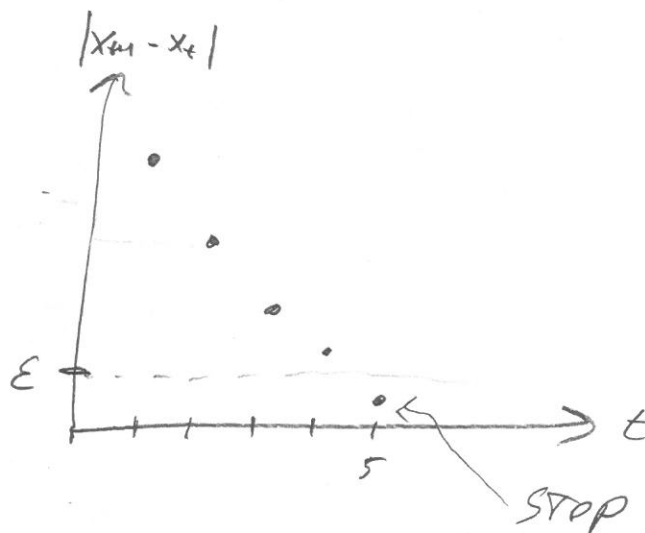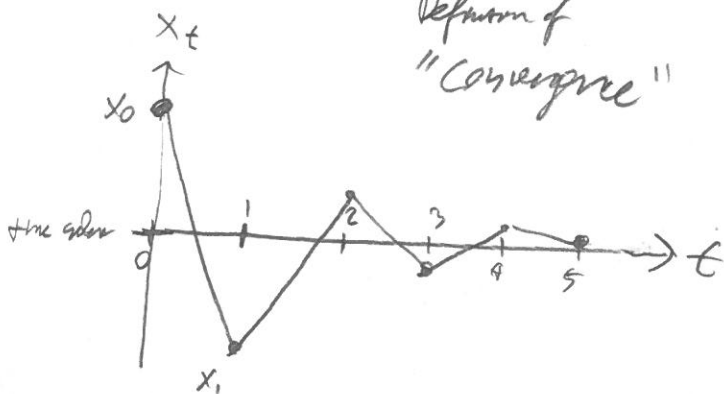N-R Method:   $f(x) = 0$  solve for $x$

① guess solution $x_0$                              } Initialize the algorithm

② Calc.  $x_1 = x_0 - \dfrac{f(x_0)}{f'(x_0)}$         } Iterate algorithm

③ Repeat Step 2 until $\underline{|x_{t+1} - x_t| < \varepsilon}$.  Return $x^*$

                                Definition of
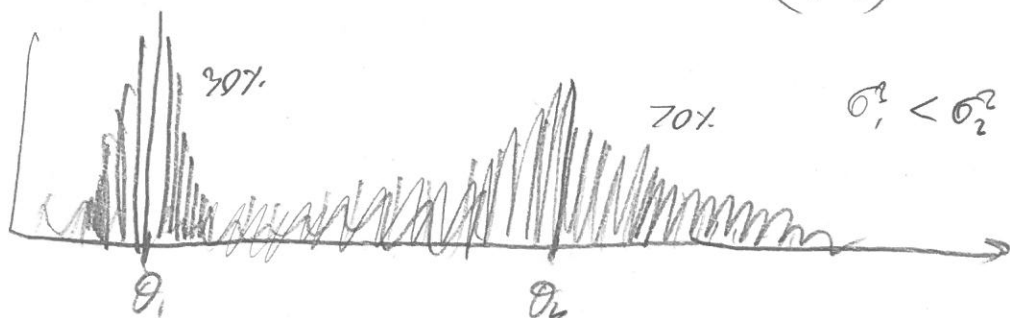                                "Convergence"



New Problem whose solution is also an iterative algorithm

$$X_1, \ldots, X_n \mid \vec{\theta}_1, \ldots, \vec{\theta}_m, \gamma_1, \ldots, \gamma_m \overset{iid}{\sim} \sum_{m=1}^{M} \gamma_m P_m(X \mid \vec{\theta}_m)$$

$\Rightarrow$ the likelihood is a mixture.  The canonical example is $M = 2$, $P_m = $ Normal:

$$X_1, \ldots, X_n \mid \theta_1, \sigma_1^2, \theta_2, \sigma_2^2, \rho \overset{iid}{\sim} \overset{\gamma_1}{\rho} N(\overset{\vec{\theta}_1}{\theta_1}, \sigma_1^2) + \overset{\gamma_2}{(1-\rho)} N(\overset{\vec{\theta}_2}{\theta_2}, \sigma_2^2)$$



30%          70%        $\sigma_1^2 < \sigma_2^2$

$\theta_1$                    $\theta_2$

Goal: inference for $\theta_1, \sigma_1^2, \theta_2, \sigma_2^2, \rho$

let $P(\theta_1, \sigma_1^2, \theta_2, \sigma_2^2, \varrho) = P(\theta_1) P(\sigma_1^2) P(\theta_2) P(\sigma_2^2) P(\varrho) = \frac{1}{\sigma_1^2} \cdot \frac{1}{\sigma_2^2}$

$\nabla_{\vec{\theta}}$

$\propto P(x|\theta) P(\theta) =$

$\Rightarrow P(\vec{\theta}|\vec{x}) \propto \left( \prod_{i=1}^{n} \varrho N(\theta_1, \sigma_1^2) + (1-\varrho) N(\theta_2, \sigma_2^2) \right) \frac{1}{\sigma_1^2} \frac{1}{\sigma_2^2} = K\left( \theta_1, \sigma_1^2 \theta_2, \sigma_2^2, \varrho | x_{1 \cdots n} \right)$

grid search?

Maybe ... but 5 dimensions

Obviously not conjugate nor of known form
$\Rightarrow \hat{\theta}_{MAP}, \hat{\theta}_{MLE}, \hat{\theta}_{MAP}$ ??? HARD

$\mathcal{E}_{\theta_1} = \langle \cdots \rangle, \quad \mathcal{E}_{\sigma_1^2} = \langle \cdots \rangle, \quad \mathcal{E}_{\theta_2} = \langle \cdots \rangle, \quad \mathcal{E}_{\sigma_2^2} = \langle \cdots \rangle, \quad \mathcal{E}_{\varrho} = \langle \cdots \rangle$

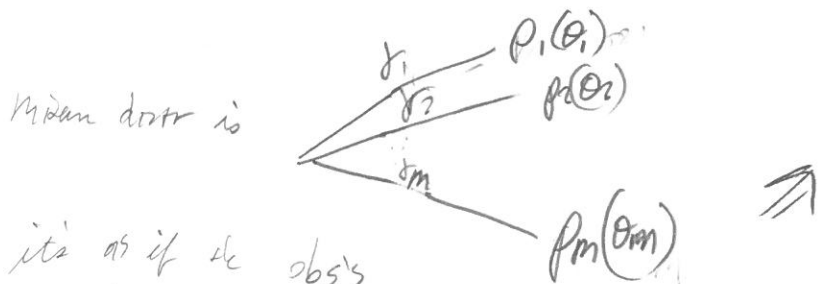You can get reasonable posterior est's of $\hat{\theta}_{MMSE}$'s.

Is there a Better way?

What if we knew that $X_1$ "belonged to" $m=2$,
$X_2$ "belonged to" $m=1$,
$X_3$ "belonged to" $m=1$,
$\vdots$
$X_n$ "belonged to" $m=1$

Define:

$I_1 = \mathbb{1}_{X_1 \text{ belongs to } m=1}$

$I_2 = \mathbb{1}_{X_2} \cdots$

$I_i = \mathbb{1}_{X_i} \cdots$

e.g.

$I_1 = 1$
$I_2 = 1$
$I_3 = 0$
$I_4 = 1$
$I_5 = 0$
$\vdots$
$I_n = 1$

called latent variables: they are important but unobserved!

Mixture distr is

$\gamma_1 \to P_1(\theta_1)$
$\gamma_2 \to P(\theta_2)$
$\gamma_m \to P_m(\theta_m)$

it's as if the obs's "choose a path"
(as long as you consider a very large # of obs)

let $I = \{ I_1 \cdots I_n \}$

Note $f(z) = \int f(z,y)\,dy = \int f(z|y)\,f(y)\,dy$

let $\theta = \{\theta_1, \sigma_1^2, \theta_2, \sigma_2^2, \rho\}$
let $I = \{I_1, \ldots, I_z\}$

$$P(X|\theta) = \int P(X, I|\theta)\,dI = \int P(X|I,\theta)\,P(I|\theta)\,dI$$

. data augmentation.. add & lower vars!

$$P(\theta|x) \propto P(X|\theta)\,P(\theta) = \int P(X|I,\theta)\,P(I|\theta)\,dI\,P(\theta) = \int P(X|I,\theta)\,P(I|\theta)\,P(\theta)\,dI = k(\theta|x)$$

$$\boxed{k(\theta|x) = \int k(\theta, I|x)\,dI}$$

$k(\theta_1, \sigma_1^2, \theta_2, \sigma_2^2, \rho|x) = \int \cdots \int P(X|I_1, \ldots, I_z, \theta_1, \sigma_1^2, \theta_2, \sigma_2^2) \, P(I_1, \ldots, I_z|\theta_1, \sigma_1^2, \theta_2, \sigma_2^2) \, P(\theta_1, \sigma_1^2, \theta_2, \sigma_2^2, \rho)\,dI_1 \cdots dI_z$

$\prod_{i=1}^{n} \frac{1}{\sqrt{2\pi\sigma_i^2}}$

All this doesn't really help us too much. But what if you cared about $\hat{\theta}_{MAP}$ = argmax $\{k(\theta|x)\}$

New!

It turns out there was an algorithm called the Expectation-Maximization Algorithm that can converge to $\hat{\theta}_{MAP}$. How?   (E-M)  (1977)

Step 1: let $\theta = \theta_0$ (a guess of $\theta_{MAP}$ just like N-R!)

Step 3: Define: $\mathcal{L}(\theta; I, X) = P(X|I,\theta)\,P(I|\theta)\,P(\theta)$

Find $\hat{\theta}_1$   un $\mathcal{L}'(\theta|I,X) = 0$. Note: $\hat{\theta}_{MLE}$ is a function of $I, X$   the max step
(ie MLE)

Step 2: Let $I_0 = E[I|X, \theta = \theta_0]$   the Expectation step

Step 4: Run Step 2-3 until $\|\hat{\theta}_{new} - \theta_{old}\| < \varepsilon$. ie until "converge"

What would the EM algorithm be in our mixture of two models case?

Step 1: $\quad \theta_{1,0} = 0, \ \sigma_{1,0}^2 = 1, \ \theta_{2,0} = 0, \ \sigma_{2,2}^2 = 0, \ \varrho = 0.5$

Step 2: $\quad I_{1,0} = E\left[I_1 \mid X_1, \dots X_n, \ \theta_1 = \theta_{1,0}, \ \sigma_1^2 = \sigma_{1,0}^2, \ \theta_2 = \theta_{2,0}, \ \sigma_2^2 = \sigma_{2,0}^2, \ \varrho = \varrho_0\right]$

If $\nearrow = P(I_1 = 1 \mid X_1, \dots) = \dfrac{P(X \mid I=1, \dots) \, P(I=1 \mid \dots)}{P(X \mid \dots)}$

$\quad P(X \mid \dots) = P(X \mid I=1, \dots) + P(X \mid I=0, \dots)$

$Q = \mathbb{1}_A \sim \text{Bern}(P(A))$

$E[Q] = P(A)$

$P(A \mid B) = \dfrac{P(B \mid A) P(A)}{P(B)}$

$= \dfrac{\dfrac{1}{\sqrt{2\pi\sigma_{1,0}^2}} e^{-\frac{1}{2\sigma_{1,0}^2}(X_1 - \theta_{1,0})^2} \, \varrho}{\dfrac{1}{\sqrt{2\pi\sigma_{1,0}^2}} e^{-\frac{1}{2\sigma_{1,0}^2}(X_1 - \theta_{1,0})^2} \varrho + \dfrac{1}{\sqrt{2\pi\sigma_{2,0}^2}} e^{-\frac{1}{2\sigma_{2,0}^2}(X_1 - \theta_{2,0})^2} (1-\varrho)}$

$I_{2,0} = \quad \dots$

$\vdots$

$I_{n,0} = \quad \dots$

Step 3: $\quad \overset{P(X \mid I, \theta)}{\overset{\|}{}} \quad \overset{P(I \mid \theta)}{\overset{\|}{}}$

$\mathcal{L}\left(\theta_1, \sigma_1^2, \theta_2, \sigma_2^2, \varrho \mid I_1 \dots I_n, X_1 \dots X_n\right) = \left(\prod_{i=1}^{n} P(X_i \mid I_i, \theta_1, \sigma_1^2, \theta_2, \sigma_2^2, \varrho)\right) \left(\prod_{i=1}^{n} \varrho^{I_i}(1-\varrho)^{1-I_i}\right) \frac{1}{\sigma_1^2} \cdot \frac{1}{\sigma_2^2}$

$= \frac{1}{\sigma_1^2} \frac{1}{\sigma_2^2} \prod_{i=1}^{n} N(\theta_1, \sigma_1^2)^{I_i} N(\theta_2, \sigma_2^2)^{1-I_i} \prod \varrho^{I_i}(1-\varrho)^{1-I_i}$

$= \frac{1}{\sigma_1^2} \frac{1}{\sigma_2^2} \prod_{i=1}^{n} \left(\frac{1}{\sqrt{2\pi\sigma_1^2}} e^{-\frac{1}{2\sigma_1^2}(X_i - \theta_1)^2}\right)^{I_i} \left(\frac{1}{\sqrt{2\pi\sigma_2^2}} e^{-\frac{1}{2\sigma_2^2}(X_i - \theta_2)^2}\right)^{1-I_i} \varrho^{I_i}(1-\varrho)^{1-I_i}$

$= \left(\frac{1}{\sqrt{2\pi}}\right)^{2n} (\sigma_1^2)^{-1} (\sigma_2^2)^{-1} (\sigma_1^2)^{-\frac{\sum I_i}{2}} (\sigma_2^2)^{\frac{(1-I_i)}{2}} e^{-\frac{1}{2\sigma_1^2}\sum I_i (X_i - \theta_1)^2} e^{-\frac{1}{2\sigma_2^2}\sum(1-I_i)(X_i - \theta_2)^2} \varrho^{\sum I_i}(1-\varrho)^{n - \sum I_i}$

$\ell(\dots \mid \dots) = 2n \ln\left(\frac{1}{\sqrt{2\pi}}\right) - \left(\frac{n}{2}+1\right) \ln(\sigma_1^2) - \left(\frac{n}{2}+1\right) \ln(\sigma_2^2) - \ \downarrow \quad \wedge \quad + \sum I_i \ln(\varrho) + (n - \sum I_i) \ln(1-\varrho)$

$\frac{1}{2\sigma_1^2} \sum I_i (X_i^2 - 2X_i\theta + \theta^2) = \frac{\sum I_i X_i^2}{2\sigma_1^2} - \frac{\theta \sum X_i I_i}{\sigma_1^2} + \frac{\theta^2 \sum I_i}{2\sigma_1^2}$

maybe $\hat{\theta}_1$ se $\frac{\partial}{\partial \theta_1}[\ell \dots] = 0$

$$\Rightarrow \frac{\sum X_i I_i}{\sigma_1^2} + \frac{\theta \sum I_i}{\sigma_1^2} = 0 \Rightarrow \hat{\theta}_1 = \frac{\sum X_i I_i}{\sum I_i} \approx \overline{X}_1 \quad \text{makes sense}$$

(→ sum of x's for m=1)

(↑ # in m=1)

$$\hat{\theta}_2 = \dots \quad \frac{\sum X_i (-I_i)}{\sum 1 - I_i} \approx \overline{X}_2$$

$\hat{\sigma}_1^2$ set $\frac{\partial}{\partial \sigma_1^2}[\ell \dots] = 0$

$$\Rightarrow \frac{-\left(\frac{I_i}{2}+1\right) 2}{\sigma_1^2} + \frac{1}{2(\sigma_1^2)^2} \sum I_i (X_i - \theta_1)^2 = 0$$

SSE for m=1

$$\Rightarrow -(\sum I_i + 2) + \frac{\sum I_i (X_i - \theta_1)^2}{\sigma_1^2} = 0 \Rightarrow \hat{\sigma}_1^2 = \frac{\sum I_i (X_i - \theta_1)^2}{\sum I_i + 2} \approx \text{sample var for } m=1$$

$$\hat{\sigma}_2^2 = \dots \quad \frac{\sum (1 - I_i)(X_i - \theta_2)^2}{\sum (1 - I_i) + 2} \approx \text{sample var for } m=2$$

$\hat{\rho}$ set $\frac{\partial}{\partial \rho}[\ell(\dots)] = 0$

$$\frac{\sum I_i}{\rho} - \frac{n - \sum I_i}{1 - \rho} = 0$$

$$\Rightarrow \frac{\sum I_i}{\rho} = \frac{n - \sum I_i}{1 - \rho} \Rightarrow \sum I_i - \rho\sum I_i = \rho n - \rho\sum I_i \Rightarrow \hat{\rho} = \frac{\sum I_i}{n} \quad \text{why? obvious...}$$

Now iterate !!!