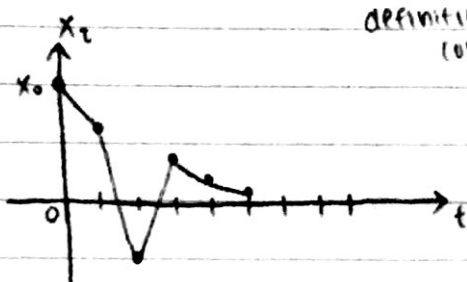


5/4/17

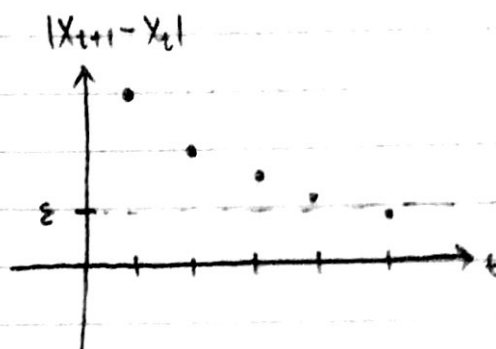
## Newton-Raphson

 $f(x) = 0$  solve for  $x$ . Given  $\epsilon$ .

1. Guess solution is  $x_0$
2. Calculate  $x_1 = x_0 - \frac{f(x_0)}{f'(x_0)}$  } iterative step is an iteration algorithm
3. Repeat step 2 until  $|x_{t+1} - x_t| \leq \epsilon$



definition of convergence



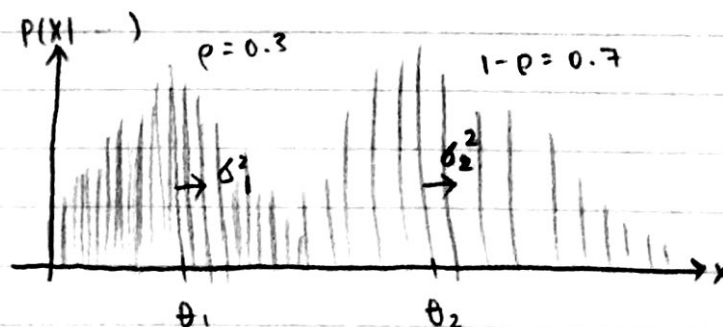
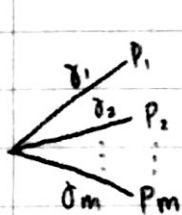
mixture

Consider the model

$$X_1, \dots, X_n | \vec{\theta}_1, \dots, \vec{\theta}_M, \gamma_1, \dots, \gamma_M \stackrel{iid}{\sim} \sum_{m=1}^M \gamma_m P_m(\vec{\theta}_m) \text{ s.t. } \gamma_1 + \gamma_2 + \dots + \gamma_M = 1$$

eg.

$$X_1, \dots, X_n | \underbrace{\theta_1, \sigma_1^2}_{\vec{\theta}_1}, \underbrace{\theta_2, \sigma_2^2}_{\vec{\theta}_2}, p \stackrel{iid}{\sim} \underbrace{p N(\theta_1, \sigma_1^2)}_{\gamma_1 P_1(\vec{\theta}_1)} + \underbrace{(1-p) N(\theta_2, \sigma_2^2)}_{\gamma_2 P_2(\vec{\theta}_2)}$$



→ ex: height of male &amp; female students where both are normally distributed.

$$P(\theta_1, \sigma_1^2, \theta_2, \sigma_2^2, p | x) \propto P(x | \theta_1, \sigma_1^2, \theta_2, \sigma_2^2, p) P(\theta_1, \sigma_1^2, \theta_2, \sigma_2^2, p)$$

$$= \underbrace{P(\theta_1)}_{\propto 1} \underbrace{P(\sigma_1^2)}_{\propto \frac{1}{\sigma_1^2}} \underbrace{P(\theta_2)}_{\propto 1} \underbrace{P(\sigma_2^2)}_{\propto \frac{1}{\sigma_2^2}} \underbrace{P(p)}_{\propto 1}$$

$$= \left( \prod_{i=1}^n e^{-\frac{1}{2\sigma_1^2} (x_i - \theta_1)^2} + (1-p) \frac{1}{\sqrt{2\pi\sigma_2^2}} e^{-\frac{1}{2\sigma_2^2} (x_i - \theta_2)^2} \right) \underbrace{\frac{1}{\sigma_1^2} \frac{1}{\sigma_2^2}}_{\text{prior}} = K(\theta_1, \sigma_1^2, \theta_2, \sigma_2^2, p | x)$$

How is our inference?

Grid Search

$$G_{\theta_1} = \langle \dots \rangle, G_{\theta_2} = \langle \dots \rangle, \dots$$

$\Rightarrow$  inaccurate in large  $M$  (ie: high dimension).

$\hookrightarrow \Delta$  is large and the grid will be spread out and lose accuracy.

What if we know which component each  $x_i$  belongs to?

Define:  $I_1 := \mathbb{1}_{x_1 \text{ is in } m=1}$

$I_2 := \mathbb{1}_{x_2 \text{ is in } m=2}$

$\vdots$

$I_n := \mathbb{1}_{x_n \text{ is in } m=n}$

Let  $I = \{I_1, \dots, I_n\}$

"Latent variable/information"

$\hookrightarrow$  the  $I_i$ 's are unobserved but important

Recall

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

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

"data augmentation" - adding more data & averaging

$$\begin{aligned} P(\theta, \delta_1^2, \theta_2, \delta_2^2, p|x) &\propto \int P(x|I, \theta, \delta_1^2, \theta_2, \delta_2^2, p) P(I|\theta, \delta_1^2, \theta_2, \delta_2^2, p) P(\theta, \delta_1^2, \theta_2, \delta_2^2, p) dI \\ &= K(\theta, \delta_1^2, \theta_2, \delta_2^2, p|x) \\ &= \int \underbrace{K(\theta, \delta_1^2, \theta_2, \delta_2^2, p|x, I)}_{\text{body of integral}} dI \end{aligned}$$

Modest goal:

$$\text{Get } \hat{\theta}_{\text{MAP}} = \text{argmax}_{\theta} \{K(\theta|x)\}$$

Expectation-Maximization Algorithm (1977)

Step 1: Guess  $\hat{\theta}_{\text{MAP}} = \theta_0$  to start

Expectation step  $\rightarrow$  Step 2: Compute  $I_0 = E[I|x, \theta = \theta_0]$

Maximization step  $\rightarrow$  Step 3: Consider likelihood:  $\mathcal{L}(\theta; I_0, x) = \underbrace{K(\theta|I=I_0, x)}_{\text{body of integral}}$

and find  $\hat{\theta}_1 = \text{argmax}_{\theta} \{\mathcal{L}(\theta; I_0, x)\}$  ie: the MLE procedure

Step 4: Repeat steps 2 & 3 until  $\|\theta_{t+1} - \theta_t\| < \epsilon$ , where  $\epsilon$  is the predefined tolerance level.

E-M Initialization for our 2-normal mixture:Step 1: Initialize

$$\theta_{1,0} = 0$$

$$\delta_{1,0}^2 = 1$$

$$\theta_{2,0} = 0$$

$$\delta_{2,0}^2 = 1$$

$$p = 0.5$$

$$\rightarrow I_1 \sim \text{Bern}(P(I_1=1|\dots))$$

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

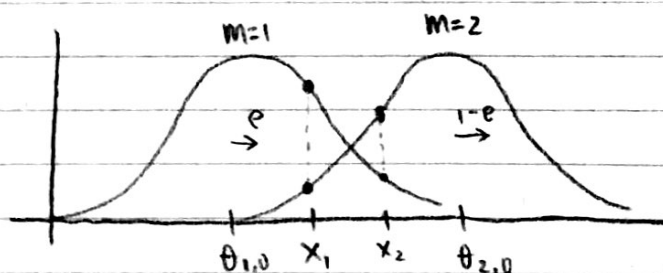
$$E[Q] = P(A)$$

Step 2:  $I_{1,0} = E[I_1 | X, \theta_1 = \theta_{1,0}, \delta_1^2 = \delta_{1,0}^2, \theta_2 = \theta_{2,0}, \delta_2^2 = \delta_{2,0}^2, p = p]$ 

$$= P(I_1=1 | X, \dots) = \frac{P(X | I_1=1, \dots) P(I_1=1 | \dots)}{P(X | \dots)} \Rightarrow \text{Bayes Theorem}$$

$$P(X | I_1=1, \dots) \cdot P(I_1=1 | \dots) + P(X | I_1=0, \dots) \cdot P(I_1=0 | \dots)$$

$$= \frac{e \cdot \frac{1}{\sqrt{2\pi\delta_{1,0}^2}} e^{-\frac{1}{2\delta_{1,0}^2} (x_1 - \theta_{1,0})^2}}{e \cdot \frac{1}{\sqrt{2\pi\delta_{1,0}^2}} e^{-\frac{1}{2\delta_{1,0}^2} (x_1 - \theta_{1,0})^2} + (1-e) \cdot \frac{1}{\sqrt{2\pi\delta_{2,0}^2}} e^{-\frac{1}{2\delta_{2,0}^2} (x_1 - \theta_{2,0})^2}}$$



$$I_{2,0} = E[I_2 | X_2, \dots]$$

$$I_{3,0} = E[I_3 | X_3, \dots]$$

$$I_{n,0} = E[I_n | X_n, \dots]$$

Step 3:  $\mathcal{L}(\theta_1, \delta_1^2, \theta_2, \delta_2^2, p; I, X) = P(X | I, \theta) P(I | \theta) P(\theta)$ 

$$= \left( \prod_{i=1}^n \left( \frac{1}{\sqrt{2\pi\delta_1^2}} e^{-\frac{1}{2\delta_1^2} (x_i - \theta_1)^2} \right)^{I_i} \left( \frac{1}{\sqrt{2\pi\delta_2^2}} e^{-\frac{1}{2\delta_2^2} (x_i - \theta_2)^2} \right)^{1-I_i} \right) \cdot \left( \prod_{i=1}^n e^{I_i (1-p)} (1-p)^{1-I_i} \right) \cdot ((\delta_1^2)^{-1} (\delta_2^2)^{-1})$$

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

$$= \ell(\theta_1, \sigma_1^2, \theta_2, \sigma_2^2, p; \mathbf{I}, \mathbf{x}) \quad \left( \frac{\sum I_i x_i^2}{2\sigma_1^2} - \frac{\theta_1 \sum x_i I_i}{\sigma_1^2} + \frac{\theta_1^2 \sum I_i}{2\sigma_1^2} \right) \left( \frac{\sum (1-I_i) x_i^2}{2\sigma_2^2} - \frac{\theta_2 \sum x_i (1-I_i)}{\sigma_2^2} + \frac{\theta_2^2 \sum (1-I_i)}{2\sigma_2^2} \right)$$

$$= n \ln \left( \frac{1}{\sqrt{2\pi}} \right) - (1 + \frac{1}{2} \sum I_i) \ln(\sigma_1^2) - (1 + \frac{1}{2} \sum (1-I_i)) \ln(\sigma_2^2) - \frac{1}{2\sigma_1^2} \sum I_i (x_i - \theta_1)^2 - \frac{1}{2\sigma_2^2} \sum (1-I_i) (x_i - \theta_2)^2 + (\sum I_i) \ln p + (\sum (1-I_i)) \ln (1-p)$$

$$\text{Get } \hat{\theta}_1 \text{ by } \frac{\partial}{\partial \theta_1} [\ell] \stackrel{\text{set}}{=} 0$$

$$\frac{\sum x_i I_i}{\sigma_1^2} - \frac{\sum \theta_1 \sum I_i}{\sigma_1^2} = 0$$

$$\hat{\theta}_1 = \frac{\sum x_i I_i}{\sum I_i} \quad \bar{x}_{\text{mix } 1}$$

$$\hat{\theta}_2 = \frac{\sum x_i (1-I_i)}{\sum (1-I_i)} \quad \bar{x}_{\text{mix } 2}$$

$$\text{Get } \hat{\sigma}_1^2,$$

$$\frac{\partial}{\partial \sigma_1^2} [\ell] = -\frac{1 + \frac{1}{2} \sum I_i}{\sigma_1^2} + \frac{1}{2(\sigma_1^2)^2} \sum I_i (x_i - \theta_1)^2 = 0$$

$$= 1 + \frac{1}{2} \sum I_i = \frac{1}{2\sigma_1^2} \sum I_i (x_i - \theta_1)^2$$

$$= 2 + \sum I_i = \frac{1}{\sigma_1^2} \sum I_i (x_i - \theta_1)^2 \Rightarrow \hat{\sigma}_1^2 = \frac{\sum I_i (x_i - \theta_1)^2}{2 + \sum I_i}$$

$$\text{likewise } \hat{\sigma}_2^2 = \frac{\sum (1-I_i) (x_i - \theta_2)^2}{2 + \sum (1-I_i)}$$

$$\hat{p} \text{ get } \frac{\partial}{\partial p} [\ell] \stackrel{\text{set}}{=} 0$$

$$= \frac{\sum I_i}{p} - \frac{\sum (1-I_i)}{1-p} = 0$$

$$= \frac{\sum I_i}{p} = \frac{\sum (1-I_i)}{1-p}$$

$$= \sum I_i - p \sum I_i = p^n - p^{\sum I_i}$$

$$= \hat{p} = \frac{\sum I_i}{n}$$