq \varepsilon$ for $T \geq \frac{1}{\log \frac{1}{\left(1 - \frac{\rho^2}{2}\right)}} \cdot \log \left(\frac{\|\theta^0\|_2}{\varepsilon} \right)$

It converges in $O(\log \frac{1}{\varepsilon})$ steps to an ε -ball around $\theta^* = 0$.

Theoretically, the updates converge at a geometric rate to the true parameter $\theta^* = 0$.

Proof of Thm 1 (a) :

$$W_{\theta}(X) = \frac{\pi \cdot \exp\left(-\frac{\|\theta - x\|^2}{2\sigma^2}\right)}{\pi \cdot \exp\left(-\frac{\|\theta - x\|^2}{2\sigma^2}\right) + (1-\pi) \cdot \exp\left(-\frac{\|\theta + x\|^2}{2\sigma^2}\right)} = \frac{\pi}{\pi + (1-\pi) \cdot \exp\left(\frac{-2\theta^T x}{\sigma^2}\right)}$$

$$\nabla_{\theta} W_{\theta}(X) = \frac{2\pi(1-\pi)x}{\sigma^2} \cdot \frac{1}{\left[\pi \cdot e^{\frac{\theta^T x}{\sigma^2}} + (1-\pi) \cdot e^{-\frac{\theta^T x}{\sigma^2}}\right]^2}$$

For a scalar $u \in [0,1]$, define function $h(u) = W_{u\theta}(X) = \frac{\pi}{\pi + (1-\pi) \cdot e^{-\frac{2u\theta^T x}{\sigma^2}}}$

$$h'(u) = \frac{2\pi(1-\pi) \cdot \frac{\theta^T x}{\sigma^2}}{\left[\pi \cdot e^{\frac{u\theta^T x}{\sigma^2}} + (1-\pi) \cdot e^{-\frac{u\theta^T x}{\sigma^2}}\right]^2}$$

$$= \nabla W_{u\theta}(X)^T \cdot \theta$$

$$\theta_u = u\theta, \quad u \in [0,1].$$

By Taylor expansion,

$$\begin{aligned}
 \|m(\theta)\|_2 &= \|E_x[2x(m_\theta(x) - m_0(x))]\|_2 \\
 &= \|E_x[2x \int_0^1 h'(u) du]\|_2 \quad \text{Taylor expansion,} \\
 &= 4\pi(1-\pi) \cdot \left\| \int_0^1 E_x \left[\frac{XX^T}{\sigma^2 \left(\pi \cdot \exp\left(\frac{\theta u^T x}{\sigma^2}\right) + (1-\pi) \exp\left(-\frac{\theta u^T x}{\sigma^2}\right) \right)^2} \right] du \right\|_2 \\
 &\quad \text{matrix } P_{\theta u}(x) \\
 &= 4\pi(1-\pi) \cdot \left\| \theta \cdot \int_0^1 E_x[P_{\theta u}(x)] du \right\|_2 \\
 &\stackrel{\text{by submultiplicativity property}}{\leq} 4\pi(1-\pi) \cdot \|\theta\|_2 \cdot \left\| \int_0^1 E_x[P_{\theta u}(x)] du \right\|_2 \quad \ell_2\text{-norm of matrix}
 \end{aligned}$$

$$m_\theta(x) = m_0(x) + \int_0^1 h'(u) du$$

$$= 4\pi(1-\pi) \cdot \|\theta\|_2 \cdot \left\| \int_0^1 E_x[\Gamma_{\theta u}(x)] du \right\|_{op}.$$

$\|\cdot\|_{op}$ is the largest singular value of a matrix.

$$\leq 4\pi(1-\pi) \cdot \|\theta\|_2 \cdot \max_{u \in [0,1]} \left\| E_x[\Gamma_{\theta u}(x)] \right\|_{op}$$

$$\leq 4\pi(1-\pi) \cdot \|\theta\|_2 \cdot \frac{1-p^2/2}{1-p^2}, \quad \pi = \frac{1}{2}(1-p)$$

$$= \left(1 - \frac{p^2}{2}\right) \cdot \|\theta\|_2.$$

$$\text{Also, we need to prove } \max_{u \in [0,1]} \left\| E[\Gamma_{\theta u}(x)] \right\|_{op} \leq \frac{1-p^2/2}{1-p^2}.$$

Thm 1 (b). \exists constants c, c' s.t. for $\forall \delta \in (0, 1)$, and a sample size $n \geq c \cdot \frac{\sigma^2}{\rho^4} \left(d + \log \frac{1}{\delta} \right)$,

the sample EM updates $\theta^{t+1} = M_n(\theta^t)$ satisfied the upper bound

sample EM.

$$(\text{lab}) \quad \|\theta^t\|_2 \leq \|\theta^0\|_2 \left(1 - \frac{\rho^2}{2}\right)^t + \frac{c' \left(\|\theta^0\|_2 \sigma^2 + \rho \sigma \right)}{\rho^2} \cdot \sqrt{\frac{d + \log(\frac{1}{\delta})}{n}},$$

with probability at least $1 - \delta$.

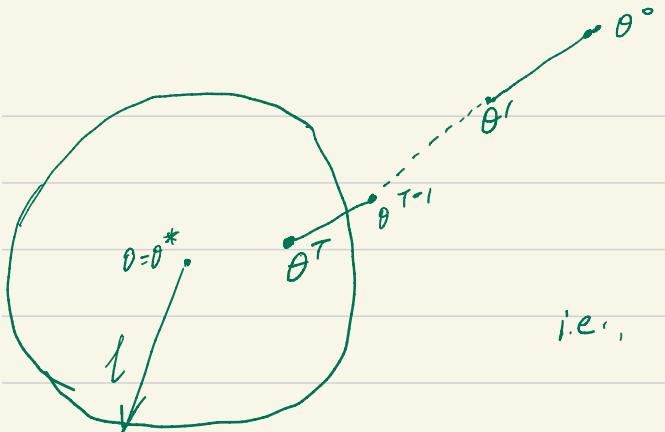
Intuition :

Our ultimate interest is in the behavior of the finite-sample EM also,
it is practical, population EM is an idealization.

The finite-sample EM updates are based on $M_n(\theta)$.

population EM updates are based on $M(\theta)$.

so important : the empirical process $\left\{ \|M_n(\theta) - M(\theta)\|_2 : \theta \in \Theta_{\epsilon=0} \right\}$



Let ℓ be an upper bound on the supremum of this empirical process that holds with probability at least $1-\delta$.

i.e., $P \left[\sup_{\|\theta - \theta^*\| \leq \ell} \|M_n(\theta) - M(\theta)\| \leq \ell \right] \geq 1-\delta.$

So we want to show that the sample-EM iterates converge to a near-optimal solution, i.e., a point whose distance from $\theta^* = 0$ is at most ℓ , or a constant multiple of ℓ .

Proof of Thm 1(b):

By Lemma 1, for \forall radius r , threshold $\delta \in (0,1)$, large enough n , \exists constants c, c' .

$$\text{we have } P \left[\sup_{\|\theta\|_2 \leq r} \|M_n(\theta) - M(\theta)\|_2 \leq \frac{c'}{2} (6^2 r + \rho s) \sqrt{\frac{d + \log(1/\delta)}{n}} \right] \geq 1 - \delta.$$

as long as $n \geq c \cdot \frac{\sigma^2}{\rho^4} (d + \log \frac{1}{\delta})$.

- $\|\theta' - \theta\|_2 = \|M_n(\theta') - \theta\|_2 = \|m(\theta') + M_n(\theta') - M(\theta')\|_2$

$$\begin{aligned} &\stackrel{\text{triangle inequality}}{\leq} \|m(\theta')\|_2 + \|M_n(\theta') - M(\theta')\|_2 \\ &\leq \underbrace{(1 - \frac{\rho^2}{2})\|\theta'\|_2}_{\text{by Thm 1(a)}} + \underbrace{\frac{c'}{2} (6^2 r + \rho s) \sqrt{\frac{d + \log(\frac{1}{\delta})}{n}}}_{\text{from previous result}} \end{aligned}$$

By mathematical induction

$$\|\theta^t - \theta^0\|_2 \leq \left(1 - \frac{\rho^2}{2}\right) \cdot \|\theta^{t-1} - \theta^0\|_2 + \underbrace{\frac{c'}{2}(G^2r + \rho G) \cdot \sqrt{\frac{d + \log(\frac{1}{\delta})}{n}}}_{\ell}.$$

- By iteration, it shows that.

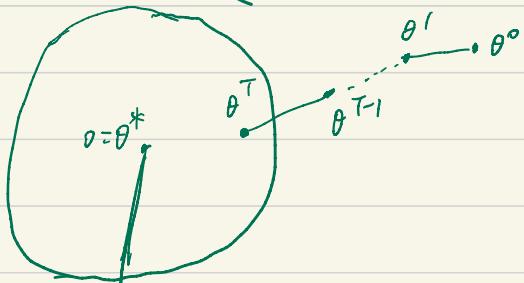
$$\begin{aligned}\|\theta^t\|_2 &\leq \left(1 - \frac{\rho^2}{2}\right) \cdot \|\theta^{t-1}\|_2 + \ell \\ &\leq \left(1 - \frac{\rho^2}{2}\right) \left[\left(1 - \frac{\rho^2}{2}\right) \cdot \|\theta^{t-2}\|_2 + \ell \right] + \ell \\ &\leq \left(1 - \frac{\rho^2}{2}\right)^t \cdot \|\theta^0\|_2 + \left[\sum_{i=0}^{t-1} \left(1 - \frac{\rho^2}{2}\right)^i \right] \cdot \ell \\ &\leq \left(1 - \frac{\rho^2}{2}\right)^t \cdot \|\theta^0\|_2 + \frac{1}{1 - \left(1 - \frac{\rho^2}{2}\right)} \cdot \ell \\ &= \left(1 - \frac{\rho^2}{2}\right)^t \cdot \|\theta^0\|_2 + \frac{c' \cdot (G^2r + \rho G)}{\rho^2} \cdot \sqrt{\frac{d + \log(\frac{1}{\delta})}{n}}\end{aligned}$$

Since it is for $\forall r > 0$, we can let $r = \|\theta^0 - \theta^0\|_2$.

So. (1^ob). $\|\theta^t\|_2 \leq \|\theta^0\|_2 \left(1 - \frac{\rho^2}{2}\right)^t + \frac{c'(\|\theta^0\|_2 \sigma^2 + \rho\sigma)}{\rho^2} \cdot \sqrt{\frac{d + \log(\frac{t}{\delta})}{n}}$ with prob. at least $1 - \delta$.

- Note that bound (1^ob) has 2 terms., the first decreases geometrically in iteration, the second is independent of t.

Thus it guarantees that the iterations converge geometrically to a ball of radius $O\left(\frac{c'}{2}(\|\theta^0\|_2 \sigma^2 + \rho\sigma) \cdot \sqrt{\frac{d + \log(\frac{t}{\delta})}{n}}\right)$.



- For a fixed sample size n , the bound

$$(10b) \|\theta^t\|_2 \leq \|\theta^0\|_2 \left(1 - \frac{\rho^2}{2}\right)^t + \underbrace{\frac{c'(\|\theta^0\|_2 \cdot \sigma^2 + \rho \cdot \sigma)}{\rho^2} \sqrt{\frac{d + \log(\frac{1}{\delta})}{n}}}_{\alpha}.$$

This bound gives a reasonable choice of iteration number T .

For iteration number T s.t. $T \geq \log_{1 - \frac{\rho^2}{2}} \frac{\alpha}{\|\theta^0\|_2}$,

$$\left(1 - \frac{\rho^2}{2}\right)^T \leq \left(1 - \frac{\rho^2}{2}\right)^{\log_{1 - \frac{\rho^2}{2}} \frac{\alpha}{\|\theta^0\|_2}} = \frac{\alpha}{\|\theta^0\|_2}.$$

$$\text{so } \|\theta^T\|_2 \leq \|\theta^0\|_2 \left(1 - \frac{\rho^2}{2}\right)^T + \alpha = \|\theta^0\|_2 \cdot \frac{\alpha}{\|\theta^0\|_2} + \alpha = 2\alpha.$$

This T -choice guarantees that the first term in (10b) is dominated by the second term, so

$$\|\theta^T\|_2 \leq 2\alpha, \text{ with prob. at least } 1 - \delta.$$

So sample-EM converges in $O\left(\log\left(\frac{n}{\alpha}\right)\right)$ steps to a ball of radius $(\frac{\alpha}{n})^{\frac{1}{2}}$.
It is also a very fast convergence.

- We see fast convergence of EM in unbalanced fit, both sample and population.

Balanced fit . $\pi = \frac{1}{2}$.

EM is worse in balanced case.

Thm 2 : convergence rate of population EM update for balance fit.

$\pi = \frac{1}{2}$. we fit model (3). $\frac{1}{2}\phi(x; \theta, \sigma^2 I_d) + \frac{1}{2}\phi(x; -\theta, \sigma^2 I_d)$ to $\pi r(\theta, \sigma^2 I_d)$

Then population EM operator

$$M(\theta) = E_x [(2w_\theta(x) - 1)x]$$

satisfies:

(Cn). for all $\theta \neq 0$, we have

$$(3a) \quad \frac{\|M(\theta)\|}{\|\theta\|_2} \leq \gamma_{up}(\theta) := \boxed{1-p + \frac{p}{1+\frac{\|\theta\|_2^2}{2\sigma^2}}} < 1.$$

$$p + \frac{1-p}{1+\frac{\|\theta\|_2^2}{2\sigma^2}}$$

p is a constant.
 $p := P(|X| \leq 1) + \frac{1}{2}P(|X| > 1)$

where $X \sim N(0, 1)$.

(b) for all $\theta \neq 0$ s.t. $\|\theta\|_2^2 \leq \frac{5\sigma^2}{8}$, we have

$$(13b) \quad \frac{\|M(\theta)\|_2}{\|\theta\|_2} \geq \gamma_{\text{low}}(\theta) := \frac{1}{1 + \frac{2\|\theta\|_2^2}{\sigma^2}}.$$

The contraction parameter is changing with $\|\theta\|$.

(13a): $\|\theta\|_2 \downarrow 0$, $\gamma_{\text{up}}(\theta) = p + \frac{1-p}{1 + \frac{\|\theta\|_2^2}{2\sigma^2}} \uparrow 1$, so convergence rate slows down.

(13b) after $\|\theta\|_2 \leq \frac{5\sigma^2}{8}$, contraction parameter lower bound $\gamma_{\text{low}}(\theta) \geq \frac{4}{3}$. and $\gamma_{\text{low}}(\theta) \uparrow 1$.

In balanced setting,
population EM still converges to the true parameter $\theta^* = 0$. from arbitrary
initialization, but the convergence gets exponentially slower when

the iteration close to the true θ^* .

It is guaranteed to have $\|\theta^T\|_2 \leq \sqrt{6}$ after running $T := \frac{\log\left(\frac{\|\theta^0\|_2^2}{2\sigma^2}\right)}{\log\left(\frac{2}{2-p}\right)}$ iterations.

After $\|\theta^T\|_2 \leq \sqrt{6}$,

the convergence rate becomes sub-geometric,

for $\varepsilon \in (0, \sqrt{6})$, $\|\theta^{T_0+t}\|_2 \leq \varepsilon$ for $t \geq \frac{C\sigma^2}{\varepsilon^2} \log\left(\frac{6}{\varepsilon}\right)$.