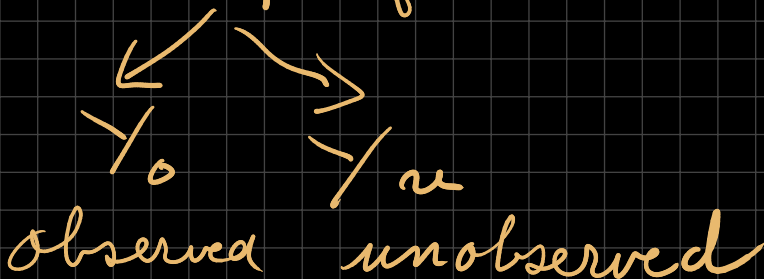


# Expectation - Maximization Algorithm

Art Dempster - Nan Laird - Donald Rubin '77

## Applications:

- Missing data:  $Y_f$  full data



"Easy" model for  $f(Y_f | \underline{\theta})$

But complex for  $f(Y_o | \underline{\theta}) = \int f(Y_o, Y_u | \underline{\theta}) dY_u$

Each step has 2 sub-steps:

Given value  $\tilde{\theta}_t$  from  $t^{\text{th}}$  step:

1) Expectation step:

$$\text{Use } l(y_f | \tilde{\theta}) = \log f(y_f | \tilde{\theta})$$

to work out

$$E(l(y_f | \tilde{\theta}) | y_o, \tilde{\theta}_t)$$

free parameter

parameter for E

$$= Q(\tilde{\theta} | \tilde{\theta}_t)$$

Note:

This is E of  
a function

of  $\tilde{\theta}$  given a value of  $\tilde{\theta}_t$

2) Maximize  $Q(\underline{\theta}, \underline{\theta}_t)$   
with respect to  $\underline{\theta}$

Let  $\underline{\theta}_{t+1} = \underset{\underline{\theta}}{\operatorname{argmax}} Q(\underline{\theta}, \underline{\theta}_t)$

---

Repeat 2 steps until convergence.

1)  $|\underline{\theta}_{t+1} - \underline{\theta}_t|$  small

2)  $|Q(\underline{\theta}_{t+1}, \underline{\theta}_t) - Q(\underline{\theta}_t, \underline{\theta}_t)|$  small

- In many applications,  $Y_u$  is not "missing data" but "random unknown parameters" as in multilevel or longitudinal modeling
- In some MLM packages (e.g. nlme), the E-M algorithm is used "under the hood".
- It's a goto method for many complex problems.
- It is generally not easy to apply but once you have an algorithm for a kind of problem then many people use it without knowing.

Example:

WARNING: LONG & TEDIOUS  
BUT INTERESTING (I THINK)  
AND GOOD PRACTICE

- Suppose you are sampling from

$$\begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \sim N_2(\underline{\mu}, \Sigma)$$

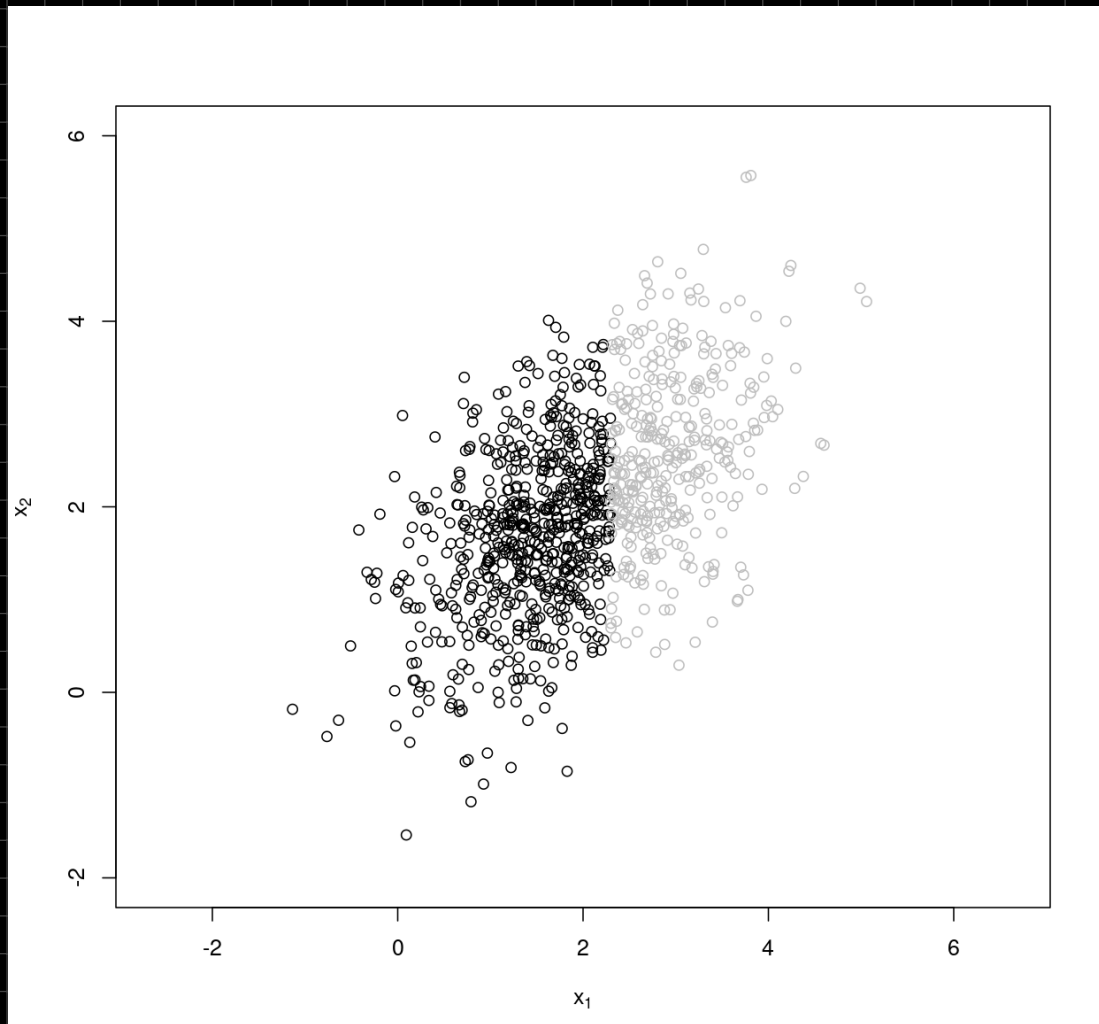
Simulation using

$$n = 1,000$$

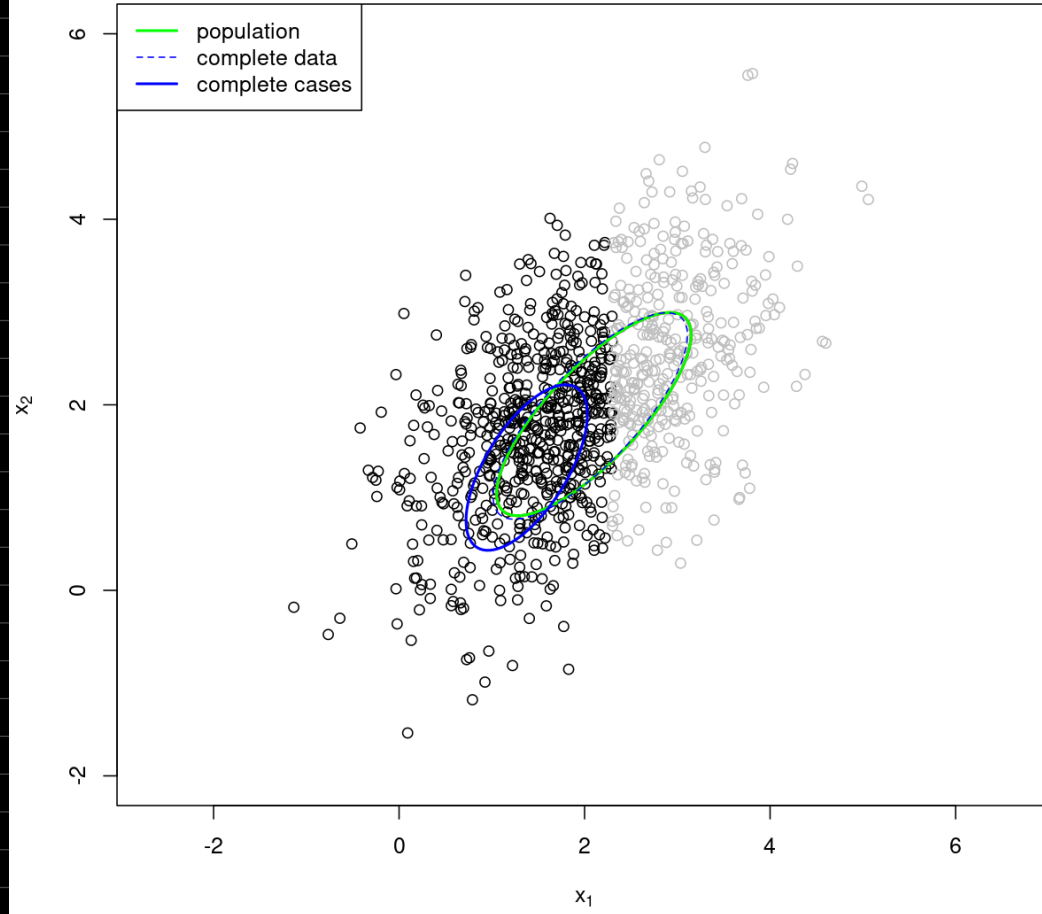
$$\underline{\mu} = \begin{pmatrix} 2.1 \\ 1.9 \end{pmatrix} \quad \Sigma = \begin{pmatrix} 1.1 & .9 \\ .9 & 1.2 \end{pmatrix}$$

"true" but unknown

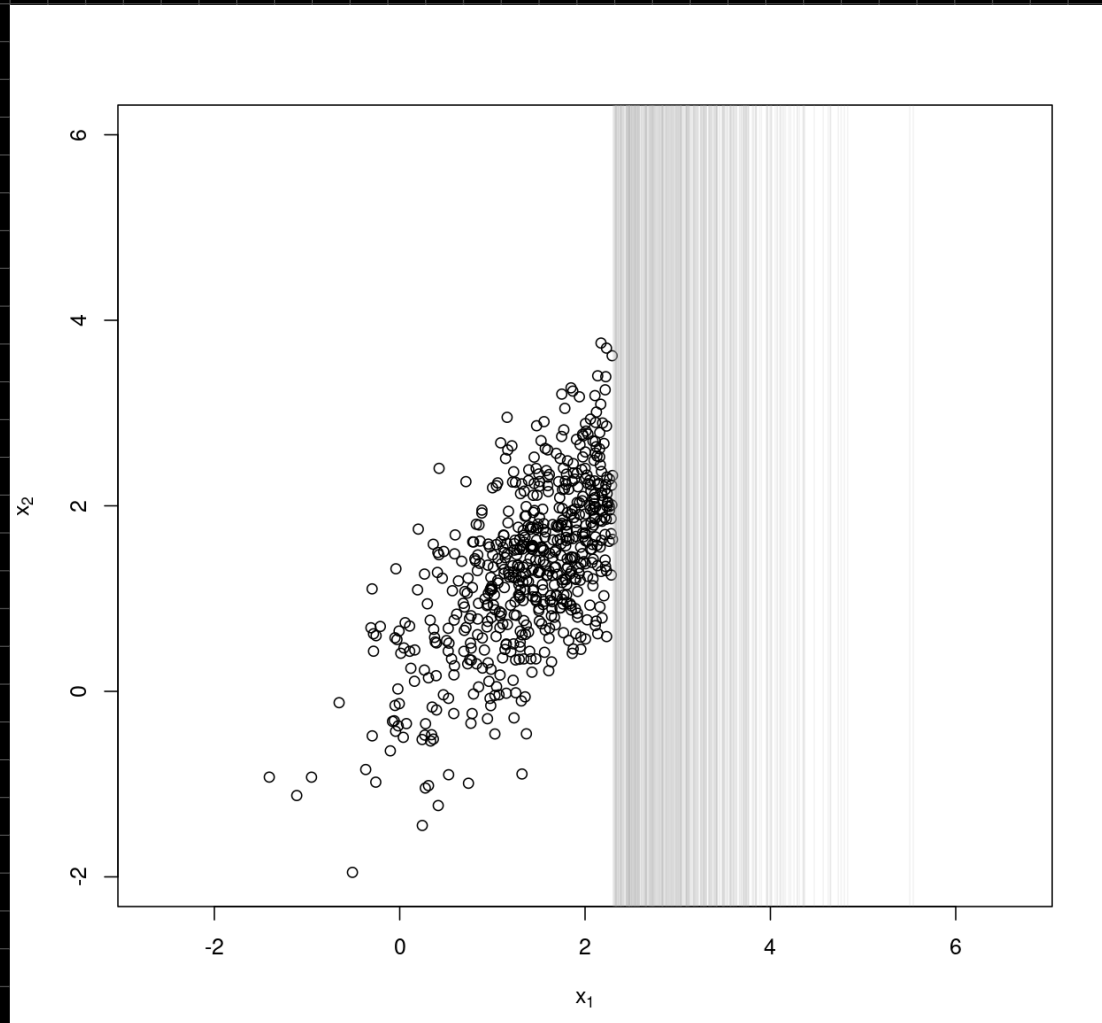
But  $x_2$  not observed if  $x_1 > 2.3$ .



$x_2$  missing  
based on  
observed  
value of  $x_1$



The data we actually have:



Is there any  
hope for  
estimating  
 $\mu_2$ ?

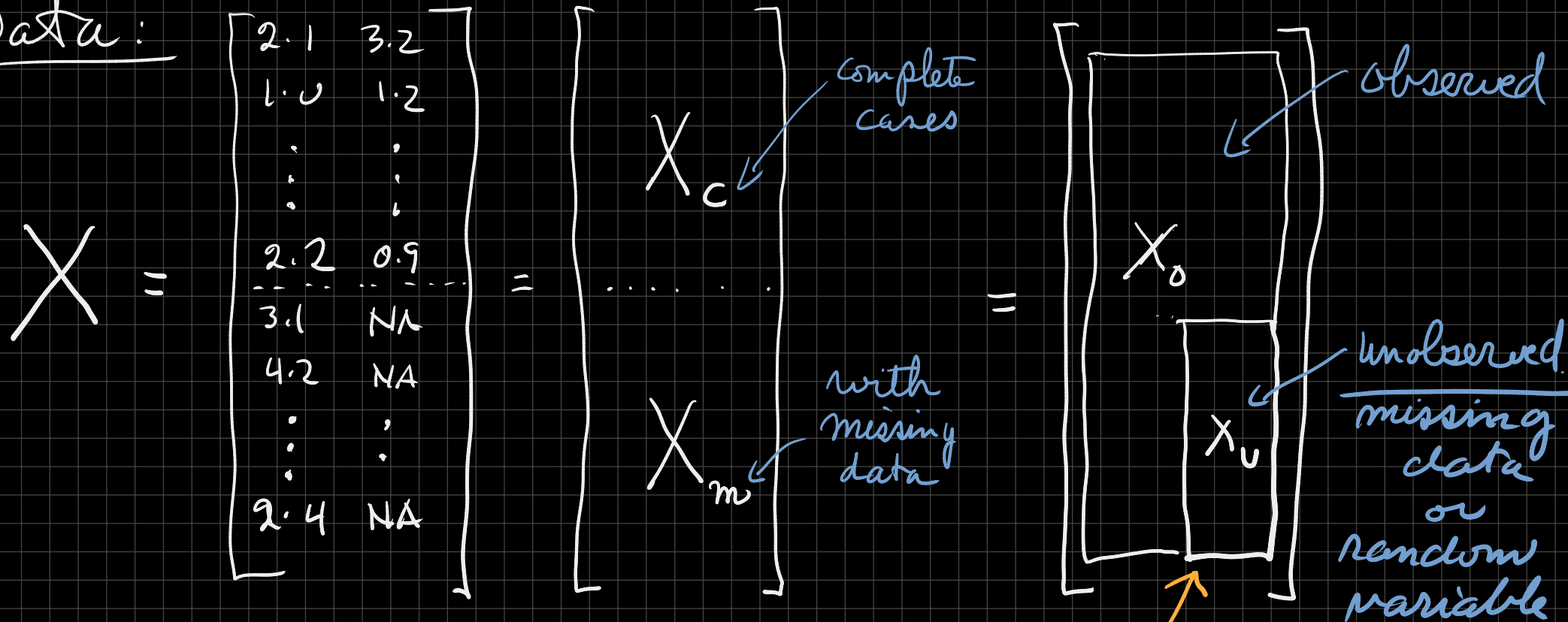
$\mu_1$  &  $\sigma_1^2 = \sigma_{11}$  easy.

What about  
 $\sigma_2^2 = \sigma_{22}$ ?

and  $\rho_{12} = \sigma_{12} / \sqrt{\sigma_{11} \sigma_{22}}$



Data:



If  $x_{ij}$  is missing,

we think of  $x_{ij}$  as a random variable.

Since  $X_u$  is unknown we think of it as a random vector.

# E-M Analysis:

Complete data model: (detail to remind some neat methods)

Of  $\underline{x} \sim N_p(\underline{\mu}, \Sigma)$

$$f(\underline{x} | \underline{\mu}, \Sigma) = \frac{1}{(2\pi)^{p/2}} \frac{1}{|\Sigma|^{1/2}} \exp\left\{-\frac{1}{2}(\underline{x} - \underline{\mu})^T \Sigma^{-1}(\underline{x} - \underline{\mu})\right\}$$

Sample of  $n$ :  $f(\underline{x}_1, \underline{x}_2, \dots, \underline{x}_n | \underline{\mu}, \Sigma)$

$$= \frac{1}{(2\pi)^{np/2}} \frac{1}{|\Sigma|^{n/2}} \exp\left\{-\frac{1}{2} \sum_{i=1}^n (\underline{x}_i - \underline{\mu})^T \Sigma^{-1} (\underline{x}_i - \underline{\mu})\right\}$$

Playing with

$$\sum_{i=1}^n (\underline{x}_i - \underline{\mu})^T \Sigma^{-1} (\underline{x}_i - \underline{\mu})$$

$$= \sum_{i=1}^n \text{tr} \left( (\underline{x}_i - \underline{\mu})^T \Sigma^{-1} (\underline{x}_i - \underline{\mu}) \right)$$

$$= \sum_{i=1}^n \text{tr} \left( \Sigma^{-1} (\underline{x}_i - \underline{\mu}) (\underline{x}_i - \underline{\mu})^T \right)$$

$$= \text{tr} \left( \Sigma^{-1} \sum_{i=1}^n (\underline{x}_i - \underline{\mu}) (\underline{x}_i - \underline{\mu})^T \right)$$

$$= \text{tr} \left( \Sigma^{-1} \sum_{i=1}^n (\underline{x}_i \underline{x}_i^T - \underline{x}_i \underline{\mu}^T - \underline{\mu} \underline{x}_i^T + \underline{\mu} \underline{\mu}^T) \right)$$

TRACE

• tr is sum of diagonal elts

•  $\text{tr}(AB) = \text{tr}(BA)$

•  $\text{tr}(A) = \text{tr}(A^T)$

•  $\text{tr}(A) = a$  if

$A = [a]$  is  $1 \times 1$

•  $\text{tr}(A+B) = \text{tr}(A) + \text{tr}(B)$

by linearity

$$= \frac{1}{n} \sum^{-1} \left( X^T X - n \bar{X} \bar{\mu}^T - n \bar{\mu} \bar{X}^T + n \bar{\mu} \bar{\mu}^T \right)$$

$$= n \bar{\mu}^T \sum^{-1} \bar{\mu} + \frac{1}{n} \sum^{-1} \left( X^T X - \left( \sum_{i=1}^n x_i \right) \bar{\mu}^T - \bar{\mu} \left( \sum_{i=1}^n x_i \right)^T \right)$$

focus on this

split  $X_c$  &  $X_m$   
 complete cases      cases with missing  $X_2$

$$X_c^T X_c - \left( \sum_{i \in X_c} x_i \right) \bar{\mu}^T - \bar{\mu} \left( \sum_{i \in X_c} x_i \right)^T$$

$$+ X_m^T X_m - \sum_{i \in X_m} x_i \bar{\mu}^T - \bar{\mu} \left( \sum_{i \in X_m} x_i \right)^T$$



We need to find  $E(* | X_{\text{observed}}, \mu_t, \Sigma_t)$

Why? Because  $\tilde{\ell}(X | \mu, \Sigma)$  is linear in  $(*)$

$$\begin{aligned} \text{So } Q(\tilde{\theta} | \tilde{\theta}_t) &= E(\tilde{\ell}(X | \tilde{\theta}) | X_0, \tilde{\theta}_t) \\ &= \tilde{\ell}(E(*) | \tilde{\theta}) | X_0, \tilde{\theta}_t \end{aligned}$$

$$\text{i.e. } E(\tilde{\ell}(*)) = \tilde{\ell}(E(*))$$

which works if  $\tilde{\ell}$  is linear in  $(*)$   
and  $(*)$  includes all random variables.

# Plan of attack:

note that  $\mu$  is an argument of  $l$   
Not a parameter for  $E$

This part is constant given  $X_0$   
why?

$$X_c' X_c - \left( \sum_{i \in X_c} x_i \right) \mu^T - \mu \left( \sum_{i \in X_c} x_i \right)^T$$

$$+ \underbrace{X_m' X_m}_{\text{ditto}} - \left[ \sum_{i \in X_m} x_i \right] \mu^T - \mu \left[ \sum_{i \in X_m} x_i \right]^T$$

This part has both constants ( $x_1$ 's) and random variables ( $x_2$ 's) given  $X_0$

We need  $E(\downarrow \downarrow | X_0, \mu^{(+)}, \Sigma^{(+)})$

Start with:

$$\textcircled{1} E\left(\sum_{i \in X_m} x_i \mid X_0, \mu^{(t)}, \Sigma^{(t)}\right)$$

$$= \sum_{i \in X_m} E(x_i \mid X_0, \mu^{(t)}, \Sigma^{(t)})$$

$$= \sum_{i \in X_m} \left( \gamma_{1i} \mu_2^{(t)} + \frac{\sigma_{12}^{(t)}}{\sigma_{11}^{(t)}} (x_{1i} - \mu_1^{(t)}) \right)$$

$$= \left( \sum_{i \in X_m} \gamma_{1i} \mu_2^{(t)} + \frac{\sigma_{12}^{(t)}}{\sigma_{11}^{(t)}} \left( \sum_{i \in X_m} x_{1i} - n_m \mu_1^{(t)} \right) \right)$$

$$= \left( S_{1m} + n_m \mu_2^{(t)} + \frac{\sigma_{12}^{(t)}}{\sigma_{11}^{(t)}} (S_{1m} - n_m \mu_1^{(t)}) \right)$$

$$(2) E \left( \sum_{i \in X_m} x_{1i} x_{2i} \mid X_0, \mu^{(t)}, \Sigma^{(t)} \right)$$

$$= \sum_{i \in X_m} x_{1i} E \left( x_{2i} \mid X_0, \mu^{(t)}, \Sigma^{(t)} \right)$$

$$= \sum_{i \in X_m} x_{1i} \left( \mu_2^{(t)} + \frac{\sigma_{12}^{(t)}}{\sigma_{11}^{(t)}} (x_{1i} - \mu_1^{(t)}) \right)$$



$$= \mu_2^{(t)} S_{1m} + \frac{\sigma_{12}^{(t)}}{\sigma_{11}^{(t)}} \left( \sum x_{1i}^2 - \mu_1^{(t)} S_{1m} \right)$$

$$(3) \quad E \left( \sum x_{2i}^2 \mid X_0, \mu^{(t)}, \Sigma^{(t)} \right)$$

$$= \sum_{i=1}^m \left\{ \left[ E(x_{2i} \mid X_0, \mu^{(t)}, \Sigma^{(t)}) \right]^2 + \text{Var}(x_{2i} \mid X_0, \mu^{(t)}, \Sigma^{(t)}) \right\}$$

$$= \sum_{i=1}^m \left\{ \left( \mu_2^{(t)} + \frac{\sigma_{12}^{(t)}}{\sigma_{11}^{(t)}} (x_{1i} - \mu_1^{(t)}) \right)^2 + \sigma_{22}^{(t)} - \frac{(\sigma_{12}^{(t)})^2}{\sigma_{11}^{(t)}} \right\}$$

As we get ready to write an algorithm to do this, we try to identify intermediate quantities to compute.

An important one is the vector

$$E_{i \in X_m}(\tilde{x}_2 | \tilde{x}_1) = \mu_2^t + \frac{\sigma_{12}^t}{\sigma_{11}^t} (\tilde{x}_1 - \mu_1^t) = \rho_{2.1}$$

"R notation"

So back to the original

$$X_c' X_c - \left( \sum_{i \in X_c} \alpha_i \right) \mu^T - \mu \left( \sum_{i \in X_c} \alpha_i \right)^T$$
$$+ E \left( X_m' X_m - \sum_{i \in X_m} \alpha_i \mu^T - \mu \left( \sum_{i \in X_m} \alpha_i \right)^T \mid X_0, \mu^t, \Sigma^t \right)$$

$$= X_c' X_c + S_m - (S_c + S_m) \mu^T - \mu (S_c + S_m)^T$$

where  $S_m = E(\cancel{X_c' X_c} \mid X_0, \mu^t, \Sigma^t)$  this should be  $X_m' X_m$

$$S_c = \sum_{i \in X_c} \alpha_i \quad S_m = E \left( \sum_{i \in X_m} \alpha_i \right)$$

$$S_m = \left[ \begin{array}{l} \sum_{i \in X_m} x_{1i}^2 \\ \sum_{i \in X_m} x_{1i} e_{2.1,i} \\ \sum_{i \in X_m} e_{2.1,i}^2 + n_m \left( \sigma_{22}^{(t)} - \frac{(\sigma_{12}^{(t)})^2}{\sigma_{11}^{(t)}} \right) \end{array} \right]$$


---

$$S_m \approx \sum_{i \in X_m} \left( x_{1i} \left( \mu_2^t + \frac{\sigma_{12}^t}{\sigma_{11}^t} (x_{1i} - \mu_1^t) \right) \right)$$

$$\text{Let } P = X_c^T X_c + S_m$$

$$\tilde{S} = \tilde{S}_c + \tilde{S}_m$$

Then  $Q(\theta, \theta^t)$  has the form:

$$\dots \text{tr} \Sigma^{-1} (P - \tilde{S} \tilde{\mu}^T - \tilde{\mu} \tilde{S}^T) \dots$$

In full:

$$Q(\theta, \theta^t) = k - \frac{n}{2} \log |\Sigma| - \frac{n}{2} \tilde{\mu}^T \Sigma^{-1} \tilde{\mu} \\ - \frac{1}{2} \text{tr} \Sigma^{-1} (P - \tilde{S} \tilde{\mu}^T - \tilde{\mu} \tilde{S}^T)$$

- This is the log-likelihood for a  $N_2(\mu, \Sigma)$  if  $P$  is the SSP matrix and  $\tilde{s}$  is the sum of  $\tilde{x}$ 's.
- Which we know is maximized if

$$\underline{\mu} = \hat{\underline{\mu}} = \frac{1}{n} \tilde{s}$$

$$\underline{\Sigma} = \hat{\underline{\Sigma}} = (P - \frac{1}{n} \tilde{s} \tilde{s}^T) / n$$

So the M-step is easy!

$$\mu^{t+1} = \hat{\underline{\mu}} = \frac{1}{n} \tilde{s} \quad \text{and} \quad \Sigma^{t+1} = \hat{\underline{\Sigma}}$$

