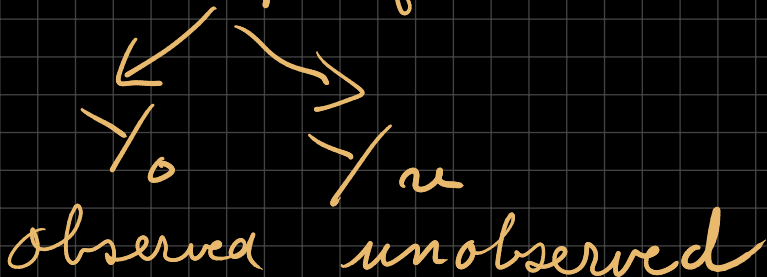


Expectation - Maximization Algorithm

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^{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

Note:

This is E of
a function

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

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

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

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

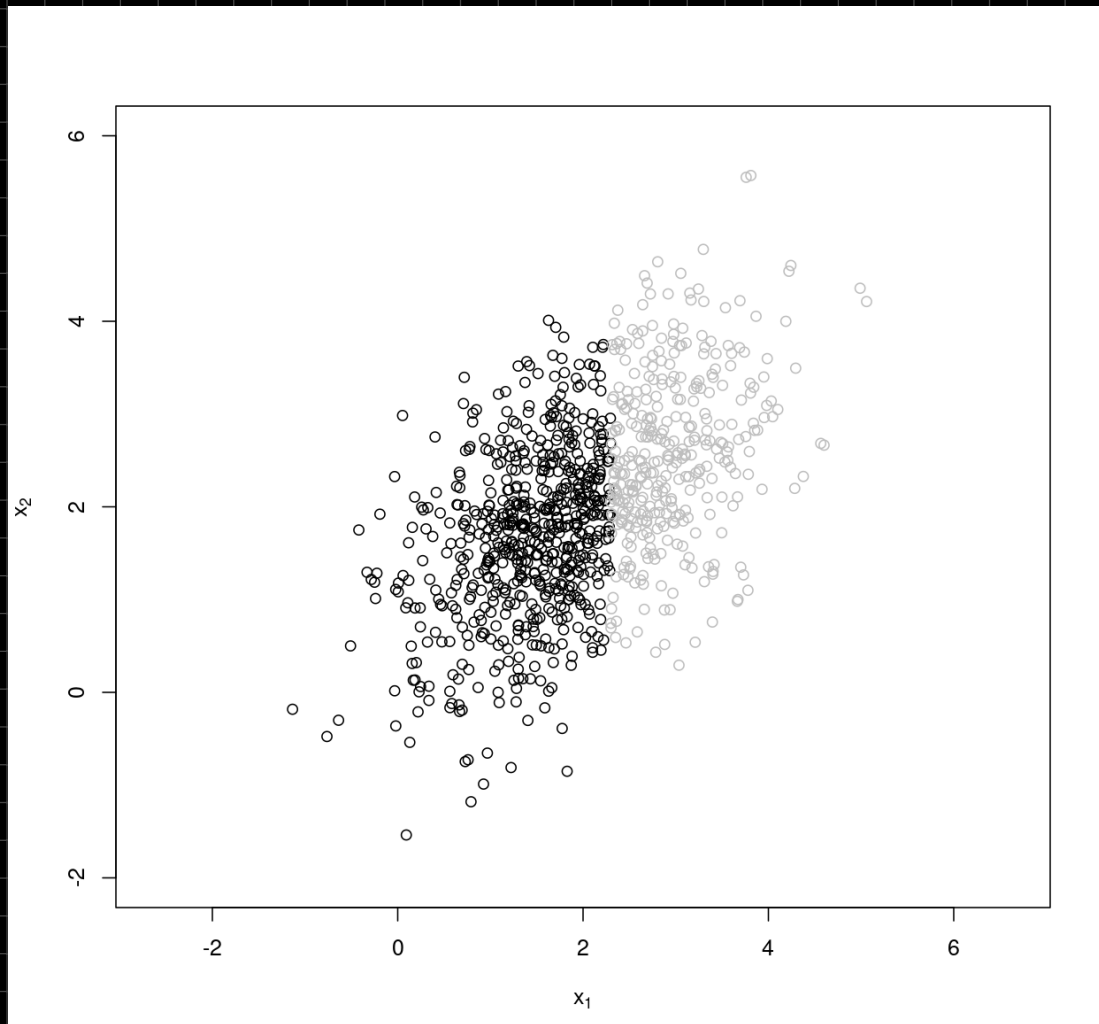
Simulation using

$n = 1,000$

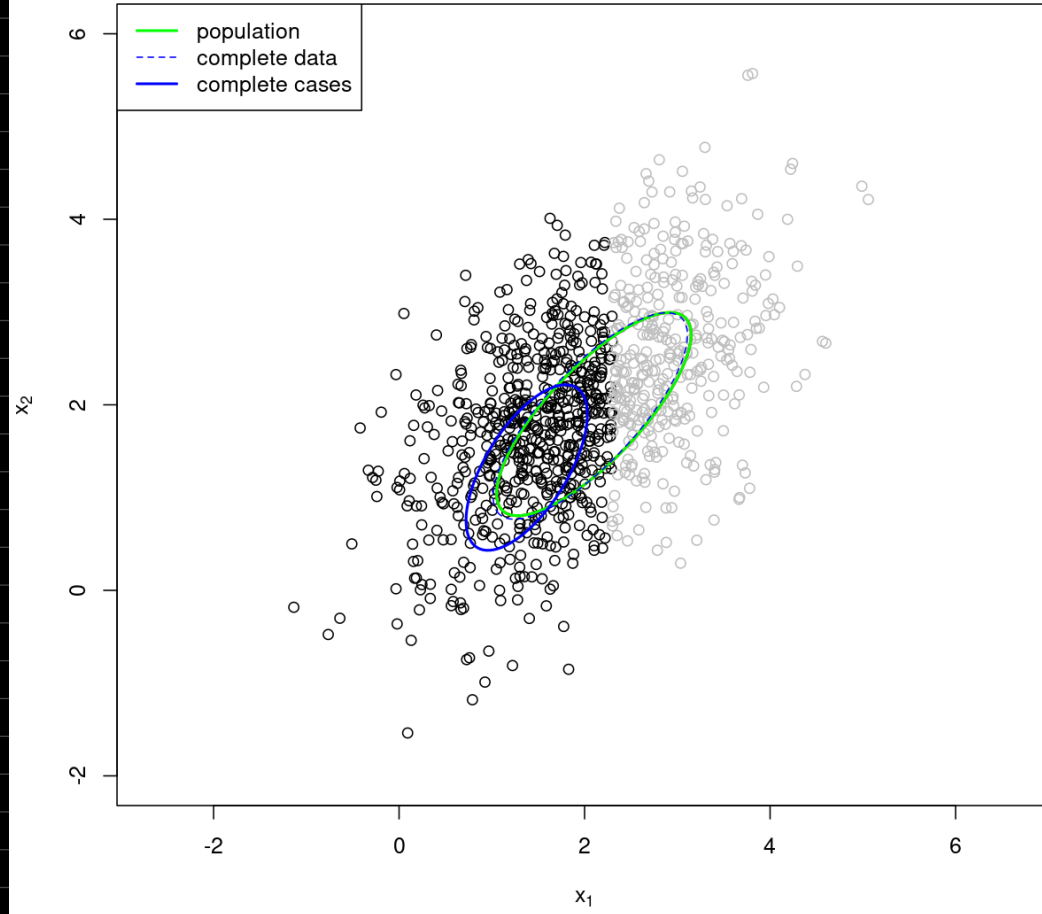
$$\underline{\mu} = \begin{pmatrix} 2.1 \\ 1.9 \end{pmatrix} \quad \Sigma = \begin{pmatrix} 1.1 & .9 \\ .9 & 1.2 \end{pmatrix}$$

"true" but unknown

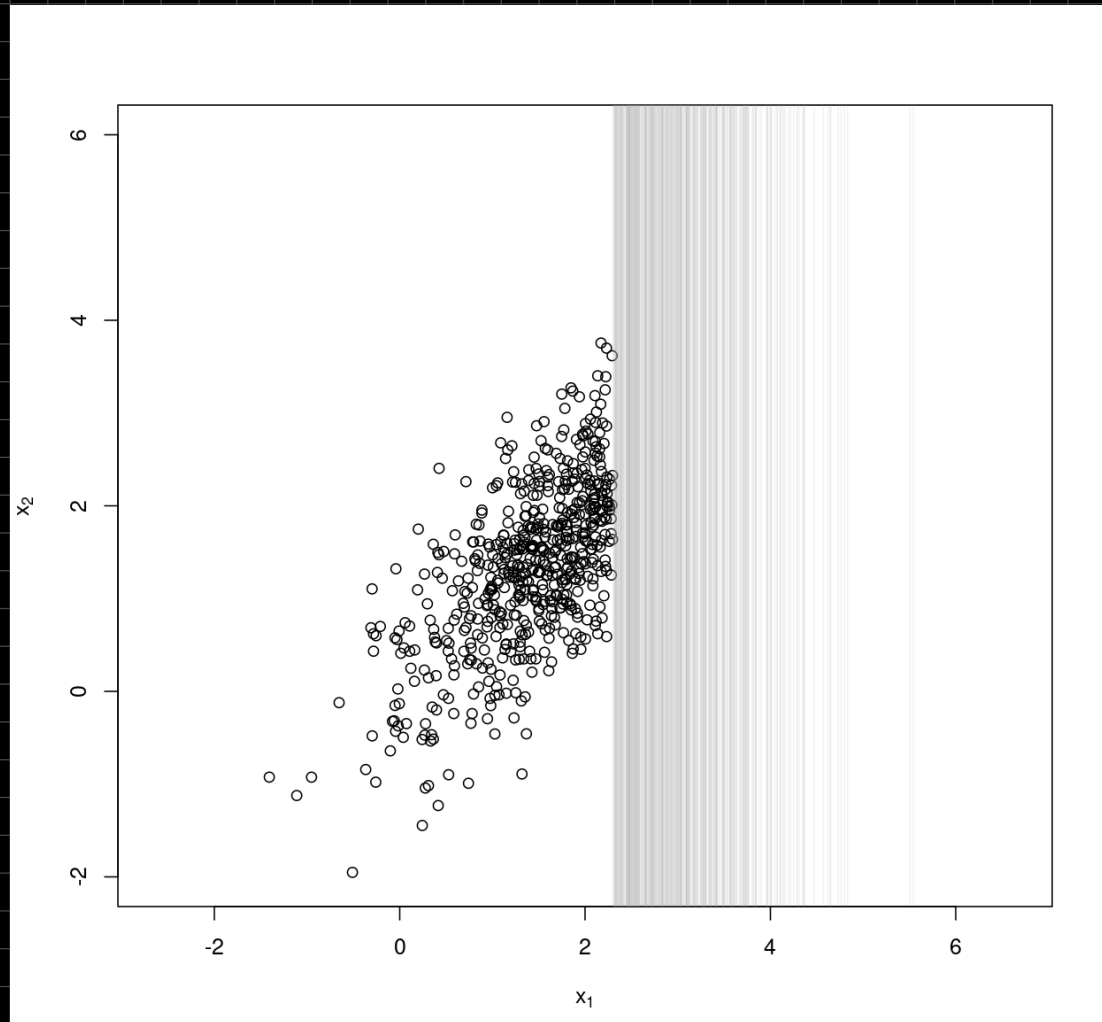
But x_2 not observed if $x_1 > 2.3$.



x_2 missing
based on
observed
value of x_1



The data we actually have:



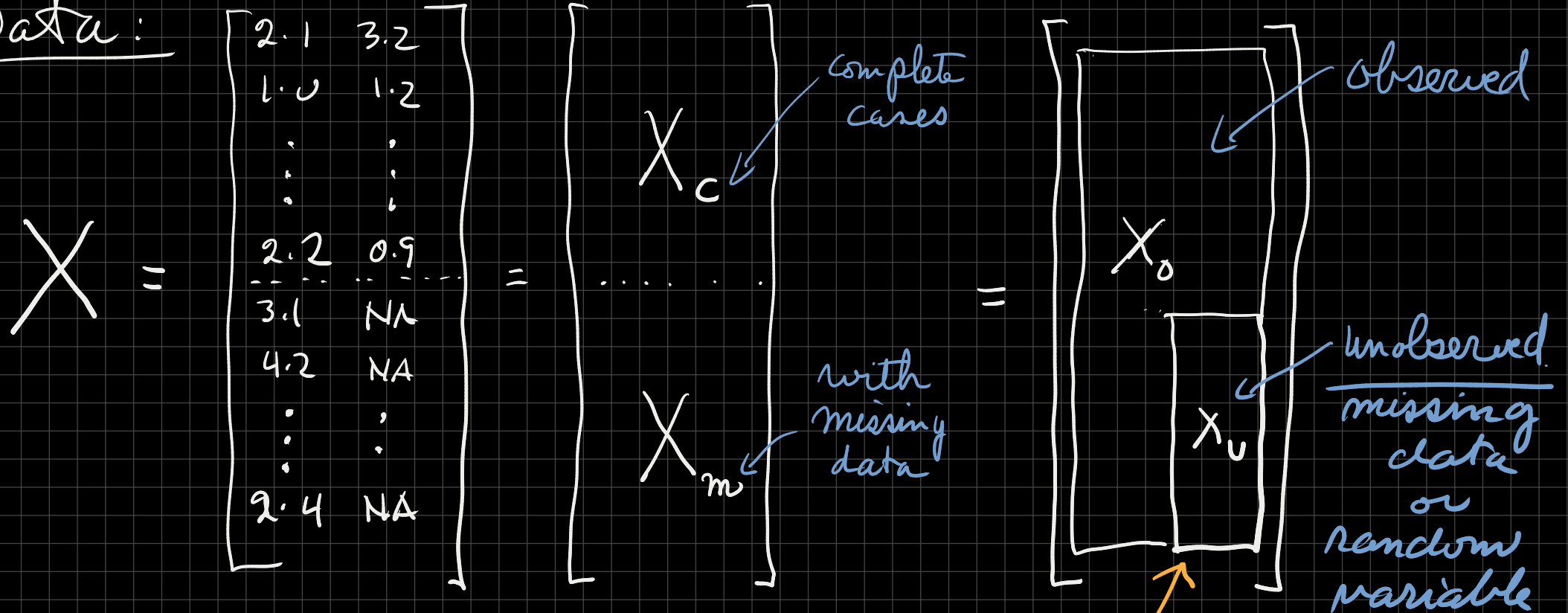
Is there any
hope for
estimating
 μ_2 ?

μ_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)

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

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

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 1×1

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

by linearity

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

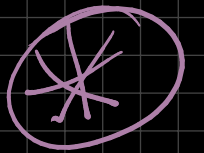
$$= n \bar{\mu}^T \sum^{-1} \bar{\mu} + \frac{1}{n} \sum^{-1} \left(X^T X - \left(\sum_{i=1}^n x_i \right) \bar{\mu}^T - \bar{\mu} \left(\sum_{i=1}^n x_i \right)^T \right)$$

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 μ 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:

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

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

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

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

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

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

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

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

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

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

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

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

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

$$E(x_2 | x_1) = \mu_2 + \frac{\sigma_{12}^t}{\sigma_{11}^t} (x_1 - \mu_1)$$

"R notation"

$= \tilde{r}_{2.1}$

Next $\tilde{x}_1 \cdot \tilde{v} = \tilde{x}_v$

$\tilde{v}_2 = \tilde{v} \times \tilde{v}$
 ↑
 element-wise

$Cvar = \sigma_{22}^t - \frac{(\sigma_{12}^t)^2}{\sigma_{11}^t}$

So: $E(X_m^T X_m | X_0, \mu^{(t)}, \Sigma^{(t)})$

$$= \left[\begin{array}{l} \text{sum}(\tilde{x}_1 \times \tilde{x}_1) \quad \cdot \quad \text{sum}(\tilde{x}_v) \\ \text{sum}(\tilde{x}_v) \quad \quad \quad \text{sum}(\tilde{v}_2) + n_m \times Cvar \end{array} \right]$$

So back to the original

$$\begin{aligned} & X_c' X_c - \left(\sum_{i \in X_c} \mathcal{X}_i \right) \mu^T - \mu \left(\sum_{i \in X_c} \mathcal{X}_i \right)^T \\ & + E \left(X_m' X_m - \sum_{i \in X_m} \mathcal{X}_i \mu^T - \mu \left(\sum_{i \in X_m} \mathcal{X}_i \right)^T \mid X_0, \mu^t, \Sigma^t \right) \\ & = X_c' X_c + S_m - (S_c + S_m) \mu^T - \mu (S_c + S_m)^T \\ & \text{where } S_m = E(X_c' X_c \mid X_0, \mu^t, \Sigma^t) \\ & \quad S_c = \sum_{i \in X_c} \mathcal{X}_i \quad S_m = E \left(\sum_{i \in X_m} \mathcal{X}_i \right) \end{aligned}$$

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

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

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

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

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

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

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

So the M-step is easy!

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

$$f(x_1, \dots, x_n | \mu, \Sigma) = \frac{1}{(2\pi)^{n/2}} \frac{1}{|\Sigma|^{n/2}} \left\{ -\frac{1}{2} \sum_{i=1}^n (x_i - \mu)^T \Sigma^{-1} (x_i - \mu) \right\}$$

$$\sum (x_i - \mu)^T \Sigma^{-1} (x_i - \mu)$$

$$\sum (x_i - \bar{x} + \bar{x} - \mu)^T \Sigma^{-1} (x_i - \bar{x} + \bar{x} - \mu)$$

$$= \sum (x_i - \bar{x})^T \Sigma^{-1} (x_i - \bar{x}) + \sum (x_i - \bar{x})^T \Sigma^{-1} (\bar{x} - \mu)$$

$$+ \sum (\bar{x} - \mu)^T \Sigma^{-1} (x_i - \bar{x}) + \sum (\bar{x} - \mu)^T \Sigma^{-1} (\bar{x} - \mu)$$

$$= \frac{1}{n} \sum_{i=1}^n \Sigma^{-1} (x_i - \bar{x})(x_i - \bar{x})^T + \underline{0} + \underline{0} + \frac{1}{n} \sum_{i=1}^n \Sigma^{-1} (\bar{x} - \mu)(\bar{x} - \mu)^T$$

$$= n t_n \Sigma^{-1} \hat{\Sigma} + m t_n \Sigma^{-1} (\bar{x} - \mu)(\bar{z} - \mu)^T$$

$$= n t_n \Sigma^{-1} \left\{ \hat{\Sigma} + (\bar{x} - \mu)(\bar{z} - \mu)^T \right\}$$

$$f = \frac{1}{(2\pi)^{p/2}} \frac{1}{|\Sigma|^{n/2}} \exp \left\{ -\frac{n}{2} t_n \Sigma^{-1} \left\{ \hat{\Sigma} + (\bar{x} - \mu)(\bar{z} - \mu)^T \right\} \right\}$$

$$\ln f = a - \frac{n}{2} \ln |\Sigma| - \frac{n}{2} t_n \Sigma^{-1} \left\{ \hat{\Sigma} + (\bar{x} - \mu)(\bar{z} - \mu)^T \right\}$$

Suppose we're missing half the x_2 - for the largest x_1 's - we have the x_1 's.

using imaginary data

$$\ln f = a - \frac{n}{2} \ln |\Sigma| - \frac{n}{2} \text{tr} \Sigma^{-1}$$

Recall Sample of n from $N(\mu, \sigma)$

$$Q(y | \mu, \sigma) = \log \left(\frac{1}{\sqrt{2\pi}} \cdot \frac{1}{\sigma} \exp \left\{ -\frac{1}{2} \left(\frac{y-\mu}{\sigma} \right)^2 \right\} \right)$$

$$Q(\underset{\in \mathbb{R}^p}{y} | \mu, \Sigma) = \log \left(\frac{1}{(2\pi)^{p/2} |\Sigma|^{1/2}} \exp \left\{ -\frac{1}{2} (y-\mu)' \Sigma^{-1} (y-\mu) \right\} \right)$$

$$= k - \frac{1}{2} \log |\Sigma| - \frac{1}{2} (x-\mu)' \Sigma^{-1} (x-\mu)$$

$$l(\underline{y}_1 \dots \underline{y}_n | \mu, \Sigma) =$$

$$= b' - \frac{n}{2} \log |\Sigma| - \frac{n}{2} t \left(\Sigma^{-1} \left(\hat{\Sigma} + (\bar{\underline{y}} - \underline{\mu})(\bar{\underline{y}} - \underline{\mu})^T \right) \right)$$

$$n_c = n \text{ complete} \quad n_m = n \text{ missing } y_2$$

$$n = n_c + n_m$$

$$= b' - \frac{n_c}{2} \log |\Sigma| - \frac{n_c}{2} t \left(\Sigma^{-1} \left(\Sigma \right) \right)$$

$$E(l(\underline{y}_1, \underline{y}_2 | \mu, \Sigma) | \underline{x}_0, \underline{\mu}, \Sigma)$$

Because l is an exp. family of y is the canonical variable

l is linear in η .

$$E(l(\eta, \Sigma, \mu | z)) = l(E(\eta | z))$$

$$E(\eta | \eta_0, \Sigma, \mu)$$

For an exponential family:

$$l(\eta | \theta) = \{ \eta(\theta) \cdot T(\eta) - A(\theta) + B(\eta) \}$$

$$E(l(\eta | \theta) | \eta_0, \theta_t)$$

at \uparrow

$$= E(\eta(\theta) \cdot T(y) - A(\theta) + b \mid y_0, \theta_t)$$

$$= \eta(\theta) E(T(y) \mid y_0, \theta_t) - A(\theta)$$

Need $T(y)$

$$f(y \mid \mu, \Sigma) = \frac{1}{(2\pi)^{p/k}} \frac{1}{|\Sigma|^{1/2}} \exp\left\{-\frac{1}{2}(y-\mu)' \Sigma^{-1}(y-\mu)\right\}$$

$$\mathcal{L}f = k - \frac{1}{2} \ln|\Sigma| - \frac{1}{2} (y-\mu)' \Sigma^{-1} (y-\mu)$$

$$= k - \frac{1}{2} \ln|\Sigma| - \frac{1}{2} \text{tr} \Sigma^{-1} (yy' - 2\mu y' + \mu\mu')$$

$$= k - \frac{1}{2} \ln |\Sigma| - \frac{1}{2} \text{tr} \left(\underbrace{\Sigma^{-1} y y^T}_{\Sigma, \mu} - 2 \underbrace{\Sigma^{-1} \mu y^T}_{\Sigma^{-1} \mu, y} + \Sigma^{-1} \mu \mu^T \right)$$

$$= k - \frac{1}{2} \ln |\Sigma| - \frac{1}{2} \text{tr} \Sigma^{-1} \mu \mu^T - \frac{1}{2} \text{tr} \Sigma^{-1} \begin{pmatrix} y y^T - 2 \mu y^T \\ \mu \mu^T \end{pmatrix}$$

When y_2 is missing we need

$$E(y_2 | y_1, \theta), \quad E(y_1 y_2 | y_1, \theta), \quad E(y_2^2 | y_1, \theta)$$

$$E(y_2 | y_1, \theta) = (y_1 - \mu_1) \frac{\sigma_{12}}{\sigma_{11}} + \mu_2$$

$$\begin{aligned} E(y_1, y_2 | y_1, \theta) &= y_1 \cdot E(y_2 | y_1, \theta) \\ &= y_1 (y_1 - \mu_1) \frac{\sigma_{12}}{\sigma_{11}} + \mu_2 \end{aligned}$$

$$E(y_2^2 | y_1, \theta)$$

$$= E(y_2 | y_1, \theta)^2 + \text{Var}(y_2 | y_1, \theta)$$

$$= \left(y_1 - \mu_1 \right)^2 \frac{\sigma_{12}^2}{\sigma_{11}^2} + \mu_2^2 + \sigma_{22} - \frac{\sigma_{12}^2}{\sigma_{11}}$$

Try it as a function of Y, Σ

where $Y = \begin{bmatrix} y_{i1} & y_{i2} \\ \vdots & \vdots \\ y_{n1} & y_{n2} \end{bmatrix}$ 1st obs
nth obs

Let $\underline{y}_i = [y_{i1}, y_{i2}]$ $\mu = (\mu_1, \mu_2)$

log lik. fn one observation:

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

$$\ln f = k - \frac{1}{2} \ln |\Sigma| - \frac{1}{2} \text{tr} \Sigma^{-1} (y - \mu)(y - \mu)^T$$

For a sample of n in Y

$$= k - \frac{n}{2} \ln |\Sigma| - \frac{n}{2} \text{tr} \Sigma^{-1} \underbrace{\sum_{i=1}^n (y_i - \mu)^T (y_i - \mu)}$$

$$\sum_{i=1}^n \left(\underbrace{y_i}_{\sim} \underbrace{y_i^T}_{\sim} - \underbrace{y_i}_{\sim} \underbrace{\mu^T}_{\sim} - \underbrace{\mu}_{\sim} \underbrace{y_i^T}_{\sim} + \underbrace{\mu}_{\sim} \underbrace{\mu^T}_{\sim} \right)$$

$$= Y^T Y - Y^T \mathbf{1} \mu^T - \mu \mathbf{1}^T Y + \mu \mu^T$$

$$= Y_0^T Y_0 + \underbrace{Y_m^T Y_m}_{\substack{- \mu \mathbf{1}^T Y_0 \\ - \mu \mathbf{1}^T Y_m}} - \underbrace{Y_0^T \mathbf{1} \mu^T}_{\substack{- \mu \mathbf{1}^T Y_0 \\ - \mu \mathbf{1}^T Y_m}} - \underbrace{Y_m^T \mathbf{1} \mu^T}_{\substack{- \mu \mathbf{1}^T Y_0 \\ - \mu \mathbf{1}^T Y_m}} + \mu \mu^T$$

Now $\underline{E(\mu \mathbf{1}^T Y_m)}$

need $E(Y_2 | \mu_1, \Sigma) = \mu_2 + \frac{\sigma_{12}}{\sigma_{11}} (z_1 - \mu_1)$

$$E\left(\begin{matrix} y_1 \\ y_2 \end{matrix} \middle| z_1, \mu_1, \Sigma\right) = \begin{pmatrix} y_1 \\ \mu_2 + \frac{\sigma_{12}}{\sigma_{11}} (z_1 - \mu_1) \end{pmatrix}$$

and summing

$$\left(\begin{matrix} \sum y_{i1} \\ n_m \mu_2 + \frac{\sigma_{12}}{\sigma_{11}} (\sum y_{i1} - n_m \mu_1) \end{matrix} \right)$$

$$\underline{E(m | Z, \mu_t, \Sigma_t)}$$

$$y_m^T y_m = \sum \left(\begin{array}{c} y_{1i}^2 \\ \frac{y_{1i} y_{2i}}{y_{2i}^2} \end{array} \right)$$

