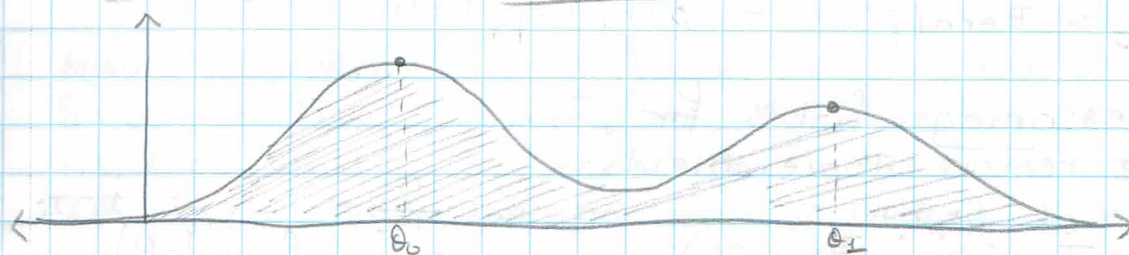


LECTURE 20

Let X_1, \dots, X_n ^{exch} $p N(\theta_0, \sigma_0^2) + (1-p) N(\theta_1, \sigma_1^2)$
 So for X_i , the density looks like



The joint density is the product of the individual densities.

$$\text{So } P(\vec{X} | \theta_0, \theta_1, \sigma_0^2, \sigma_1^2, p) = \prod_{i=1}^n p \frac{1}{\sqrt{2\pi\sigma_0^2}} e^{-\frac{1}{2\sigma_0^2}(x_i - \theta_0)^2} + (1-p) \frac{1}{\sqrt{2\pi\sigma_1^2}} e^{-\frac{1}{2\sigma_1^2}(x_i - \theta_1)^2}$$

This is also our likelihood.

We cannot compute $\hat{\theta}_{MLE}$ by hand (i.e. take \ln of the likelihood, compute the derivative with respect to $\theta_0, \theta_1, \sigma_0^2, \sigma_1^2$, & p & set it equal to zero). So what is our goal?

GOAL: ESTIMATE $\theta_0, \theta_1, \sigma_0^2, \sigma_1^2, p$

How?

① compute the MLE for θ_0, \dots, p which is hard.

② Grid sample $\theta_0, \theta_1, \sigma_0^2, \sigma_1^2, p$ but this would not be very accurate.

Numerical method ③ 5 Dimensional Newton-Raphson BUT

YOU STILL NEED TO TAKE THE DERIVATIVE AS IN ①.

Let us explore a different option...

What if we knew X_1 belongs to $N(\theta_0, \sigma_0^2)$ &
 X_2 belongs to $N(\theta_1, \sigma_1^2)$ &
 X_n belongs to $N(\theta_1, \sigma_1^2)$

Define $I_i := \mathbb{1}_{X_i} = \begin{cases} 1 & \text{if } X_i \text{ belongs to } N(\theta_0, \sigma_0^2) \\ 0 & \text{if } X_i \text{ is NOT } N(\theta_0, \sigma_0^2) \end{cases}$ ^{i.e. must be $N(\theta_1, \sigma_1^2)$}

Ex: Based off what we defined above, $I_1 = 1, I_2 = 0$ & $I_n = 0$

Now what is the distribution of I_1 ? That is,

What is the probability X_1 belongs to $N(\theta_0, \sigma_0^2)$?
Well,

$$X_1 \sim p N(\theta_0, \sigma_0^2) + (1-p) N(\theta_1, \sigma_1^2)$$

Hence X_1 is $N(\theta_0, \sigma_0^2)$ with probability p .

$$\text{Hence } I_1 \sim \text{Bern}(p) = p^{I_1} (1-p)^{1-I_1}$$

The same reasoning holds for I_2 .

Repeating the reason above yields

$$I_1, I_2, \dots, I_n \stackrel{\text{i.i.d.}}{\sim} \text{Bern}(p) \quad \prod_{i=1}^n (p^{I_i} (1-p)^{1-I_i})$$

$$\begin{aligned} \text{Hence } p(\vec{I} | p) &= p(I_1, \dots, I_n | p) = p(I_1 | p) p(I_2 | p) \dots p(I_n | p) \\ &= \text{Product of Bernoulli's} \\ &= \text{Bin}(n, p) \end{aligned}$$

$$\text{So } \vec{I} | p \sim \text{Bin}(n, p)$$

$$\begin{aligned} p(\vec{X}, \vec{I} | \theta_0, \theta_1, \sigma_0^2, \sigma_1^2, p) &= p(\vec{X} | \theta_0, \theta_1, \sigma_0^2, \sigma_1^2, p, \vec{I}) p(\vec{I} | \theta_0, \theta_1, \sigma_0^2, \sigma_1^2, p) \\ &= p(\vec{X} | \theta_0, \theta_1, \sigma_0^2, \sigma_1^2, \vec{I}) p(\vec{I} | p) \end{aligned}$$

Note $p(\vec{X} | \theta_0, \theta_1, \sigma_0^2, \sigma_1^2, p, \vec{I}) = p(\vec{X} | \theta_0, \theta_1, \sigma_0^2, \sigma_1^2, \vec{I})$. Why?

B/c you are given \vec{I} . Remember, \vec{I} is the "INDICATOR RANDOM VARIABLE". It tells you if $X_i \sim N(\theta_0, \sigma_0^2)$ OR if $X_i \sim N(\theta_1, \sigma_1^2)$. So given \vec{I} , you do NOT need know p .

Also note $p(\vec{I} | \theta_0, \theta_1, \sigma_0^2, \sigma_1^2, p) = p(\vec{I} | p)$. Why? Because I_j , (the indicator r.v. of X_j) is either 1 w.p. p or 0 w.p. $1-p$.

Hence \vec{I} depends ONLY on p .

Now remember $X_1 \sim p N(\theta_0, \sigma_0^2) + (1-p) N(\theta_1, \sigma_1^2)$. What is $p(X_1 | \theta_0, \theta_1, \sigma_0^2, \sigma_1^2, I)$? We are given I . $I=0$ OR $I=1$. We need $p(X_1 | \theta_0, \theta_1, \sigma_0^2, \sigma_1^2, I)$ to "switch" from $N(\theta_0, \sigma_0^2)$ if $I=1$ to $N(\theta_1, \sigma_1^2)$ if $I=0$. Hence

$$\begin{aligned} p(X_1 | \theta_0, \theta_1, \sigma_0^2, \sigma_1^2, I) &= (N(\theta_0, \sigma_0^2))^{I_1} (N(\theta_1, \sigma_1^2))^{1-I_1} \\ &= \left(\frac{1}{\sqrt{2\pi\sigma_0^2}} e^{-\frac{1}{2\sigma_0^2}(X_1 - \theta_0)^2} \right)^{I_1} \left(\frac{1}{\sqrt{2\pi\sigma_1^2}} e^{-\frac{1}{2\sigma_1^2}(X_1 - \theta_1)^2} \right)^{1-I_1} \end{aligned}$$

and in general

$$P(X_j | \theta_0, \theta_1, \sigma_0^2, \sigma_1^2, I) = \left(\frac{1}{\sqrt{2\pi\sigma_0^2}} e^{-\frac{1}{2\sigma_0^2}(x_j - \theta_0)^2} \right)^{I_j} \left(\frac{1}{\sqrt{2\pi\sigma_1^2}} e^{-\frac{1}{2\sigma_1^2}(x_j - \theta_1)^2} \right)^{1-I_j}$$

Again, why? Well $X_j \sim pN(\theta_0, \sigma_0^2) + (1-p)N(\theta_1, \sigma_1^2)$. Since in $P(X_j | \theta_0, \theta_1, \sigma_0^2, \sigma_1^2, I)$, we are given \vec{I} , we know what I_j equals to. $I_j = 0$ or $I_j = 1$. If $I_j = 0$, then based on how we defined I_j on page 1, we X_j is $N(\theta_1, \sigma_1^2)$. If $I_j = 1$, then we want $P(X_j | \theta_0, \theta_1, \sigma_0^2, \sigma_1^2, I)$ to be $N(\theta_0, \sigma_0^2)$ & the way we switch between them is to use the above equation b/c if $I_j = 0$ you get $N(\theta_1, \sigma_1^2)$ & if $I_j = 1$ you get $N(\theta_0, \sigma_0^2)$.

Hence

$$\begin{aligned} P(\vec{X}, \vec{I}) &= P(\vec{X} | \theta_0, \theta_1, \sigma_0^2, \sigma_1^2, \vec{I}) P(\vec{I} | p) \\ &= \prod_{i=1}^n \underbrace{\left(\frac{1}{\sqrt{2\pi\sigma_0^2}} e^{-\frac{1}{2\sigma_0^2}(x_i - \theta_0)^2} \right)^{I_i}}_{\alpha} \underbrace{\left(\frac{1}{\sqrt{2\pi\sigma_1^2}} e^{-\frac{1}{2\sigma_1^2}(x_i - \theta_1)^2} \right)^{1-I_i}}_{\beta} (p^{I_i})(1-p)^{1-I_i} \\ &= \prod_{i=1}^n (\rho \alpha)^{I_i} ((1-p)\beta)^{1-I_i} \end{aligned}$$

We then showed what we did is mathematically valid.

So originally, our likelihood would have been $P(\vec{X} | \theta_0, \theta_1, \sigma_0^2, \sigma_1^2, p)$. We then "added \vec{I} " in & we computed our new likelihood, $P(\vec{X}, \vec{I} | \theta_0, \theta_1, \sigma_0^2, \sigma_1^2, p)$. [This process is called DATA AUGMENTATION.] However, the question remains, WHY DID WE DO THIS? The answer is that if we wanted $\theta_0^{MLE}, \theta_1^{MLE}, \sigma_0^{2MLE}, \sigma_1^{2MLE}, p^{MLE}$, we would need to take the derivative of the log likelihood, $\ln(P(\vec{X} | \theta_0, \theta_1, \sigma_0^2, \sigma_1^2, p))$ & set it equal to 0 & solve for $\theta_0, \theta_1, \sigma_0^2, \sigma_1^2, p$. However, this is too difficult to do by hand.

It turns out, we can take the derivative of the log likelihood, $\ln(P(\vec{X}, \vec{I} | \theta_0, \theta_1, \sigma_0^2, \sigma_1^2, p))$ & get our MLE'S. We have:

$$\textcircled{1} \hat{\theta}_0^{\text{MLE}} = \frac{\sum x_i I_i}{n_0}$$

$$\textcircled{2} \hat{\theta}_1^{\text{MLE}} = \frac{\sum x_i (1 - I_i)}{n_1}$$

$$\textcircled{3} \hat{p}^{\text{MLE}} = \frac{n_0}{n} = \frac{n_0}{n_0 + n_1}$$

$$\textcircled{4} \hat{\sigma}_0^2{}^{\text{MLE}} = \frac{\sum (x_i - \theta_0)^2 I_i}{n_0}$$

$$\textcircled{5} \hat{\sigma}_1^2{}^{\text{MLE}} = \frac{\sum (x_i - \theta_1)^2 (1 - I_i)}{n_1}$$

Where $n_0 := \sum I_i$

$n_1 := \sum 1 - I_i$

Remark: Note that $\hat{\theta}_0^{\text{MLE}}$, $\hat{\theta}_1^{\text{MLE}}$, \hat{p}^{MLE} are like our " \bar{X} ". These MLE's are "natural averages" in the following sense:

- ① For $\hat{\theta}_0^{\text{MLE}}$, add up the # of x_i 's that belong to the $N(\theta_0, \sigma_0^2)$ mixture & divide by the number of such terms.
- ② Similar interpretation for $\hat{\theta}_1^{\text{MLE}}$ & \hat{p}^{MLE} .
- ③ $\hat{\sigma}_k^2{}^{\text{MLE}}$ is the variance if θ_k is known for $k=0, 1$.

Okay, let's recap...

We augmented in some data which allowed us to get the MLE's, which was something we were NOT able to do before. However, looking at our MLE's, there appears to be a problem: Our MLE's REQUIRES US TO KNOW \vec{I} . BUT WE DO NOT KNOW \vec{I} AHEAD OF TIME! Hence, the above equations, as it stands now, are useless.

BUT, what if we can estimate \vec{I} ? That is, what if we can estimate I_i ?

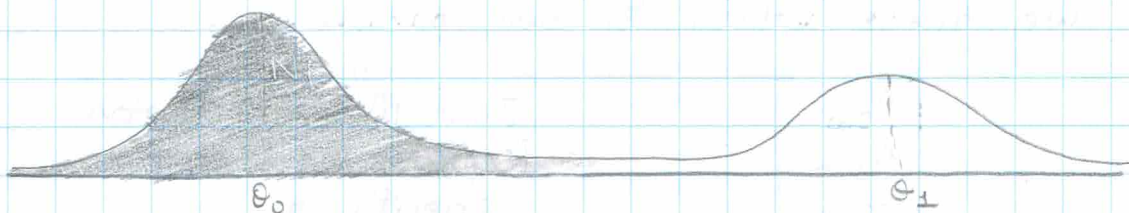
Before answering this question, we note for an INDICATOR R.V. X , we have $P(X) = E(X)$. We now return to the question.

Define $\hat{I}_1 := E[I_1 | \theta_0, \theta_1, \sigma_0^2, \sigma_1^2, p]$
 $\rightarrow = P(I_1 | \theta_0, \theta_1, \sigma_0^2, \sigma_1^2, p)$
 since \hat{I}_1 is an indicator r.v.

What is $P(I_1 | \theta_0, \theta_1, \sigma_0^2, \sigma_1^2, p)$?

Remember,

$$I_1 = \begin{cases} 1 & \text{if } X_1 \text{ belongs to } N(\theta_0, \sigma_0^2) \\ 0 & \text{if } X_1 \text{ belongs to } N(\theta_1, \sigma_1^2) \end{cases}$$

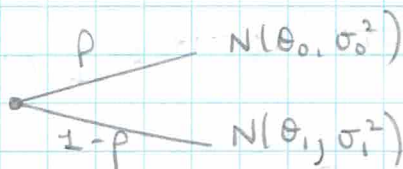


That is, how likely is it X_1 belongs to $N(\theta_0, \sigma_0^2)$ vs belonging to $N(\theta_1, \sigma_1^2)$?

From the picture, it is clear

$$\hat{I}_1 = P(I_1 | \theta_0, \theta_1, \sigma_0^2, \sigma_1^2, p) = \frac{p N(\theta_0, \sigma_0^2)}{p N(\theta_0, \sigma_0^2) + (1-p) N(\theta_1, \sigma_1^2)} = \frac{\text{frequency}}{\text{sum}}$$

i.e., the probability of being in the $N(\theta_0, \sigma_0^2)$ distribution over our entire sample space. Note we multiply by p since



Hence

$$\hat{I}_1 = P(I_1 | \theta_0, \theta_1, \sigma_0^2, \sigma_1^2, p) = \frac{p \left(\frac{1}{\sqrt{2\pi\sigma_0^2}} e^{-\frac{1}{2\sigma_0^2}(x_1 - \theta_0)^2} \right)}{p \left(\frac{1}{\sqrt{2\pi\sigma_0^2}} e^{-\frac{1}{2\sigma_0^2}(x_1 - \theta_0)^2} \right) + (1-p) \left(\frac{1}{\sqrt{2\pi\sigma_1^2}} e^{-\frac{1}{2\sigma_1^2}(x_1 - \theta_1)^2} \right)}$$

AND MORE GENERALLY, we call the following equation (6)

$$\hat{I}_i = \frac{p \left(\frac{1}{\sqrt{2\pi\sigma_0^2}} e^{-\frac{1}{2\sigma_0^2}(x_i - \theta_0)^2} \right)}{p \left(\frac{1}{\sqrt{2\pi\sigma_0^2}} e^{-\frac{1}{2\sigma_0^2}(x_i - \theta_0)^2} \right) + (1-p) \left(\frac{1}{\sqrt{2\pi\sigma_1^2}} e^{-\frac{1}{2\sigma_1^2}(x_i - \theta_1)^2} \right)}$$

BUT wait... to compute \hat{I}_i , we need $\theta_0, \theta_1, \sigma_0^2, \sigma_1^2, p$.

This is bad because these are the parameters we WANT

TO ESTIMATE! We do NOT have them ahead of time. So again it seems like this is useless... But wait! Here is a solution to this called the E-M algorithm

STEP ONE

EXPECTATION-MAXIMIZATION

Guess values for $\theta_0, \theta_1, \sigma_0^2, \sigma_1^2, \rho$. We denote these initial guesses by indexing via 0 (i.e. guess 0). That is, we guess values for our parameters:

$$\begin{bmatrix} \theta_{0,0} \\ \theta_{1,0} \\ \sigma_{0,0}^2 \\ \sigma_{1,0}^2 \\ \rho_0 \end{bmatrix}$$

(Just like in Newton Raphson we first specify a guess/starting point)

STEP TWO (called E-step)

Calculate \hat{I}_i , the augmented data, by using equation (6). So we compute

$$\hat{I}_i = \frac{\frac{1}{\rho_0 \sqrt{2\pi\sigma_{0,0}^2}} e^{-\frac{1}{2\sigma_{0,0}^2}(X_i - \theta_{0,0})^2}}{\frac{1}{\rho_0 \sqrt{2\pi\sigma_{0,0}^2}} e^{-\frac{1}{2\sigma_{0,0}^2}(X_i - \theta_{0,0})^2} + (1 - \rho_0) e^{-\frac{1}{2\sigma_{1,0}^2}(X_i - \theta_{1,0})^2}}$$

$\forall i = 1, 2, 3, \dots, n$. So we get $\hat{I}_1, \hat{I}_2, \dots, \hat{I}_n$ which we can really think of as $\hat{\hat{I}}_1, \hat{\hat{I}}_2, \dots, \hat{\hat{I}}_n$ since we used estimates for $\theta_0, \theta_1, \sigma_0^2, \sigma_1^2, \rho$.

STEP THREE (called m-step)

Compute $\hat{\theta}_{0,MLE}, \hat{\theta}_{1,MLE}, \hat{\sigma}_{0,MLE}^2, \hat{\sigma}_{1,MLE}^2, \hat{\rho}_{MLE}$ (i.e., equations (1)-(5) by using the $\hat{\hat{I}}_i$'s from step two.) we really

get $\hat{\hat{\theta}}_{0,MLE}, \hat{\hat{\theta}}_{1,MLE}, \hat{\hat{\sigma}}_{0,MLE}^2, \hat{\hat{\sigma}}_{1,MLE}^2, \hat{\hat{\rho}}_{MLE}$ since we used estimates of \hat{I}_i 's.

STEP FOUR: REPEAT STEP TWO WITH PARAMETER GUESSES FROM STEP THREE AND ITERATE UNTIL $\|\theta_{t+1} - \theta_t\| < \epsilon$, a predetermined TOLERANCE LEVEL.

Note that one may show this process does converge.

At the end of the E-M, we get estimates of our parameters $\theta_0, \theta_1, \sigma_0^2, \sigma_1^2$, & ρ . Thus we have achieved our goal listed on page 1.

As a last remark, note that the E-M is useful for when the typical MLE case breaks down b/c of the complexity of taking the derivative of the log likelihood.