

# Statistical Inference

## Lecture 06a

ANU - RSFAS

Last Updated: Wed Mar 28 11:39:40 2018

# MLE Computation: Expectation - Maximization (EM) Algorithm

- Presentation adapted from CB & *Computational Statistics*.
- The EM algorithm is a general algorithm to find MLEs when some of the data are missing (or the problem can be set in a manner that there are missing data).
- Suppose we observe all of the data  $\mathbf{y} = \{y_1, \dots, y_n\}$ , then all we do to find the MLE is maximize:

$$\ell(\boldsymbol{\theta}; \mathbf{y})$$

- Suppose we don't observe all the  $\mathbf{y}$ s then based on the notation by Donald Rubin we have  $\mathbf{y} = (\mathbf{y}_{obs}, \mathbf{y}_{miss})$ .

$$\begin{aligned} f(\mathbf{y}; \theta) &= f(\mathbf{y}_{obs}, \mathbf{y}_{miss}; \theta) \\ &= k(\mathbf{y}_{miss} | \mathbf{y}_{obs}, \theta) g(\mathbf{y}_{obs}; \theta) \end{aligned}$$

- This leads to:  $g(\mathbf{y}_{obs}; \theta) = \frac{f(\mathbf{y}; \theta)}{k(\mathbf{y}_{miss} | \mathbf{y}_{obs}, \theta)}$

$$\log [g(\mathbf{y}_{obs}; \theta)] = \log [f(\mathbf{y}_{obs}, \mathbf{y}_{miss}; \theta)] - \log [k(\mathbf{y}_{miss} | \mathbf{y}_{obs}, \theta)]$$

$$\ell_{obs}(\theta; \mathbf{y}_{obs}) = \ell_{comp}(\theta; \mathbf{y}_{obs}, \mathbf{y}_{miss}) - \log [k(\mathbf{y}_{miss} | \mathbf{y}_{obs}, \theta)]$$

- As  $\mathbf{y}_{miss}$  is missing, we replace the right side of the equation with its expectation:

$$\begin{aligned}\ell_{obs}(\boldsymbol{\theta}; \mathbf{y}_{obs}) &= E \left\{ \ell_{comp}(\boldsymbol{\theta}; \mathbf{y}_{obs}, \mathbf{y}_{miss}) \middle| \boldsymbol{\theta}', \mathbf{y}_{obs} \right\} \\ &\quad - E \left\{ \log [k(\mathbf{y}_{miss} | \mathbf{y}_{obs}, \boldsymbol{\theta})] \middle| \boldsymbol{\theta}', \mathbf{y}_{obs} \right\}\end{aligned}$$


- The EM algorithm seeks to maximize  $\ell(\theta; \mathbf{y}_{obs})$  with respect to  $\theta$  through the following process:

1. **E step:** Calculate the expectation of the complete likelihood conditional on the observed data and the current value of  $\theta$ :

$$\begin{aligned} Q(\theta | \theta^{(r)}) &= E \left\{ \ell_{comp}(\theta; \mathbf{y}_{obs}, \mathbf{y}_{miss}) \middle| \theta^{(r)}, \mathbf{y}_{obs} \right\} \\ &= \int [\ell_{comp}(\theta; \mathbf{y}_{obs}, \mathbf{y}_{miss})] k(\mathbf{y}_{miss} | \mathbf{y}_{obs}, \theta) d\mathbf{y}_{miss} \end{aligned}$$

2. **M step:** Maximize  $Q(\theta | \theta^{(r)})$  with respect to  $\theta$ . Set  $\theta^{(r+1)}$  equal to the maximizer of  $Q$ .
  3. Return to the E step unless a stopping criterion has been reached.
- I will not present a proof that the EM algorithm maximizes  $\ell(\theta; \mathbf{y}_{obs})$ . If you are interested please see either Casella and Berger Exercise 7.31 or *Computational Statistics* Section 4.2.1.

# EM Example

- Let  $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} f(y; \theta)$ .  
*prob*  $\downarrow$  *bimodal two component mixture model*  *whole density must integrate to 1*

$$f(y; \theta) = p \text{ normal}(\mu_0, \sigma_0^2) + (1 - p) \text{ normal}(\mu_1, \sigma_1^2)$$

- For this problem generally we have  $\theta = (\mu_0, \sigma_0^2, \mu_1, \sigma_1^2, p)$ .
- For our example let's simplify the problem and assume  $p = \frac{1}{2}, \sigma_0^2 = \sigma_1^2 = 1$
- We have the following likelihood:

$$L(\theta) = \prod_{i=1}^n \left[ \frac{1}{2} \text{ normal}(\mu_0, 1) + \frac{1}{2} \text{ normal}(\mu_1, 1) \right]$$

- Directly optimizing this is hard.

- To make things easier, we can introduce latent (missing) variables  $Z_1, \dots, Z_n$ .
  - Where  $Z_i = 0$  if  $Y_i$  is from  $\text{normal}(\mu_0, 1)$ .
  - Where  $Z_i = 1$  if  $Y_i$  is from  $\text{normal}(\mu_1, 1)$ .
- Why is this easier? If we know what normal distribution each  $Y$  comes from then it is easy to determine the MLEs for  $\mu_0$  and  $\mu_1$ .
- We can use the EM algorithm, where  $\mathbf{y}_{\text{miss}} = \mathbf{z}$ .
- Note:  $Z_i$  is a Bernoulli random variable.  $P(Z_i = 1) = \frac{1}{2}$ .

$$L(\theta)_{\text{comp}} = \prod_{i=1}^n \text{normal}(y_i; \mu_0, 1)^{1-z_i} \text{normal}(y_i; \mu_1, 1)^{z_i}$$



$$\begin{aligned}
 \ell_{comp}(\theta) &= \sum_{i=1}^n (1 - z_i) \log \left( \frac{1}{\sqrt{2\pi}} \exp \left( -\frac{1}{2} (y_i - \mu_0)^2 \right) \right) \\
 &\quad + \sum_{i=1}^n z_i \log \left( \frac{1}{\sqrt{2\pi}} \exp \left( -\frac{1}{2} (y_i - \mu_1)^2 \right) \right) \\
 &\quad + \text{constants} \\
 &= -\frac{1}{2} \sum_{i=1}^n (1 - z_i) (y_i - \mu_0)^2 + -\frac{1}{2} \sum_{i=1}^n z_i (y_i - \mu_1)^2 + \text{constants}
 \end{aligned}$$

1. Let's determine  $Q(\theta|\theta^{(r)})$ . Note that  $z$  is linear in the log likelihood!  
Makes our job much easier!

$$\begin{aligned} E[\ell_{comp}(\theta)] &= -\frac{1}{2} \sum_{i=1}^n (1 - E[Z_i|\mathbf{y}_{obs}, \theta^r]) (y_i - \mu_0)^2 \\ &\quad - \frac{1}{2} \sum_{i=1}^n E[Z_i|\mathbf{y}_{obs}, \theta^r] (y_i - \mu_1)^2 + constants \end{aligned}$$

- We need to determine:  $E[Z_i|\mathbf{y}_{obs}, \theta^r]$ .
- Note:  $E[Z_i|\mathbf{y}_{obs}; \theta^r] = Pr(Z_i = 1|\mathbf{y}_{obs}; \theta^r)$

$Z=0,1$

$$\begin{aligned} E(Z) &= P(Z=1) \cdot 1 + 0 \cdot P(Z=0) \\ &= P(Z=1) \end{aligned}$$

- We will use Bayes' rule.

$$P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)} = \frac{P(A \cap B)}{P(B)}$$

$$\begin{aligned} Pr(Z_i = 1 | \mathbf{y}_{obs}; \theta^r) &= \frac{f(\mathbf{y}_{obs} | Z_i = 1; \theta^r) Pr(Z_i = 1)}{f(\mathbf{y}_{obs} | Z_i = 1; \theta^r) Pr(Z_i = 1) + f(\mathbf{y}_{obs} | Z_i = 0; \theta^r) Pr(Z_i = 0)} \\ &= \frac{\text{normal}(y_i; \mu_1, 1)^{\frac{1}{2}}}{\text{normal}(y_i; \mu_1, 1)^{\frac{1}{2}} + \text{normal}(y_i; \mu_0, 1)^{\frac{1}{2}}} \\ &= \frac{\text{normal}(y_i; \mu_1, 1)}{\text{normal}(y_i; \mu_1, 1) + \text{normal}(y_i; \mu_0, 1)} \\ &= p_i \end{aligned}$$

- So we have:

$$Q(\theta | \theta^{(r)}) = -\frac{1}{2} \sum_{i=1}^n (1 - p_i)(y_i - \mu_0)^2 - \frac{1}{2} \sum_{i=1}^n p_i(y_i - \mu_1)^2 + \text{constants}$$

2. Let's maximize  $Q(\boldsymbol{\theta}|\boldsymbol{\theta}^{(r)})$  with respect to  $\mu_0, \mu_1$ . We find:

$$\hat{\mu}_0^{(r+1)} = \frac{\sum_{i=1}^n (1 - p_i) y_i}{\sum_{i=1}^n (1 - p_i)}$$

$$\hat{\mu}_1^{(r+1)} = \frac{\sum_{i=1}^n p_i y_i}{\sum_{i=1}^n p_i}$$

- Recompute  $p_i$  with  $\hat{\mu}_0^{(r+1)}$  and  $\hat{\mu}_1^{(r+1)}$ . Thus iterate between the E and M steps till convergence.

```
## generate data from a bivariate normal:  
set.seed(1001)  
n <- 1000  
z <- rbinom(n, 1, 1/2)  
  
y <- rep(NA, n)  
y[z==0] <- rnorm(length(z[z==0]), -2, 1)  
y[z==1] <- rnorm(length(z[z==1]), 2, 1)
```

- Because we generated the data, we know  $z$ , so we know the MLEs:

$$\hat{\mu}_0 = \bar{y}|_{z=0}$$

```
mean(y[z==0])
```

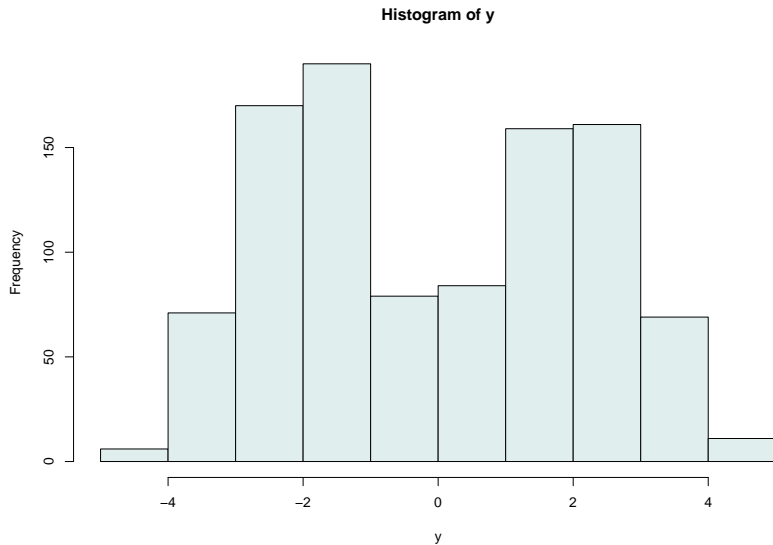
```
## [1] -1.957452
```

$$\hat{\mu}_1 = \bar{y}|_{z=1}$$

```
mean(y[z==1])
```

```
## [1] 1.997553
```

```
hist(y, col="azure2")
```



# E-M

```
## starting values
mu.0 <- -0.5
mu.1 <- 0.5

##
check <- 10
eps <- 1e-10

##
while(check > eps){

  # vector  $E[z/y]$  - E step
  rho <- dnorm(y, mu.1, 1)/(dnorm(y, mu.1, 1) + dnorm(y, mu.0, 1))

  # M step
  mu.0.new <- sum((1-rho)*y)/(sum((1-rho)))
  mu.1.new <- sum((rho)*y)/(sum((rho)))

  check <- sum( c(abs(mu.0.new - mu.0), abs(mu.1.new - mu.1)))
  mu.0 <- mu.0.new
  mu.1 <- mu.1.new
}

mu.0.hat <- mu.0
mu.1.hat <- mu.1
```



## MLEs based on E-M algorithm

```
mu.0.hat
```

```
## [1] -1.942764
```

```
mu.1.hat
```

```
## [1] 2.007483
```

# Invariance Property of MLEs

**Lemma 3.2:** Suppose that  $\theta$  and  $\eta$  represent two alternative parameterizations for some probability distribution and that  $\eta$  is a (1-1) function of  $\theta$ , so that we can write  $\eta = \mathbf{g}(\theta)$ ,  $\theta = \mathbf{h}(\eta)$  for appropriate functions  $\mathbf{g}(\cdot)$ ,  $\mathbf{h}(\cdot)$ .

- If  $\hat{\theta}$  is the MLE of  $\theta$  then  $\hat{\eta} = \mathbf{g}(\hat{\theta})$  is the MLE for  $\eta$
- If the mapping is (1-1) we simply note:

$$\eta = \tau(\theta) \rightarrow \tau^{-1}(\eta) = \theta$$

# Invariance Property of MLEs

- Define our likelihood based on the reparameterization ( $\theta = \tau^{-1}(\eta)$ ):

$$L^*(\eta; \mathbf{x}) = \prod_{i=1}^n f(x_i; \tau^{-1}(\eta)) = L(\tau^{-1}(\eta); \mathbf{x}) = L(\theta; \mathbf{x})$$

- We find the supremum of likelihood

$$\overset{\text{max}}{\sup}_{\eta} L^*(\eta; \mathbf{x}) = \overset{\text{max}}{\sup}_{\eta} L(\tau^{-1}(\eta); \mathbf{x}) = \overset{\text{max}}{\sup}_{\theta} L(\theta; \mathbf{x})$$

to see that the maximum of  $L^*(\eta; \mathbf{x})$  is when  $\eta = \tau(\theta) = \tau(\hat{\theta})$ .

# Invariance Property of MLEs

$\hat{\theta}$  is the MLE for  $\theta$   
 $\tau = (\theta)^2 \Rightarrow \hat{\tau} = (\hat{\theta})^2$

- However, many functions of interest are not one-to-one:  $\theta \rightarrow \theta^2$ .
- We proceed by defining the **induced likelihood function** of  $L^*$  for  $\tau(\theta)$

$$L^*(\eta; \mathbf{x}) = \sup_{\theta: \tau(\theta) = \eta} L(\theta; \mathbf{x})$$

$\eta = \tau(\theta)$

- The value  $\hat{\eta}$  that maximizes  $L^*(\eta; \mathbf{x})$  will be called the MLE of  $\eta$ .

not 1-1.

$$\tau(-5) = \tau(5) = 10$$

# Invariance Property of MLEs

**Proof:** Let  $\hat{\eta}$  denote the value that maximizes  $L^*(\eta; \mathbf{x})$ . Let's show for all values of  $\eta$  that

$$L^*(\hat{\eta} = \tau(\hat{\theta}); \mathbf{x}) \geq L^*(\eta; \mathbf{x}) \quad L(\hat{\theta}) \geq L(\theta)$$

• Steve Stern's note  
• Casella & Berger

$$\begin{aligned} \underbrace{L^*(\eta; \mathbf{x})}_{\text{constrained}} &= \sup_{\theta: \tau(\theta) = \eta} L(\theta; \mathbf{x}) \\ &\leq \sup_{\theta} L(\theta; \mathbf{x}) \quad \uparrow \text{unconstrained} \\ &= L(\hat{\theta}; \mathbf{x}) \\ &= \sup_{\theta: \tau(\theta) = \tau(\hat{\theta})} L(\theta; \mathbf{x}) \\ &= L^*(\tau(\hat{\theta}); \mathbf{x}) \\ &\quad \hat{\eta} \parallel \end{aligned}$$

# Invariance Property of MLEs

Eg. Normal:  $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} \text{normal}(\mu, \sigma^2)$ .

- If we want the MLE of  $\mu^2$  it is  $\widehat{\mu^2} = (\hat{\mu})^2$ .
- If we want the MLE of  $\sigma$  it is  $\hat{\sigma} = \sqrt{\hat{\sigma}^2}$ .
- This tends to be helpful in a computational sense as well (we can remove bounds on parameters):

$$\hat{\sigma}^2 = \exp(\hat{\theta}) \quad -\infty < \theta < \infty$$