

The Expectation-Maximization Algorithm

version: 2020-02-25 · 00:40:12

Motivation: Inference with Missing Data

► **Regression Model:** $y_i = \alpha x_i + \beta z_i + \sigma \varepsilon_i$, $\varepsilon_i \stackrel{\text{iid}}{\sim} \mathcal{N}(0, 1)$.

► Suppose some of the z_i are **missing**:

► Let $\delta_i = 0$ if z_i missing and $\delta_i = 1$ if z_i observed.

► Observed data:

$$\mathcal{D} = \begin{bmatrix} y_1 & x_1 & z_1 & \delta_1 = 1 \\ y_2 & x_2 & \text{NA} & \delta_2 = 0 \\ y_3 & x_3 & \text{NA} & \delta_3 = 0 \\ \vdots & \vdots & \vdots & \vdots \\ y_n & x_n & z_n & \delta_n = 1 \end{bmatrix}.$$

► **Problem:** How to estimate $\theta = (\alpha, \beta, \sigma)$ from \mathcal{D} ?

Inference with Missing Data

Solution 1: Use only complete observations $\mathcal{S}_1 = \{i : \delta_i = 1\}$.

Inference with Missing Data

Solution 1: Use only complete observations $\mathcal{S}_1 = \{i : \delta_i = 1\}$.

1. **Inefficient** (throws out data)

Inference with Missing Data

Solution 1: Use only complete observations $\mathcal{S}_1 = \{i : \delta_i = 1\}$.

1. **Inefficient** (throws out data)

2. Potentially **misleading**, as in the following example:

▶ $x, z \stackrel{\text{iid}}{\sim} \mathcal{N}(0, 1)$

▶ True parameters $\alpha = \beta = \sigma = 1$.

▶ Missing data mechanism: $P(\delta = 0 | y \leq 2) = 5\%$, $P(\delta = 0 | y > 2) = 90\%$
(overall 15% missing)

▶ Parameter estimates for $n = 10^6$:

	$\hat{\alpha}(\text{se})$	$\hat{\beta}(\text{se})$
No missing data ($\delta \equiv 1$)	.997(.002)	1.001(.002)
Using only \mathcal{S}_1 (85% of sample)	.901(.001)	.900(.001)

Inference with Missing Data

Solution 2: Maximize likelihood over θ and $\mathbf{z}_0 = \{z_i : \delta_i = 0\}$

$$(\hat{\theta}, \hat{\mathbf{z}}_0) = \arg \max_{(\theta, \mathbf{z}_0)} \left\{ -\frac{n}{2} \log \sigma^2 - \frac{1}{2} \sum_{i=1}^n \frac{(y_i - \alpha x_i - \beta z_i)^2}{\sigma^2} \right\}.$$

Inference with Missing Data

Solution 2: Maximize likelihood over θ and $\mathbf{z}_0 = \{z_i : \delta_i = 0\}$

$$(\hat{\theta}, \hat{\mathbf{z}}_0) = \arg \max_{(\theta, \mathbf{z}_0)} \left\{ -\frac{n}{2} \log \sigma^2 - \frac{1}{2} \sum_{i=1}^n \frac{(y_i - \alpha x_i - \beta z_i)^2}{\sigma^2} \right\}.$$

► **Profile likelihood:**

$$\begin{aligned} \hat{z}_i(\theta) = \beta^{-1}(y_i - \alpha x_i) &\implies (y_i - \alpha x_i - \beta \hat{z}_i(\theta))^2 = 0 \\ \implies \hat{\theta} = \arg \max_{\theta} &\left\{ -\frac{n}{2} \log(\sigma^2) - \frac{1}{2} \sum_{i \in \mathcal{S}_1} \frac{(y_i - \alpha x_i - \beta z_i)^2}{\sigma^2} \right\} \end{aligned}$$

Inference with Missing Data

Solution 2: Maximize likelihood over θ and $\mathbf{z}_0 = \{z_i : \delta_i = 0\}$

$$(\hat{\theta}, \hat{\mathbf{z}}_0) = \arg \max_{(\theta, \mathbf{z}_0)} \left\{ -\frac{n}{2} \log \sigma^2 - \frac{1}{2} \sum_{i=1}^n \frac{(y_i - \alpha x_i - \beta z_i)^2}{\sigma^2} \right\}.$$

► **Profile likelihood:**

$$\begin{aligned} \hat{z}_i(\theta) = \beta^{-1}(y_i - \alpha x_i) &\implies (y_i - \alpha x_i - \beta \hat{z}_i(\theta))^2 = 0 \\ \implies \hat{\theta} = \arg \max_{\theta} &\left\{ -\frac{n}{2} \log(\sigma^2) - \frac{1}{2} \sum_{i \in \mathcal{S}_1} \frac{(y_i - \alpha x_i - \beta z_i)^2}{\sigma^2} \right\} \end{aligned}$$

- $(\hat{\alpha}, \hat{\beta})$ exactly the same as using complete data \mathcal{S}_1 only!
- $\hat{\sigma} = \hat{\sigma}_1 \cdot n_1/n$, where $\hat{\sigma}_1$ is the estimator from \mathcal{S}_1 , so confidence intervals even narrower!

Inference with Missing Data

Solution 2: Maximize likelihood over θ and $\mathbf{z}_0 = \{z_i : \delta_i = 0\}$

$$(\hat{\theta}, \hat{\mathbf{z}}_0) = \arg \max_{(\theta, \mathbf{z}_0)} \left\{ -\frac{n}{2} \log \sigma^2 - \frac{1}{2} \sum_{i=1}^n \frac{(y_i - \alpha x_i - \beta z_i)^2}{\sigma^2} \right\}.$$

► **Profile likelihood:**

$$\begin{aligned} \hat{z}_i(\theta) = \beta^{-1}(y_i - \alpha x_i) &\implies (y_i - \alpha x_i - \beta \hat{z}_i(\theta))^2 = 0 \\ \implies \hat{\theta} = \arg \max_{\theta} &\left\{ -\frac{n}{2} \log(\sigma^2) - \frac{1}{2} \sum_{i \in \mathcal{S}_1} \frac{(y_i - \alpha x_i - \beta z_i)^2}{\sigma^2} \right\} \end{aligned}$$

- $(\hat{\alpha}, \hat{\beta})$ exactly the same as using complete data \mathcal{S}_1 only!
- $\hat{\sigma} = \hat{\sigma}_1 \cdot n_1/n$, where $\hat{\sigma}_1$ is the estimator from \mathcal{S}_1 , so confidence intervals even narrower!
- **Problem:** \mathbf{z}_0 is a random variable, not a parameter.

Inference with Missing Data

Solution 3: Model the missing data.

Inference with Missing Data

Solution 3: Model the missing data.

If nothing is known about the missing data mechanism, consider the following model:

$$x \sim p(x | \eta)$$

$\theta = (\alpha, \beta, \sigma)$: original parameters

$$z | x \sim p(z | x, \varphi)$$

φ : nuisance parameters

$$y | z, x \sim \mathcal{N}(\alpha x + \beta z, \sigma^2)$$

η : ignorable parameters

$$\delta | y, z, x \sim \text{Bernoulli}\{r(y, x, \eta)\}$$

$\Theta = (\theta, \varphi, \eta)$: all parameters

Inference with Missing Data

Solution 3: Model the missing data.

$$x \sim p(x | \eta)$$

$\theta = (\alpha, \beta, \sigma)$: original parameters

$$z | x \sim p(z | x, \varphi)$$

φ : nuisance parameters

$$y | z, x \sim \mathcal{N}(\alpha x + \beta z, \sigma^2)$$

η : ignorable parameters

$$\delta | y, z, x \sim \text{Bernoulli}\{r(y, x, \eta)\}$$

$\Theta = (\theta, \varphi, \eta)$: all parameters

► **Likelihood:**

$$\begin{aligned} \mathcal{L}(\Theta | \mathcal{D}) &= \prod_{i \in \mathcal{S}_1} p(\delta_i = 1, y_i, z_i, x_i | \Theta) \times \prod_{i \in \mathcal{S}_0} p(\delta_i = 0, y_i, x_i | \Theta) \\ &= \prod_{i \in \mathcal{S}_1} r(y_i, x_i, \eta) \cdot p(y_i | z_i, x_i, \theta) \cdot p(z_i | x_i, \varphi) \cdot p(x_i | \eta) \\ &\quad \times \prod_{i \in \mathcal{S}_0} [1 - r(y_i, x_i, \eta)] \cdot \underbrace{p(y_i | x_i, \theta)}_{\int p(y_i | x_i, z_i, \theta) \cdot p(z_i | x_i, \varphi) dz_i} \cdot p(x_i | \eta) \end{aligned}$$

Inference with Missing Data

Solution 3: Model the missing data.

$$x \sim p(x | \eta)$$

$$z | x \sim p(z | x, \varphi)$$

$$y | z, x \sim \mathcal{N}(\alpha x + \beta z, \sigma^2)$$

$$\delta | y, z, x \sim \text{Bernoulli}\{r(y, x, \eta)\}$$

$\theta = (\alpha, \beta, \sigma)$: original parameters

φ : nuisance parameters

η : ignorable parameters

$\Theta = (\theta, \varphi, \eta)$: all parameters

► Likelihood:

$$\mathcal{L}(\Theta | \mathcal{D}) = \prod_{i=1}^n r(y_i, x_i, \eta)_i^\delta \cdot [1 - r(y_i, x_i, \eta)]^{1-\delta_i} \cdot p(x_i | \eta)$$

$$\times \prod_{i \in \mathcal{S}_1} p(y_i | z_i, x_i, \theta) \cdot p(z_i | x_i, \varphi) \times \prod_{i \in \mathcal{S}_0} \int p(y_i | x_i, z_i, \theta) \cdot p(z_i | x_i, \varphi) dz_i$$

$$= \mathcal{L}(\eta | \mathcal{D}) \cdot \mathcal{L}(\theta, \varphi | \mathcal{D}) \implies \boxed{\max_{\Theta} \mathcal{L}(\Theta | \mathcal{D}) = \max_{\eta} \mathcal{L}(\eta | \mathcal{D}) \cdot \max_{\theta, \varphi} \mathcal{L}(\theta, \varphi | \mathcal{D})}$$

Inference with Missing Data

Solution 3: Model the missing data.

$$x \sim p(x | \eta)$$

$\theta = (\alpha, \beta, \sigma)$: original parameters

$$z | x \sim p(z | x, \varphi)$$

φ : nuisance parameters

$$y | z, x \sim \mathcal{N}(\alpha x + \beta z, \sigma^2)$$

η : ignorable parameters

$$\delta | y, z, x \sim \text{Bernoulli}\{r(y, x, \eta)\}$$

$\Theta = (\theta, \varphi, \eta)$: all parameters

► Likelihood:

$$\mathcal{L}(\Theta | \mathcal{D}) = \prod_{i=1}^n r(y_i, x_i, \eta)_i^\delta \cdot [1 - r(y_i, x_i, \eta)]^{1-\delta_i} \cdot p(x_i | \eta)$$

$$\times \prod_{i \in \mathcal{S}_1} p(y_i | z_i, x_i, \theta) \cdot p(z_i | x_i, \varphi) \times \prod_{i \in \mathcal{S}_0} \int p(y_i | x_i, z_i, \theta) \cdot p(z_i | x_i, \varphi) dz_i$$

$$= \mathcal{L}(\eta | \mathcal{D}) \cdot \mathcal{L}(\theta, \varphi | \mathcal{D}) \implies \boxed{\max_{\Theta} \mathcal{L}(\Theta | \mathcal{D}) = \max_{\eta} \mathcal{L}(\eta | \mathcal{D}) \cdot \max_{\theta, \varphi} \mathcal{L}(\theta, \varphi | \mathcal{D})}$$

$\implies \hat{\theta}_{\text{ML}}$ does not depend on $p(x | \eta)$ and $p(\delta | y, z, x, \eta)$. The missing data mechanism and covariate distribution of x are thus said to be **ignorable**.

Ignorable vs Nuisance Parameters

▶ **True Data-Generating Process:** $(\mathbf{Y}, \mathbf{X}) \sim p_0(\mathbf{Y}, \mathbf{X})$.

▶ **Conditional Inference Model:**

$$M_C : \mathbf{Y} | \mathbf{X} \sim p(\mathbf{Y} | \mathbf{X}, \theta).$$

▶ Suppose M_C is *correct*, i.e., exists $\theta = \theta_0$ such that $p(\mathbf{Y} | \mathbf{X}, \theta_0) = p_0(\mathbf{Y} | \mathbf{X})$.

▶ Let $\hat{\theta}_C = \arg \max_{\theta} p(\mathbf{Y} | \mathbf{X}, \theta)$. If M_C is correct, then $\hat{\theta}_C \rightarrow \theta_0$ as sample size $N \rightarrow \infty$.

Ignorable vs Nuisance Parameters

- ▶ **True Data-Generating Process:** $(\mathbf{Y}, \mathbf{X}) \sim p_0(\mathbf{Y}, \mathbf{X})$.
- ▶ **Conditional Inference Model:** $M_C : \mathbf{Y} | \mathbf{X} \sim p(\mathbf{Y} | \mathbf{X}, \theta)$.
- ▶ **Full Inference Model:**

$$M_F : (\mathbf{Y}, \mathbf{X}) \sim p(\mathbf{Y} | \mathbf{X}, \theta) \times p(\mathbf{X} | \theta, \eta).$$

- ▶ Let $(\hat{\theta}_F, \hat{\eta}_F) = \arg \max_{(\theta, \eta)} p(\mathbf{Y} | \mathbf{X} | \theta) \cdot p(\mathbf{X} | \theta, \eta)$.
- ▶ If marginal model $M_X : \mathbf{X} \sim p(\mathbf{X} | \eta)$ **does not** depend on θ , then $\hat{\theta}_F = \hat{\theta}_C$, i.e., $p(\mathbf{X} | \eta)$ is **ignorable**.
- ▶ If $M_X : \mathbf{X} \sim p(\mathbf{X} | \theta, \eta)$ **does** depend on θ , then $\hat{\theta}_F \neq \hat{\theta}_C$.
 - ▶ If M_X is correct, then $\hat{\theta}_F \rightarrow \theta_0$, and $\text{var}(\hat{\theta}_F) < \text{var}(\hat{\theta}_C)$. Since we need to maximize over η to get $\hat{\theta}_F$, η are called **nuisance parameters**.
 - ▶ If M_X is incorrect, then generally $\hat{\theta}_F \not\rightarrow \theta_0$, even if M_C is correct.

Inference with Missing Data (Continued)

- ▶ **Missing Data Setup:** $y_i = \alpha x_i + \beta z_i + \sigma \varepsilon_i$, $\varepsilon_i \stackrel{\text{iid}}{\sim} \mathcal{N}(0, 1)$.
 $\delta_i = 1$ if z_i is observed and $\delta_i = 0$ if it is missing.
- ▶ **Complete Data Model:** $M : (y, x, z) \sim \mathcal{N}(\mathbf{0}, \Sigma)$
 $\delta \mid y, x, z \sim \text{Bernoulli}\{r(y, x, \eta)\}$

Note that under M , we have $y \mid x, z \sim \mathcal{N}(\alpha x + \beta z, \sigma^2)$, where

$$\begin{bmatrix} \alpha \\ \beta \end{bmatrix} = \begin{bmatrix} \Sigma_{xx} & \Sigma_{xz} \\ \Sigma_{zx} & \Sigma_{zz} \end{bmatrix}^{-1} \begin{bmatrix} \Sigma_{xy} \\ \Sigma_{zy} \end{bmatrix}, \quad \sigma^2 = \Sigma_{yy} - \begin{bmatrix} \Sigma_{yx} & \Sigma_{yz} \end{bmatrix} \begin{bmatrix} \alpha \\ \beta \end{bmatrix}.$$

Inference with Missing Data (Continued)

- ▶ **Missing Data Setup:** $y_i = \alpha x_i + \beta z_i + \sigma \varepsilon_i$, $\varepsilon_i \stackrel{\text{iid}}{\sim} \mathcal{N}(0, 1)$.
 $\delta_i = 1$ if z_i is observed and $\delta_i = 0$ if it is missing.
- ▶ **Complete Data Model:** $M : (y, x, z) \sim \mathcal{N}(\mathbf{0}, \Sigma)$
(and **ignorable** missing data), with $M : y \mid x, z \sim \mathcal{N}(\alpha x + \beta z, \sigma^2)$.

- ▶ **Observed Data Likelihood:**

$$\begin{aligned} \ell(\Sigma \mid \mathcal{D}) = & -\frac{1}{2} \sum_{i \in \mathcal{S}_1} \left\{ \begin{bmatrix} y_i & x_i & z_i \end{bmatrix} \Sigma^{-1} \begin{bmatrix} y_i \\ x_i \\ z_i \end{bmatrix} + \log |\Sigma| \right\} \\ & -\frac{1}{2} \sum_{i \in \mathcal{S}_0} \left\{ \begin{bmatrix} y_i & x_i \end{bmatrix} \begin{bmatrix} \Sigma_{yy} & \Sigma_{yx} \\ \Sigma_{xy} & \Sigma_{xx} \end{bmatrix}^{-1} \begin{bmatrix} y_i \\ x_i \end{bmatrix} + \log \left| \begin{array}{cc} \Sigma_{yy} & \Sigma_{yx} \\ \Sigma_{xy} & \Sigma_{xx} \end{array} \right| \right\} \end{aligned}$$

Inference with Missing Data (Continued)

- ▶ **Missing Data Setup:** $y_i = \alpha x_i + \beta z_i + \sigma \varepsilon_i$, $\varepsilon_i \stackrel{\text{iid}}{\sim} \mathcal{N}(0, 1)$.
 $\delta_i = 1$ if z_i is observed and $\delta_i = 0$ if it is missing.

- ▶ **Complete Data Model:** $M : (y, x, z) \sim \mathcal{N}(\mathbf{0}, \Sigma)$

(and ignorable missing data), with $M : y | x, z \sim \mathcal{N}(\alpha x + \beta z, \sigma^2)$.

- ▶ **Observed Data Likelihood:**

$$\begin{aligned} \ell(\Sigma | \mathcal{D}) = & -\frac{1}{2} \sum_{i \in \mathcal{S}_1} \left\{ [y_i \ x_i \ z_i] \Sigma^{-1} \begin{bmatrix} y_i \\ x_i \\ z_i \end{bmatrix} + \log |\Sigma| \right\} \\ & -\frac{1}{2} \sum_{i \in \mathcal{S}_0} \left\{ [y_i \ x_i] \begin{bmatrix} \Sigma_{yy} & \Sigma_{yx} \\ \Sigma_{xy} & \Sigma_{xx} \end{bmatrix}^{-1} \begin{bmatrix} y_i \\ x_i \end{bmatrix} + \log \left| \begin{bmatrix} \Sigma_{yy} & \Sigma_{yx} \\ \Sigma_{xy} & \Sigma_{xx} \end{bmatrix} \right| \right\} \end{aligned}$$

- ▶ **Inference:** $\hat{\Sigma} = \arg \max_{\Sigma} \ell(\Sigma | \mathcal{D})$

Difficult to calculate directly, but simple when $\mathbf{z}_0 = \{z_i : \delta_i = 0\}$ is observed!

That is, if $\mathbf{Y}_{n \times 3} = (\mathbf{y}, \mathbf{x}, \mathbf{z})$, then $\hat{\Sigma} = \frac{1}{n} \mathbf{Y}' \mathbf{Y}$.

Inference with Missing Data (Continued)

► **Observed Data Likelihood:**

$$\begin{aligned} \ell(\boldsymbol{\Sigma} | \mathcal{D}) = & -\frac{1}{2} \sum_{i \in \mathcal{S}_1} \left\{ [y_i \ x_i \ z_i] \boldsymbol{\Sigma}^{-1} \begin{bmatrix} y_i \\ x_i \\ z_i \end{bmatrix} + \log |\boldsymbol{\Sigma}| \right\} \\ & -\frac{1}{2} \sum_{i \in \mathcal{S}_0} \left\{ [y_i \ x_i] \begin{bmatrix} \boldsymbol{\Sigma}_{yy} & \boldsymbol{\Sigma}_{yx} \\ \boldsymbol{\Sigma}_{xy} & \boldsymbol{\Sigma}_{xx} \end{bmatrix}^{-1} \begin{bmatrix} y_i \\ x_i \end{bmatrix} + \log \left| \begin{array}{cc} \boldsymbol{\Sigma}_{yy} & \boldsymbol{\Sigma}_{yx} \\ \boldsymbol{\Sigma}_{xy} & \boldsymbol{\Sigma}_{xx} \end{array} \right| \right\} \end{aligned}$$

► **Inference:** $\hat{\boldsymbol{\Sigma}} = \arg \max_{\boldsymbol{\Sigma}} \ell(\boldsymbol{\Sigma} | \mathcal{D})$

Difficult to calculate directly, but simple when $\mathbf{z}_0 = \{z_i : \delta_i = 0\}$ is observed. That is, if $\mathbf{Y}_{n \times 3} = (\mathbf{y}, \mathbf{x}, \mathbf{z})$, then $\hat{\boldsymbol{\Sigma}} = \frac{1}{n} \mathbf{Y}' \mathbf{Y}$.

► **Strategy:** Iterative algorithm $(\hat{\boldsymbol{\Sigma}}^{(1)}, \hat{\mathbf{z}}_0^{(1)}), \dots, (\hat{\boldsymbol{\Sigma}}^{(m)}, \hat{\mathbf{z}}_0^{(m)})$

► $\hat{\boldsymbol{\Sigma}}^{(m+1)} = \hat{\boldsymbol{\Sigma}}(\mathbf{y}, \mathbf{x}, \hat{\mathbf{z}}_0^{(m)}, \mathbf{z}_1)$

► $\hat{\mathbf{z}}_0^{(m+1)} = ???$

Inference with Missing Data (Continued)

► Observed Data Likelihood:

$$\begin{aligned} \ell(\Sigma | \mathcal{D}) = & -\frac{1}{2} \sum_{i \in \mathcal{S}_1} \left\{ [y_i \ x_i \ z_i] \Sigma^{-1} \begin{bmatrix} y_i \\ x_i \\ z_i \end{bmatrix} + \log |\Sigma| \right\} \\ & -\frac{1}{2} \sum_{i \in \mathcal{S}_0} \left\{ [y_i \ x_i] \begin{bmatrix} \Sigma_{yy} & \Sigma_{yx} \\ \Sigma_{xy} & \Sigma_{xx} \end{bmatrix}^{-1} \begin{bmatrix} y_i \\ x_i \end{bmatrix} + \log \left| \begin{array}{cc} \Sigma_{yy} & \Sigma_{yx} \\ \Sigma_{xy} & \Sigma_{xx} \end{array} \right| \right\} \end{aligned}$$

► Inference: $\hat{\Sigma} = \arg \max_{\Sigma} \ell(\Sigma | \mathcal{D})$

Difficult to calculate directly, but simple when $\mathbf{z}_0 = \{z_i : \delta_i = 0\}$ is observed. That is, if $\mathbf{Y}_{n \times 3} = (\mathbf{y}, \mathbf{x}, \mathbf{z})$, then $\hat{\Sigma} = \frac{1}{n} \mathbf{Y}' \mathbf{Y}$.

► Strategy: Iterative algorithm $(\hat{\Sigma}^{(1)}, \hat{\mathbf{z}}_0^{(1)}), \dots, (\hat{\Sigma}^{(m)}, \hat{\mathbf{z}}_0^{(m)})$

► $\hat{\Sigma}^{(m+1)} = \hat{\Sigma}(\mathbf{y}, \mathbf{x}, \hat{\mathbf{z}}_0^{(m)}, \mathbf{z}_1)$

► $\hat{\mathbf{z}}_0^{(m+1)} = \arg \max_{\mathbf{z}_0} p(\mathbf{z}_0 | \mathcal{D}, \hat{\Sigma}^{(m+1)})?$

Inference with Missing Data (Continued)

► **Observed Data Likelihood:**

$$\begin{aligned} \ell(\Sigma | \mathcal{D}) = & -\frac{1}{2} \sum_{i \in \mathcal{S}_1} \left\{ [y_i \ x_i \ z_i] \Sigma^{-1} \begin{bmatrix} y_i \\ x_i \\ z_i \end{bmatrix} + \log |\Sigma| \right\} \\ & -\frac{1}{2} \sum_{i \in \mathcal{S}_0} \left\{ [y_i \ x_i] \begin{bmatrix} \Sigma_{yy} & \Sigma_{yx} \\ \Sigma_{xy} & \Sigma_{xx} \end{bmatrix}^{-1} \begin{bmatrix} y_i \\ x_i \end{bmatrix} + \log \left| \begin{array}{cc} \Sigma_{yy} & \Sigma_{yx} \\ \Sigma_{xy} & \Sigma_{xx} \end{array} \right| \right\} \end{aligned}$$

► **Inference:** $\hat{\Sigma} = \arg \max_{\Sigma} \ell(\Sigma | \mathcal{D})$

Difficult to calculate directly, but simple when $\mathbf{z}_0 = \{z_i : \delta_i = 0\}$ is observed. That is, if

$\mathbf{Y}_{n \times 3} = (\mathbf{y}, \mathbf{x}, \mathbf{z})$, then $\hat{\Sigma} = \frac{1}{n} \mathbf{Y}' \mathbf{Y}$.

► **Strategy:** Iterative algorithm $(\hat{\Sigma}^{(1)}, \hat{\mathbf{z}}_0^{(1)}), \dots, (\hat{\Sigma}^{(m)}, \hat{\mathbf{z}}_0^{(m)})$

► $\hat{\Sigma}^{(m+1)} = \hat{\Sigma}(\mathbf{y}, \mathbf{x}, \hat{\mathbf{z}}_0^{(m)}, \mathbf{z}_1)$

► $\hat{\mathbf{z}}_0^{(m+1)} = \arg \max_{\mathbf{z}_0} p(\mathbf{z}_0 | \mathcal{D}, \hat{\Sigma}^{(m+1)})?$

No, converges to $\arg \max_{\Sigma, \mathbf{z}_0} \ell(\Sigma | \mathcal{D}, \mathbf{z}_0) \neq \arg \max_{\Sigma} \ell(\Sigma | \mathcal{D})$.

Inference with Missing Data (Continued)

► **Observed Data Likelihood:**

$$\begin{aligned} \ell(\boldsymbol{\Sigma} | \mathcal{D}) = & -\frac{1}{2} \sum_{i \in \mathcal{S}_1} \left\{ [y_i \ x_i \ z_i] \boldsymbol{\Sigma}^{-1} \begin{bmatrix} y_i \\ x_i \\ z_i \end{bmatrix} + \log |\boldsymbol{\Sigma}| \right\} \\ & -\frac{1}{2} \sum_{i \in \mathcal{S}_0} \left\{ [y_i \ x_i] \begin{bmatrix} \boldsymbol{\Sigma}_{yy} & \boldsymbol{\Sigma}_{yx} \\ \boldsymbol{\Sigma}_{xy} & \boldsymbol{\Sigma}_{xx} \end{bmatrix}^{-1} \begin{bmatrix} y_i \\ x_i \end{bmatrix} + \log \left| \begin{array}{cc} \boldsymbol{\Sigma}_{yy} & \boldsymbol{\Sigma}_{yx} \\ \boldsymbol{\Sigma}_{xy} & \boldsymbol{\Sigma}_{xx} \end{array} \right| \right\} \end{aligned}$$

► **Inference:** $\hat{\boldsymbol{\Sigma}} = \arg \max_{\boldsymbol{\Sigma}} \ell(\boldsymbol{\Sigma} | \mathcal{D})$

Difficult to calculate directly, but simple when $\mathbf{z}_0 = \{z_i : \delta_i = 0\}$ is observed. That is, if $\mathbf{Y}_{n \times 3} = (\mathbf{y}, \mathbf{x}, \mathbf{z})$, then $\hat{\boldsymbol{\Sigma}} = \frac{1}{n} \mathbf{Y}' \mathbf{Y}$.

► **Strategy:** Iterative algorithm $(\hat{\boldsymbol{\Sigma}}^{(1)}, \hat{\mathbf{z}}_0^{(1)}), \dots, (\hat{\boldsymbol{\Sigma}}^{(m)}, \hat{\mathbf{z}}_0^{(m)})$

► $\hat{\boldsymbol{\Sigma}}^{(m+1)} = \hat{\boldsymbol{\Sigma}}(\mathbf{y}, \mathbf{x}, \hat{\mathbf{z}}_0^{(m)}, \mathbf{z}_1)$

► $\hat{\mathbf{z}}_0^{(m+1)} \sim p(\mathbf{z}_0 | \mathcal{D}, \hat{\boldsymbol{\Sigma}}^{(m+1)})?$

Inference with Missing Data (Continued)

► Observed Data Likelihood:

$$\begin{aligned} \ell(\boldsymbol{\Sigma} | \mathcal{D}) = & -\frac{1}{2} \sum_{i \in \mathcal{S}_1} \left\{ \begin{bmatrix} y_i & x_i & z_i \end{bmatrix} \boldsymbol{\Sigma}^{-1} \begin{bmatrix} y_i \\ x_i \\ z_i \end{bmatrix} + \log |\boldsymbol{\Sigma}| \right\} \\ & -\frac{1}{2} \sum_{i \in \mathcal{S}_0} \left\{ \begin{bmatrix} y_i & x_i \end{bmatrix} \begin{bmatrix} \boldsymbol{\Sigma}_{yy} & \boldsymbol{\Sigma}_{yx} \\ \boldsymbol{\Sigma}_{xy} & \boldsymbol{\Sigma}_{xx} \end{bmatrix}^{-1} \begin{bmatrix} y_i \\ x_i \end{bmatrix} + \log \left| \begin{bmatrix} \boldsymbol{\Sigma}_{yy} & \boldsymbol{\Sigma}_{yx} \\ \boldsymbol{\Sigma}_{xy} & \boldsymbol{\Sigma}_{xx} \end{bmatrix} \right| \right\} \end{aligned}$$

► Inference: $\hat{\boldsymbol{\Sigma}} = \arg \max_{\boldsymbol{\Sigma}} \ell(\boldsymbol{\Sigma} | \mathcal{D})$

Difficult to calculate directly, but simple when $\mathbf{z}_0 = \{z_i : \delta_i = 0\}$ is observed. That is, if $\mathbf{Y}_{n \times 3} = (\mathbf{y}, \mathbf{x}, \mathbf{z})$, then $\hat{\boldsymbol{\Sigma}} = \frac{1}{n} \mathbf{Y}' \mathbf{Y}$.

► Strategy: Iterative algorithm $(\hat{\boldsymbol{\Sigma}}^{(1)}, \hat{\mathbf{z}}_0^{(1)}), \dots, (\hat{\boldsymbol{\Sigma}}^{(m)}, \hat{\mathbf{z}}_0^{(m)})$

► $\hat{\boldsymbol{\Sigma}}^{(m+1)} = \hat{\boldsymbol{\Sigma}}(\mathbf{y}, \mathbf{x}, \hat{\mathbf{z}}_0^{(m)}, \mathbf{z}_1)$

► $\hat{\mathbf{z}}_0^{(m+1)} \sim p(\mathbf{z}_0 | \mathcal{D}, \hat{\boldsymbol{\Sigma}}^{(m+1)})?$

This produces a stationary stochastic process $\hat{\boldsymbol{\Sigma}}^{(1)}, \hat{\boldsymbol{\Sigma}}^{(2)}, \dots$, for which the expectation $\tilde{\boldsymbol{\Sigma}} = E[\hat{\boldsymbol{\Sigma}}^{(t)}] \rightarrow \boldsymbol{\Sigma}_0$ as $n \rightarrow \infty$. However, $\tilde{\boldsymbol{\Sigma}}$ is less efficient than the MLE...

Inference with Missing Data (Continued)

► **Observed Data Likelihood:**

$$\begin{aligned} \ell(\boldsymbol{\Sigma} | \mathcal{D}) = & -\frac{1}{2} \sum_{i \in \mathcal{S}_1} \left\{ [y_i \ x_i \ z_i] \boldsymbol{\Sigma}^{-1} \begin{bmatrix} y_i \\ x_i \\ z_i \end{bmatrix} + \log |\boldsymbol{\Sigma}| \right\} \\ & -\frac{1}{2} \sum_{i \in \mathcal{S}_0} \left\{ [y_i \ x_i] \begin{bmatrix} \boldsymbol{\Sigma}_{yy} & \boldsymbol{\Sigma}_{yx} \\ \boldsymbol{\Sigma}_{xy} & \boldsymbol{\Sigma}_{xx} \end{bmatrix}^{-1} \begin{bmatrix} y_i \\ x_i \end{bmatrix} + \log \left| \begin{array}{cc} \boldsymbol{\Sigma}_{yy} & \boldsymbol{\Sigma}_{yx} \\ \boldsymbol{\Sigma}_{xy} & \boldsymbol{\Sigma}_{xx} \end{array} \right| \right\} \end{aligned}$$

► **Inference:** $\hat{\boldsymbol{\Sigma}} = \arg \max_{\boldsymbol{\Sigma}} \ell(\boldsymbol{\Sigma} | \mathcal{D})$

Difficult to calculate directly, but simple when $\mathbf{z}_0 = \{z_i : \delta_i = 0\}$ is observed. That is, if $\mathbf{Y}_{n \times 3} = (\mathbf{y}, \mathbf{x}, \mathbf{z})$, then $\hat{\boldsymbol{\Sigma}} = \frac{1}{n} \mathbf{Y}' \mathbf{Y}$.

► **Strategy:** Iterative algorithm $(\hat{\boldsymbol{\Sigma}}^{(1)}, \hat{\mathbf{z}}_0^{(1)}), \dots, (\hat{\boldsymbol{\Sigma}}^{(m)}, \hat{\mathbf{z}}_0^{(m)})$

► $\hat{\boldsymbol{\Sigma}}^{(m+1)} = \hat{\boldsymbol{\Sigma}}(\mathbf{y}, \mathbf{x}, \hat{\mathbf{z}}_0^{(m)}, \mathbf{z}_1)$

► $\hat{\mathbf{z}}_0^{(m+1)} = E[\mathbf{z}_0 | \mathcal{D}, \hat{\boldsymbol{\Sigma}}^{(m+1)}]?$

Inference with Missing Data (Continued)

► **Observed Data Likelihood:**

$$\begin{aligned} \ell(\boldsymbol{\Sigma} | \mathcal{D}) = & -\frac{1}{2} \sum_{i \in \mathcal{S}_1} \left\{ [y_i \ x_i \ z_i] \boldsymbol{\Sigma}^{-1} \begin{bmatrix} y_i \\ x_i \\ z_i \end{bmatrix} + \log |\boldsymbol{\Sigma}| \right\} \\ & -\frac{1}{2} \sum_{i \in \mathcal{S}_0} \left\{ [y_i \ x_i] \begin{bmatrix} \boldsymbol{\Sigma}_{yy} & \boldsymbol{\Sigma}_{yx} \\ \boldsymbol{\Sigma}_{xy} & \boldsymbol{\Sigma}_{xx} \end{bmatrix}^{-1} \begin{bmatrix} y_i \\ x_i \end{bmatrix} + \log \left| \begin{bmatrix} \boldsymbol{\Sigma}_{yy} & \boldsymbol{\Sigma}_{yx} \\ \boldsymbol{\Sigma}_{xy} & \boldsymbol{\Sigma}_{xx} \end{bmatrix} \right| \right\} \end{aligned}$$

► **Inference:** $\hat{\boldsymbol{\Sigma}} = \arg \max_{\boldsymbol{\Sigma}} \ell(\boldsymbol{\Sigma} | \mathcal{D})$

Difficult to calculate directly, but simple when $\mathbf{z}_0 = \{z_i : \delta_i = 0\}$ is observed. That is, if

$\mathbf{Y}_{n \times 3} = (\mathbf{y}, \mathbf{x}, \mathbf{z})$, then $\hat{\boldsymbol{\Sigma}} = \frac{1}{n} \mathbf{Y}' \mathbf{Y}$.

► **Strategy:** Iterative algorithm $(\hat{\boldsymbol{\Sigma}}^{(1)}, \hat{\mathbf{z}}_0^{(1)}), \dots, (\hat{\boldsymbol{\Sigma}}^{(m)}, \hat{\mathbf{z}}_0^{(m)})$

► $\hat{\boldsymbol{\Sigma}}^{(m+1)} = \hat{\boldsymbol{\Sigma}}(\mathbf{y}, \mathbf{x}, \hat{\mathbf{z}}_0^{(m)}, \mathbf{z}_1)$

► $\hat{\mathbf{z}}_0^{(m+1)} = E[\mathbf{z}_0 | \mathcal{D}, \hat{\boldsymbol{\Sigma}}^{(m+1)}]?$

Almost!

The Expectation-Maximization Algorithm (EM)

▶ **Setup:**

▶ \mathbf{y}_{obs} : observed data

▶ \mathbf{y}_{miss} : missing data

▶ $\mathbf{y}_{\text{comp}} = \mathbf{y}_{\text{obs}} \cup \mathbf{y}_{\text{miss}}$: complete data

▶ **Goal:** Find $\hat{\theta} = \arg \max_{\theta} \ell(\theta | \mathbf{y}_{\text{obs}})$.

▶ **Problem:** $\mathcal{L}(\theta | \mathbf{y}_{\text{comp}})$ is tractable but

$$\mathcal{L}(\theta | \mathbf{y}_{\text{obs}}) = \int p(\mathbf{y}_{\text{obs}}, \mathbf{y}_{\text{miss}} | \theta) d\mathbf{y}_{\text{miss}}$$

is **not**.

The Expectation-Maximization Algorithm (EM)

► Setup:

► \mathbf{y}_{obs} : observed data

► \mathbf{y}_{miss} : missing data

► $\mathbf{y}_{\text{comp}} = \mathbf{y}_{\text{obs}} \cup \mathbf{y}_{\text{miss}}$: complete data

► **Goal:** Find $\hat{\theta} = \arg \max_{\theta} \ell(\theta | \mathbf{y}_{\text{obs}})$.

► **EM Algorithm:** An *iterative* algorithm $\hat{\theta}^{(1)}, \hat{\theta}^{(2)}, \dots$ alternating between two steps:

► **E-Step:** Construct function $Q_t(\theta) = E[\ell(\theta | \mathbf{y}_{\text{comp}}) | \mathbf{y}_{\text{obs}}, \hat{\theta}^{(t)}]$
$$= \int \ell(\theta | \mathbf{y}_{\text{obs}}, \mathbf{y}_{\text{miss}}) \cdot p(\mathbf{y}_{\text{miss}} | \mathbf{y}_{\text{obs}}, \hat{\theta}^{(t)}) d\mathbf{y}_{\text{miss}}$$

The Expectation-Maximization Algorithm (EM)

► Setup:

► \mathbf{y}_{obs} : observed data

► \mathbf{y}_{miss} : missing data

► $\mathbf{y}_{\text{comp}} = \mathbf{y}_{\text{obs}} \cup \mathbf{y}_{\text{miss}}$: complete data

► **Goal:** Find $\hat{\theta} = \arg \max_{\theta} \ell(\theta | \mathbf{y}_{\text{obs}})$.

► **EM Algorithm:** An *iterative* algorithm $\hat{\theta}^{(1)}, \hat{\theta}^{(2)}, \dots$ alternating between two steps:

► **E-Step:** Construct function $Q_t(\theta) = E[\ell(\theta | \mathbf{y}_{\text{comp}}) | \mathbf{y}_{\text{obs}}, \hat{\theta}^{(t)}]$
$$= \int \ell(\theta | \mathbf{y}_{\text{obs}}, \mathbf{y}_{\text{miss}}) \cdot p(\mathbf{y}_{\text{miss}} | \mathbf{y}_{\text{obs}}, \hat{\theta}^{(t)}) d\mathbf{y}_{\text{miss}}$$

► **M-Step:** Maximize to find next value of θ : $\hat{\theta}^{(t+1)} = \arg \max_{\theta} Q_t(\theta)$.

The EM Algorithm: Exponential Families

► **Model:**

$$p(\mathbf{y}_{\text{comp}} | \boldsymbol{\eta}) = \exp \{ \mathbf{T}' \boldsymbol{\eta} - \Psi(\boldsymbol{\eta}) \} h(\mathbf{y}_{\text{comp}})$$

► **E-step:**

$$\begin{aligned} Q_t(\boldsymbol{\eta}) &= E[\ell(\boldsymbol{\eta} | \mathbf{y}_{\text{comp}}) | \mathbf{y}_{\text{obs}}, \hat{\boldsymbol{\eta}}^{(t)}] \\ &= \bar{\mathbf{T}}_t' \boldsymbol{\eta} - \Psi(\boldsymbol{\eta}), \end{aligned}$$

$$\bar{\mathbf{T}}_t = E[\mathbf{T} | \mathbf{y}_{\text{obs}}, \hat{\boldsymbol{\eta}}^{(t)}],$$

and $\bar{\mathbf{T}}_t$ often easy to compute.

► **M-step:** Convex optimization!

The EM Algorithm: Monotonicity

Theorem. If $\hat{\theta}^{(t)}$ and $\hat{\theta}^{(t+1)}$ are successive steps of the EM algorithm, then

$$\ell(\hat{\theta}^{(t)} | \mathbf{y}_{\text{obs}}) \leq \ell(\hat{\theta}^{(t+1)} | \mathbf{y}_{\text{obs}}).$$

The EM Algorithm: Rate of Convergence

- ▶ Convergence of EM to (local) mode θ^* is *linear*:

$$|\hat{\theta}^{(t+1)} - \theta^*| < K \times |\hat{\theta}^{(t)} - \theta^*|.$$

- ▶ Convergence of Newton-Raphson to (local) model is *quadratic*:

$$|\hat{\theta}^{(t+1)} - \theta^*| < K \times |\hat{\theta}^{(t)} - \theta^*|^2.$$

- ▶ In practice:

- ▶ Whichever is easier to implement will work better.
- ▶ EM can be used to find a better starting point for NR.

Example: Multivariate Normal

► **Model:** $\mathbf{y} = (\mathbf{x}, \mathbf{z}) \sim \mathcal{N} \left\{ \mathbf{0}, \begin{bmatrix} \boldsymbol{\Sigma}_{xx} & \boldsymbol{\Sigma}_{xz} \\ \boldsymbol{\Sigma}_{zx} & \boldsymbol{\Sigma}_{zz} \end{bmatrix} \right\},$ $\mathbf{x} = (x_1, \dots, x_p),$
 $\mathbf{z} = (z_1, \dots, z_q).$

► **Missing Data:** x_i always observed, but $\delta_i = \begin{cases} 1 & z_i \text{ observed} \\ 0 & z_i \text{ missing} \end{cases}$

► **Observed Data:**

► Let $\mathcal{S}_k = \{i : \delta_i = k\},$ $\mathbf{Z}_k = \{z_i : i \in \mathcal{S}_k\},$ $k = 0, 1.$

► $\mathbf{y}_{\text{obs}} = \mathcal{D} = (\mathbf{X}, \mathbf{Z}_1, \boldsymbol{\delta}).$

► **Complete Data:**

$$\mathbf{X}_{n \times p} = (\mathbf{x}_1, \dots, \mathbf{x}_n),$$

► $\mathbf{y}_{\text{comp}} = (\mathbf{Y}, \boldsymbol{\delta}),$ $\mathbf{Y}_{n \times (p+q)} = (\mathbf{X}, \mathbf{Z}),$ $\mathbf{Z}_{n \times q} = (\mathbf{z}_1, \dots, \mathbf{z}_n).$

► Previous example $y \sim \mathcal{N}(\alpha x + \beta z, \sigma^2)$ is a special case with $p = 2$ and $q = 1:$

$$\mathbf{x} \leftarrow (\mathbf{y}, \mathbf{x}), \quad \mathbf{z} \leftarrow \mathbf{z}.$$

Example: Multivariate Normal

- ▶ **Model:** $\mathbf{y} = (\mathbf{x}, \mathbf{z}) \sim \mathcal{N} \left\{ \mathbf{0}, \begin{bmatrix} \boldsymbol{\Sigma}_{xx} & \boldsymbol{\Sigma}_{xz} \\ \boldsymbol{\Sigma}_{zx} & \boldsymbol{\Sigma}_{zz} \end{bmatrix} \right\}, \quad \begin{aligned} \mathbf{x} &= (x_1, \dots, x_p), \\ \mathbf{z} &= (z_1, \dots, z_q). \end{aligned}$
- ▶ **Missing Data:** x_i always observed, but $\delta_i = 1$ (0) if z_i is observed (missing)
- ▶ **Observed Data:** $\mathbf{y}_{\text{obs}} = \mathcal{D} = (\mathbf{X}, \mathbf{Z}_1, \boldsymbol{\delta}), \quad \mathbf{Z}_1 = \{z_i : \delta_i = 1\}.$
- ▶ **Complete Data:** $\mathbf{y}_{\text{comp}} = (\mathbf{Y}, \boldsymbol{\delta}), \quad \mathbf{Y} = (\mathbf{X}, \mathbf{Z}).$
- ▶ **Complete Data Likelihood:** Assuming an **ignorable** missing data mechanism $\boldsymbol{\delta} \mid \mathbf{x}, \mathbf{z} \sim \text{Bernoulli}\{r(\mathbf{x}, \boldsymbol{\eta})\},$

$$\begin{aligned} \ell(\boldsymbol{\Sigma} \mid \mathbf{y}_{\text{comp}}) &= -\frac{1}{2} \left\{ n \log |\boldsymbol{\Sigma}| + \sum_{i=1}^n \mathbf{y}_i' \boldsymbol{\Sigma}^{-1} \mathbf{y}_i \right\} \\ &= -\frac{1}{2} \left\{ n \log |\boldsymbol{\Sigma}| + \sum_{i=1}^n \text{tr}(\boldsymbol{\Sigma}^{-1} \mathbf{y}_i \mathbf{y}_i') \right\}. \end{aligned}$$

Multivariate Normal: EM Algorithm

► **Model:** $\mathbf{y} = (\mathbf{x}, \mathbf{z}) \sim \mathcal{N} \left\{ \mathbf{0}, \begin{bmatrix} \boldsymbol{\Sigma}_{xx} & \boldsymbol{\Sigma}_{xz} \\ \boldsymbol{\Sigma}_{zx} & \boldsymbol{\Sigma}_{zz} \end{bmatrix} \right\}$, $\mathbf{x} = (x_1, \dots, x_p)$,
 $\mathbf{z} = (z_1, \dots, z_q)$.

► **Observed Data:** $\mathbf{y}_{\text{obs}} = \mathcal{D} = (\mathbf{X}, \mathbf{Z}_1, \delta)$, $\mathbf{Z}_1 = \{\mathbf{z}_i : \delta_i = 1\}$.

► **Complete Data Likelihood:** For $\mathbf{y}_{\text{comp}} = (\mathbf{Y}, \delta)$, $\mathbf{Y} = (\mathbf{X}, \mathbf{Z})$,

$$\ell(\boldsymbol{\Sigma} | \mathbf{y}_{\text{comp}}) = -\frac{1}{2} \left\{ n \log |\boldsymbol{\Sigma}| + \sum_{i=1}^n \text{tr}(\boldsymbol{\Sigma}^{-1} \mathbf{y}_i \mathbf{y}_i') \right\}.$$

► **E-Step:**

► **Q-Function:**

$$Q_t(\boldsymbol{\Sigma}) = -\frac{1}{2} \left\{ n \log |\boldsymbol{\Sigma}| + \text{tr}(\boldsymbol{\Sigma}^{-1} \mathbf{Y}'_1 \mathbf{Y}_1) + \sum_{i \in \mathcal{S}_0} \text{tr} \left(\boldsymbol{\Sigma}^{-1} E \left[\mathbf{y}_i \mathbf{y}_i' | \mathbf{x}_i, \hat{\boldsymbol{\Sigma}}^{(t)} \right] \right) \right\}, \text{ where}$$

$$\mathbf{Y}_1 = \{\mathbf{y}_i : i \in \mathcal{S}_1\}.$$

Multivariate Normal: EM Algorithm

- ▶ **Complete Data Likelihood:** For $\mathbf{y}_{\text{comp}} = (\mathbf{Y}, \delta)$, $\mathbf{Y} = (\mathbf{X}, \mathbf{Z})$,

$$\ell(\boldsymbol{\Sigma} | \mathbf{y}_{\text{comp}}) = -\frac{1}{2} \left\{ n \log |\boldsymbol{\Sigma}| + \sum_{i=1}^n \text{tr}(\boldsymbol{\Sigma}^{-1} \mathbf{y}_i \mathbf{y}_i') \right\}.$$

- ▶ **E-Step:**

- ▶ **Q-Function:**

$$Q_t(\boldsymbol{\Sigma}) = -\frac{1}{2} \left\{ n \log |\boldsymbol{\Sigma}| + \text{tr}(\boldsymbol{\Sigma}^{-1} \mathbf{Y}'_1 \mathbf{Y}_1) + \sum_{i \in \mathcal{S}_0} \text{tr}(\boldsymbol{\Sigma}^{-1} E[\mathbf{y}_i \mathbf{y}_i' | \mathbf{x}_i, \hat{\boldsymbol{\Sigma}}^{(t)}]) \right\}, \text{ where}$$

$$\mathbf{Y}_1 = \{\mathbf{y}_i : i \in \mathcal{S}_1\}.$$

- ▶ **Conditional Expectation:**

$$\mathbf{z}_i | \mathbf{x}_i, \hat{\boldsymbol{\Sigma}}^{(t)} \sim \mathcal{N}(\hat{\boldsymbol{\mu}}_i^{(t)}, \hat{\boldsymbol{\Omega}}^{(t)}), \quad \begin{aligned} \hat{\boldsymbol{\mu}}_i^{(t)} &= \hat{\boldsymbol{\Sigma}}_{\text{zx}}^{(t)} [\hat{\boldsymbol{\Sigma}}_{\text{xx}}^{(t)}]^{-1} \mathbf{x}_i \\ \hat{\boldsymbol{\Omega}}^{(t)} &= \hat{\boldsymbol{\Sigma}}_{\text{zz}}^{(t)} - \hat{\boldsymbol{\Sigma}}_{\text{zx}}^{(t)} [\hat{\boldsymbol{\Sigma}}_{\text{xx}}^{(t)}]^{-1} \hat{\boldsymbol{\Sigma}}_{\text{zx}}^{(t)'} \end{aligned}$$

$$\Rightarrow E[\mathbf{y}_i \mathbf{y}_i' | \mathbf{x}_i, \hat{\boldsymbol{\Sigma}}^{(t)}] = E \left\{ \begin{bmatrix} \mathbf{x}_i \mathbf{x}_i' & \mathbf{x}_i \mathbf{z}_i' \\ \mathbf{z}_i \mathbf{x}_i' & \mathbf{z}_i \mathbf{z}_i' \end{bmatrix} \middle| \mathbf{x}_i, \hat{\boldsymbol{\Sigma}}^{(t)} \right\} = \underbrace{\begin{bmatrix} \mathbf{x}_i \mathbf{x}_i' & \mathbf{x}_i [\hat{\boldsymbol{\mu}}_i^{(t)}]' \\ [\hat{\boldsymbol{\mu}}_i^{(t)}]_{\mathbf{x}_i}' \boldsymbol{\Omega}^{(t)} + [\hat{\boldsymbol{\mu}}_i^{(t)}][\hat{\boldsymbol{\mu}}_i^{(t)}]' \end{bmatrix}}_{\hat{\boldsymbol{\tau}}_i^{(t)}}$$

$$\Rightarrow Q_t(\boldsymbol{\Sigma}) = -\frac{1}{2} \left\{ n \log |\boldsymbol{\Sigma}| + \text{tr}[\boldsymbol{\Sigma}^{-1} (\mathbf{Y}'_1 \mathbf{Y}_1 + \sum_{i \in \mathcal{S}_0} \hat{\boldsymbol{\tau}}_i^{(t)})] \right\}.$$

Multivariate Normal: EM Algorithm

► **Model:** $\mathbf{y} = (\mathbf{x}, \mathbf{z}) \sim \mathcal{N} \left\{ \mathbf{0}, \begin{bmatrix} \boldsymbol{\Sigma}_{xx} & \boldsymbol{\Sigma}_{xz} \\ \boldsymbol{\Sigma}_{zx} & \boldsymbol{\Sigma}_{zz} \end{bmatrix} \right\}, \quad \begin{aligned} \mathbf{x} &= (x_1, \dots, x_p), \\ \mathbf{z} &= (z_1, \dots, z_q). \end{aligned}$

► **Observed Data:** $\mathbf{y}_{\text{obs}} = \mathcal{D} = (\mathbf{X}, \mathbf{Z}_1, \delta), \quad \mathbf{Z}_1 = \{\mathbf{z}_i : \delta_i = 1\}.$

► **E-Step:**

$$Q_t(\boldsymbol{\Sigma}) = -\frac{1}{2} \left\{ n \log |\boldsymbol{\Sigma}| + \text{tr} \left[\boldsymbol{\Sigma}^{-1} (\mathbf{Y}'_1 \mathbf{Y}_1 + \sum_{i \in \mathcal{S}_0} \hat{\mathbf{T}}_i^{(t)}) \right] \right\}.$$

► **M-Step:**

$$\hat{\boldsymbol{\Sigma}}^{(t+1)} = \frac{1}{n} \left(\mathbf{Y}'_1 \mathbf{Y}_1 + \sum_{i \in \mathcal{S}_0} \hat{\mathbf{T}}_i^{(t)} \right)$$

(Since $Q_t(\boldsymbol{\Sigma})$ has same shape as loglikelihood of $\mathbf{y}_i \stackrel{\text{iid}}{\sim} \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma})$)

Example: Mixture of Exponential Families

- ▶ **Exponential Family:** $\mathbf{y} \sim g(\mathbf{y} | \boldsymbol{\eta}) = \exp\{\mathbf{T}'\boldsymbol{\eta} - \Psi(\boldsymbol{\eta})\} \cdot h(\mathbf{y})$.
- ▶ **Mixture Model:** The K -component mixture model is

$$f(\mathbf{y} | \boldsymbol{\theta}) = \sum_{k=1}^K \rho_k \cdot g(\mathbf{y} | \boldsymbol{\eta}_k),$$

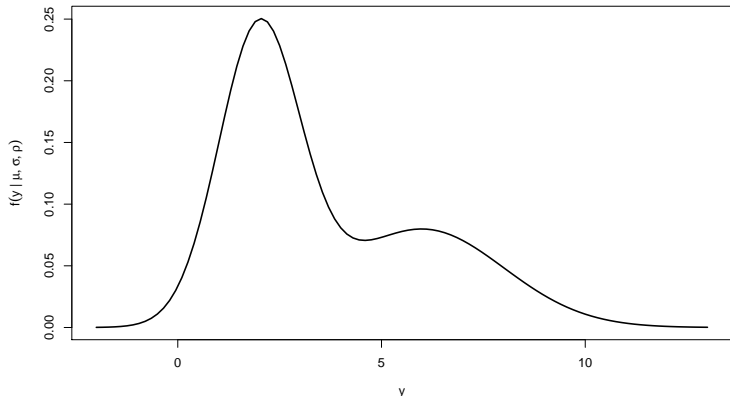
where $\boldsymbol{\Lambda} = (\boldsymbol{\eta}_1, \dots, \boldsymbol{\eta}_K)$, $\boldsymbol{\rho} = (\rho_1, \dots, \rho_K)$, and $\rho_k \geq 0$, $\sum_{k=1}^K \rho_k = 1$.

Example: Mixture of Exponential Families

► **Model:** $f(y | \Lambda, \rho) = \sum_{k=1}^K \rho_k \cdot g(y | \eta_k), \quad g(y | \eta) = \exp\{T'\eta - \Psi(\eta)\}h(y).$

► **Example:** $K = 2, \quad g(y | \eta_k) \cong \mathcal{N}(\mu_k, \sigma_k^2),$

$$\mu = (2, 6), \quad \sigma = (1, 2), \quad \rho = (.6, .4)$$



Example: Mixture of Exponential Families

► **Model:**
$$f(\mathbf{y} | \boldsymbol{\Lambda}, \boldsymbol{\rho}) = \sum_{k=1}^K \rho_k \cdot g(\mathbf{y} | \boldsymbol{\eta}_k), \quad g(\mathbf{y} | \boldsymbol{\eta}) = \exp\{\mathbf{T}'\boldsymbol{\eta} - \Psi(\boldsymbol{\eta})\}h(\mathbf{y}).$$

► **Applications:**

1. **Density Estimation:** For large enough K , mixture model is arbitrarily accurate approximate to any data-generating process $\mathbf{y} \sim f_0(\mathbf{y})$ with same support.
2. **Classification:** To simulate $\mathbf{y} \sim f(\mathbf{y} | \boldsymbol{\Lambda}, \boldsymbol{\rho})$:

$$\mathbf{z} = (z_1, \dots, z_K) \stackrel{\text{iid}}{\sim} \text{Multinomial}(1, \boldsymbol{\rho})$$

$$\mathbf{y} | \mathbf{z} \stackrel{\text{ind}}{\sim} g(\mathbf{y} | \boldsymbol{\eta}_z), \quad \boldsymbol{\eta}_z \text{ is } \boldsymbol{\eta}_k \text{ for which } z_k = 1$$

$$\implies \Pr(\mathbf{y} \text{ is in group } k | \mathbf{y}, \boldsymbol{\Lambda}, \boldsymbol{\rho}) = \Pr(z_k = 1 | \mathbf{y}, \boldsymbol{\Lambda}, \boldsymbol{\rho})$$

$$\begin{aligned} \text{by Bayes Formula: } \Pr(A | B) &= \frac{\Pr(B | A) \Pr(A)}{\Pr(B)} &= \frac{\Pr(\mathbf{y} | z_k = 1, \boldsymbol{\Lambda}, \boldsymbol{\rho}) \Pr(z_k = 1, \boldsymbol{\Lambda}, \boldsymbol{\rho})}{f(\mathbf{y} | \boldsymbol{\Lambda}, \boldsymbol{\rho})} \\ & &= \frac{\rho_k \cdot g(\mathbf{y} | \boldsymbol{\eta}_k)}{\sum_{j=1}^K \rho_j \cdot g(\mathbf{y} | \boldsymbol{\eta}_j)} \end{aligned}$$

Example: Mixture of Exponential Families

- ▶ **Model:** $f(\mathbf{y} | \boldsymbol{\Lambda}, \boldsymbol{\rho}) = \sum_{k=1}^K \rho_k \cdot g(\mathbf{y} | \boldsymbol{\eta}_k), \quad g(\mathbf{y} | \boldsymbol{\eta}) = \exp\{\mathbf{T}'\boldsymbol{\eta} - \Psi(\boldsymbol{\eta})\}h(\mathbf{y}).$
- ▶ **Inference:** Estimate $\boldsymbol{\theta} = (\boldsymbol{\Lambda}, \boldsymbol{\rho})$ given $\mathbf{Y} = (\mathbf{y}_1, \dots, \mathbf{y}_n), \mathbf{y}_i \stackrel{\text{iid}}{\sim} f(\mathbf{y} | \boldsymbol{\Lambda}, \boldsymbol{\rho}).$

Example: Mixture of Exponential Families

- ▶ **Model:** $f(\mathbf{y} | \boldsymbol{\Lambda}, \boldsymbol{\rho}) = \sum_{k=1}^K \rho_k \cdot g(\mathbf{y} | \boldsymbol{\eta}_k), \quad g(\mathbf{y} | \boldsymbol{\eta}) = \exp\{\mathbf{T}'\boldsymbol{\eta} - \Psi(\boldsymbol{\eta})\}h(\mathbf{y}).$
- ▶ **Inference:** Estimate $\boldsymbol{\theta} = (\boldsymbol{\Lambda}, \boldsymbol{\rho})$ given $\mathbf{Y} = (\mathbf{y}_1, \dots, \mathbf{y}_n), \mathbf{y}_i \stackrel{\text{iid}}{\sim} f(\mathbf{y} | \boldsymbol{\Lambda}, \boldsymbol{\rho}).$
- ▶ **Simulation:**

$$\mathbf{z}_i = (z_{i1}, \dots, z_{iK}) \stackrel{\text{iid}}{\sim} \text{Multinomial}(\mathbf{1}, \boldsymbol{\rho})$$

$$\mathbf{y}_i | \mathbf{z}_i \stackrel{\text{ind}}{\sim} g(\mathbf{y} | \boldsymbol{\eta}_{\mathbf{z}_i}), \quad \boldsymbol{\eta}_{\mathbf{z}_i} \text{ is } \boldsymbol{\eta}_k \text{ for which } z_{ik} = 1$$

Example: Mixture of Exponential Families

- ▶ **Model:** $f(\mathbf{y} | \boldsymbol{\Lambda}, \boldsymbol{\rho}) = \sum_{k=1}^K \rho_k \cdot g(\mathbf{y} | \boldsymbol{\eta}_k), \quad g(\mathbf{y} | \boldsymbol{\eta}) = \exp\{\mathbf{T}'\boldsymbol{\eta} - \Psi(\boldsymbol{\eta})\}h(\mathbf{y}).$
- ▶ **Inference:** Estimate $\boldsymbol{\theta} = (\boldsymbol{\Lambda}, \boldsymbol{\rho})$ given $\mathbf{Y} = (\mathbf{y}_1, \dots, \mathbf{y}_n), \mathbf{y}_i \stackrel{\text{iid}}{\sim} f(\mathbf{y} | \boldsymbol{\Lambda}, \boldsymbol{\rho}).$
- ▶ **Simulation:**

$$\mathbf{z}_i = (z_{i1}, \dots, z_{iK}) \stackrel{\text{iid}}{\sim} \text{Multinomial}(\mathbf{1}, \boldsymbol{\rho})$$

$$\mathbf{y}_i | \mathbf{z}_i \stackrel{\text{ind}}{\sim} g(\mathbf{y} | \boldsymbol{\eta}_{\mathbf{z}_i}), \quad \boldsymbol{\eta}_{\mathbf{z}_i} \text{ is } \boldsymbol{\eta}_k \text{ for which } z_{ik} = 1$$

- ▶ Suggests that the EM setup would be
 - ▶ $\mathbf{y}_{\text{comp}} = (\mathbf{Y}, \mathbf{Z})$
 - ▶ $\mathbf{y}_{\text{obs}} = \mathbf{Y}$
 - ▶ $\mathbf{y}_{\text{miss}} = \mathbf{Z} = (\mathbf{z}_1, \dots, \mathbf{z}_n).$

Mixture of EFs: EM Algorithm

► **Model:** $y_i \stackrel{\text{iid}}{\sim} f(y | \boldsymbol{\Lambda}, \boldsymbol{\rho}) = \sum_{k=1}^K \rho_k \cdot g(y | \boldsymbol{\eta}_k), \quad g(y | \boldsymbol{\eta}) = \exp\{\mathbf{T}'\boldsymbol{\eta} - \Psi(\boldsymbol{\eta})\}h(y).$

► **Complete Data:** $y_i | z_i \stackrel{\text{ind}}{\sim} g(y | \boldsymbol{\eta}_{z_i}), \quad z_i \stackrel{\text{iid}}{\sim} \text{Multinomial}(1, \boldsymbol{\rho}).$

► **Complete Data Log-likelihood:**

$$\begin{aligned} \ell(\boldsymbol{\theta} | \mathbf{Y}, \mathbf{Z}) &= \sum_{i=1}^n \underbrace{[\mathbf{T}'_i \boldsymbol{\eta}_{z_i} - \Psi(\boldsymbol{\eta}_{z_i})]}_{\text{Exponential Family}} + \underbrace{\sum_{i=1}^n \sum_{k=1}^K z_{ik} \log(\rho_k)}_{\text{Multinomial}} \\ &= \sum_{i=1}^n \sum_{k=1}^K z_{ik} \left[\mathbf{T}'_i \boldsymbol{\eta}_k - \Psi(\boldsymbol{\eta}_k) + \log(\rho_k) \right] \\ &= \sum_{k=1}^K \sum_{i=1}^n z_{ik} \left[\mathbf{T}'_i \boldsymbol{\eta}_k - \Psi(\boldsymbol{\eta}_k) + \log(\rho_k) \right]. \end{aligned}$$

Mixture of EFs: EM Algorithm

► **Model:** $\mathbf{y}_i \stackrel{\text{iid}}{\sim} f(\mathbf{y} | \boldsymbol{\Lambda}, \boldsymbol{\rho}) = \sum_{k=1}^K \rho_k \cdot g(\mathbf{y} | \boldsymbol{\eta}_k), \quad g(\mathbf{y} | \boldsymbol{\eta}) = \exp\{\mathbf{T}'\boldsymbol{\eta} - \Psi(\boldsymbol{\eta})\}h(\mathbf{y}).$

► **Complete Data:** $\mathbf{y}_i | \mathbf{z}_i \stackrel{\text{ind}}{\sim} g(\mathbf{y} | \boldsymbol{\eta}_{\mathbf{z}_i}), \quad \mathbf{z}_i \stackrel{\text{iid}}{\sim} \text{Multinomial}(1, \boldsymbol{\rho}).$

► **Complete Data Log-likelihood:**

$$\ell(\boldsymbol{\theta} | \mathbf{Y}, \mathbf{Z}) = \sum_{k=1}^K \sum_{i=1}^n z_{ik} \left[\mathbf{T}'_i \boldsymbol{\eta}_k - \Psi(\boldsymbol{\eta}_k) + \log(\rho_k) \right].$$

► **E-Step:** $Q_t(\boldsymbol{\theta}) = \sum_{k=1}^K \sum_{i=1}^n E[z_{ik} | \mathbf{y}_i, \hat{\boldsymbol{\theta}}^{(t)}] \left[\mathbf{T}'_i \boldsymbol{\eta}_k - \Psi(\boldsymbol{\eta}_k) + \log(\rho_k) \right].$

To calculate the expectation, note that $z_{ik} \in \{0, 1\}$, such that

$$E[z_{ik} | \mathbf{y}_i, \hat{\boldsymbol{\theta}}^{(t)}] = \Pr(z_{ik} = 1 | \mathbf{y}_i, \hat{\boldsymbol{\theta}}^{(t)}) = \frac{\hat{\rho}_k^{(t)} g(\mathbf{y}_i | \hat{\boldsymbol{\eta}}_k^{(t)})}{\sum_{j=1}^K \hat{\rho}_j^{(t)} g(\mathbf{y}_i | \hat{\boldsymbol{\eta}}_j^{(t)})} = \hat{q}_{ik}^{(t)}.$$

Mixture of EFs: EM Algorithm

► **Model:** $y_i \stackrel{\text{iid}}{\sim} f(y | \boldsymbol{\Lambda}, \boldsymbol{\rho}) = \sum_{k=1}^K \rho_k \cdot g(\mathbf{y} | \boldsymbol{\eta}_k), \quad g(\mathbf{y} | \boldsymbol{\eta}) = \exp\{\mathbf{T}'\boldsymbol{\eta} - \Psi(\boldsymbol{\eta})\}h(\mathbf{y}).$

► **Complete Data:** $\mathbf{y}_i | z_i \stackrel{\text{iid}}{\sim} g(\mathbf{y} | \boldsymbol{\eta}_{z_i}), \quad z_i \stackrel{\text{iid}}{\sim} \text{Multinomial}(1, \boldsymbol{\rho}).$

► **Complete Data Log-likelihood:**

$$\ell(\boldsymbol{\theta} | \mathbf{Y}, \mathbf{Z}) = \sum_{k=1}^K \sum_{i=1}^n z_{ik} \left[\mathbf{T}'_i \boldsymbol{\eta}_k - \Psi(\boldsymbol{\eta}_k) + \log(\rho_k) \right].$$

► **E-Step:** $Q_t(\boldsymbol{\theta}) = \sum_{k=1}^K \sum_{i=1}^n \hat{q}_{ik}^{(t)} \left[\mathbf{T}'_i \boldsymbol{\eta}_k - \Psi(\boldsymbol{\eta}_k) + \log(\rho_k) \right]$
 $= \sum_{k=1}^K \left[\hat{\mathbf{T}}_k^{(t)'} \boldsymbol{\eta}_k - q_k^{(t)} \Psi(\boldsymbol{\eta}_k) + q_k^{(t)} \log(\rho_k) \right],$

where $\hat{\mathbf{T}}_k^{(t)} = \sum_{i=1}^n \hat{q}_{ik}^{(t)} \mathbf{T}_i$ and $q_k^{(t)} = \sum_{i=1}^n \hat{q}_{ik}^{(t)}$.

Mixture of EFs: EM Algorithm

► **Model:** $y_i \stackrel{\text{iid}}{\sim} f(y | \Lambda, \rho) = \sum_{k=1}^K \rho_k \cdot g(y | \eta_k), \quad g(y | \eta) = \exp\{\mathbf{T}'\eta - \Psi(\eta)\}h(y).$

► **Complete Data:** $y_i | z_i \stackrel{\text{iid}}{\sim} g(y | \eta_{z_i}), \quad z_i \stackrel{\text{iid}}{\sim} \text{Multinomial}(1, \rho).$

► **E-Step:** $Q_t(\theta) = \sum_{k=1}^K \left[\hat{\mathbf{T}}_k^{(t)'} \boldsymbol{\eta}_k - q_k^{(t)} \Psi(\boldsymbol{\eta}_k) + q_k^{(t)} \log(\rho_k) \right].$

► **M-Step:**

► **EF Parameters:** $\boldsymbol{\eta}_k^{(t+1)} = \arg \max_{\boldsymbol{\eta}} \left[\hat{\mathbf{T}}_k^{(t)'} \boldsymbol{\eta} - q_k^{(t)} \Psi(\boldsymbol{\eta}) \right]$, i.e., separable convex optimization problems.

► **Mixing Parameters:** $\hat{\rho}^{(t+1)} = \arg \max_{\rho} \sum_{k=1}^K q_k^{(t)} \log(\rho_k).$

Actually a $K - 1$ dimensional optimization since $\rho_K = 1 - \sum_{k=1}^{K-1} \rho_k$. Similarly, by definition $q_K^{(t)} = 1 - \sum_{k=1}^{K-1} q_k^{(t)}$, such that with $\mathbf{q}^{(t)} = (q_1^{(t)}, \dots, q_K^{(t)})$,

$$\left. \frac{\partial}{\partial \rho_j} \sum_{k=1}^K q_k^{(t)} \log(\rho_k) \right|_{\rho = \hat{\rho}^{(t+1)}} = \frac{q_j^{(t)}}{\hat{\rho}_j^{(t+1)}} - \frac{1 - \sum_{k=1}^{K-1} q_k^{(t)}}{1 - \sum_{k=1}^{K-1} \hat{\rho}_k^{(t+1)}} = 0 \iff \hat{\rho}^{(t+1)} = \mathbf{q}^{(t)}.$$

Example: Probit Regression

► **Logistic Regression:**

$$y_i | \mathbf{x}_i \stackrel{\text{ind}}{\sim} \text{Bernoulli}(\rho_i), \quad \rho = \frac{1}{1 + \exp(-\mathbf{x}'_i \boldsymbol{\beta})}.$$

► **Probit Regression:** $\rho_i = \Phi(\mathbf{x}'_i \boldsymbol{\beta})$, where Φ is the CDF of $\mathcal{N}(0, 1)$.

Can think of this as $z_i \stackrel{\text{ind}}{\sim} \mathcal{N}(\mathbf{x}'_i \boldsymbol{\beta}, 1)$, and $y_i = \mathbb{1}\{z_i > 0\}$ since

$$\Pr(y = 1 | \mathbf{x}) = \Pr(z > 0 | \mathbf{x}) = \Pr(\underbrace{z - \mathbf{x}' \boldsymbol{\beta}}_{\mathcal{N}(0,1)} > -\mathbf{x}' \boldsymbol{\beta} | \mathbf{x}) = \Phi(\mathbf{x}' \boldsymbol{\beta}).$$

⇒ the EM setup is

$$\mathbf{y}_{\text{obs}} = (\mathbf{y}, \mathbf{X}), \quad \mathbf{y}_{\text{comp}} = (\mathbf{z}, \mathbf{y}, \mathbf{X}), \quad \mathbf{y}_{\text{miss}} = \mathbf{z}.$$

Probit Regression: EM Algorithm

- ▶ **Probit Regression:** $y_i | \mathbf{x}_i \stackrel{\text{ind}}{\sim} \text{Bernoulli}(\rho_i),$ $\rho_i = \Pr(Z < \mathbf{x}'_i \boldsymbol{\beta}),$
 $Z \sim \mathcal{N}(0, 1).$
- ▶ **Complete Data:** $z_i | \mathbf{x}_i \stackrel{\text{ind}}{\sim} \mathcal{N}(\mathbf{x}'_i \boldsymbol{\beta}),$ $y_i = \mathbb{1}\{z_i > 0\}.$
- ▶ **Complete Data Likelihood:** $\ell(\boldsymbol{\beta} | \mathbf{z}, \mathbf{y}, \mathbf{X}) = -\frac{1}{2} \sum_{i=1}^n (z_i - \mathbf{x}'_i \boldsymbol{\beta})^2$
- ▶ **E-Step:** $Q_t(\boldsymbol{\beta}) = E[\ell(\boldsymbol{\beta} | \mathbf{z}, \mathbf{y}, \mathbf{X}) | \mathbf{y}, \mathbf{X}, \hat{\boldsymbol{\beta}}^{(t)}]$
 $= -\frac{1}{2} \sum_{i=1}^n E \left[(z_i - \mathbf{x}'_i \boldsymbol{\beta})^2 | y_i, \mathbf{x}_i, \hat{\boldsymbol{\beta}}^{(t)} \right].$

Probit Regression: EM Algorithm

- ▶ **Probit Regression:** $y_i | \mathbf{x}_i \stackrel{\text{ind}}{\sim} \text{Bernoulli}(\rho_i),$ $\rho_i = \Pr(Z < \mathbf{x}'_i \beta),$
 $Z \sim \mathcal{N}(0, 1).$
- ▶ **Complete Data:** $z_i | \mathbf{x}_i \stackrel{\text{ind}}{\sim} \mathcal{N}(\mathbf{x}'_i \beta),$ $y_i = \mathbb{1}\{z_i > 0\}.$
- ▶ **Complete Data Likelihood:** $\ell(\beta | \mathbf{z}, \mathbf{y}, \mathbf{X}) = -\frac{1}{2} \sum_{i=1}^n (z_i - \mathbf{x}'_i \beta)^2$
- ▶ **E-Step:** $Q_t(\beta) = -\frac{1}{2} \sum_{i=1}^n E \left[(z_i - \mathbf{x}'_i \beta)^2 | y_i, \mathbf{x}_i, \hat{\beta}^{(t)} \right].$

To calculate the expectation, note that $\pm(z_i - \mathbf{x}'_i \hat{\beta}^{(t)}) \sim \mathcal{N}(0, 1)$, such that for $y_i = 0$,

$$\begin{aligned} E \left[(z_i - \mathbf{x}'_i \beta)^2 | y_i, \mathbf{x}_i, \hat{\beta}^{(t)} \right] &= E \left[(z_i - \mathbf{x}'_i \beta)^2 | z_i < 0, \mathbf{x}_i, \hat{\beta}^{(t)} \right] \\ &= E \left[\{z_i - \mathbf{x}'_i \hat{\beta}^{(t)} - \mathbf{x}'_i (\beta - \hat{\beta}^{(t)})\}^2 | z_i - \mathbf{x}'_i \hat{\beta}^{(t)} < -\mathbf{x}'_i \hat{\beta}^{(t)} \right] \\ &= E \left[\{Z - \mathbf{x}'_i (\beta - \hat{\beta}^{(t)})\}^2 | Z < -\mathbf{x}'_i \hat{\beta}^{(t)} \right], \quad Z \sim \mathcal{N}(0, 1). \end{aligned}$$

Probit Regression: EM Algorithm

- ▶ **Probit Regression:** $y_i | \mathbf{x}_i \stackrel{\text{ind}}{\sim} \text{Bernoulli}(\rho_i), \quad \rho_i = \Pr(Z < \mathbf{x}'_i \boldsymbol{\beta}),$
 $Z \sim \mathcal{N}(0, 1).$
- ▶ **Complete Data:** $z_i | \mathbf{x}_i \stackrel{\text{ind}}{\sim} \mathcal{N}(\mathbf{x}'_i \boldsymbol{\beta}), \quad y_i = \mathbb{1}\{z_i > 0\}.$
- ▶ **Complete Data Likelihood:** $\ell(\boldsymbol{\beta} | \mathbf{z}, \mathbf{y}, \mathbf{X}) = -\frac{1}{2} \sum_{i=1}^n (z_i - \mathbf{x}'_i \boldsymbol{\beta})^2$
- ▶ **E-Step:** $Q_t(\boldsymbol{\beta}) = -\frac{1}{2} \sum_{i=1}^n E \left[(z_i - \mathbf{x}'_i \boldsymbol{\beta})^2 | y_i, \mathbf{x}_i, \hat{\boldsymbol{\beta}}^{(t)} \right].$

Similarly for $y_i = 1,$

$$\begin{aligned} E \left[(z_i - \mathbf{x}'_i \boldsymbol{\beta})^2 | y_i, \mathbf{x}_i, \hat{\boldsymbol{\beta}}^{(t)} \right] &= E \left[(z_i - \mathbf{x}'_i \boldsymbol{\beta})^2 | z_i > 0, \mathbf{x}_i, \hat{\boldsymbol{\beta}}^{(t)} \right] \\ &= E \left[\{Z - \mathbf{x}'_i (\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}^{(t)})\}^2 | Z > -\mathbf{x}'_i \hat{\boldsymbol{\beta}}^{(t)} \right] \\ &= E \left[\{Z - \mathbf{x}'_i (\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}^{(t)})\}^2 | Z < \mathbf{x}'_i \hat{\boldsymbol{\beta}}^{(t)} \right], \quad Z \sim \mathcal{N}(0, 1), \end{aligned}$$

where we can replace Z by $-1 \times Z$ since $\mathcal{N}(0, 1)$ is symmetric.

Probit Regression: EM Algorithm

► **Probit Regression:** $y_i | \mathbf{x}_i \stackrel{\text{ind}}{\sim} \text{Bernoulli}(\rho_i),$ $\rho_i = \Pr(Z < \mathbf{x}'_i \boldsymbol{\beta}),$
 $Z \sim \mathcal{N}(0, 1).$

► **Complete Data:** $z_i | \mathbf{x}_i \stackrel{\text{ind}}{\sim} \mathcal{N}(\mathbf{x}'_i \boldsymbol{\beta}),$ $y_i = \mathbb{1}\{z_i > 0\}.$

► **Complete Data Likelihood:** $\ell(\boldsymbol{\beta} | \mathbf{z}, \mathbf{y}, \mathbf{X}) = -\frac{1}{2} \sum_{i=1}^n (z_i - \mathbf{x}'_i \boldsymbol{\beta})^2$

► **E-Step:**

$$\begin{aligned} Q_t(\boldsymbol{\beta}) &= -\frac{1}{2} \sum_{i=1}^n E \left[(z_i - \mathbf{x}'_i \boldsymbol{\beta})^2 \mid y_i, \mathbf{x}_i, \hat{\boldsymbol{\beta}}^{(t)} \right] \\ &= -\frac{1}{2} \sum_{i=1}^n \mathcal{G} \left(\mathbf{x}'_i (\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}^{(t)}), (2 \times \mathbb{1}\{y_i = 1\} - 1) \cdot \mathbf{x}'_i \hat{\boldsymbol{\beta}}^{(t)} \right) \end{aligned}$$

where $\mathcal{G}(a, b) = E[(Z - a)^2 \mid Z < b].$

Probit Regression: EM Algorithm

- ▶ **Probit Regression:** $y_i | \mathbf{x}_i \stackrel{\text{ind}}{\sim} \text{Bernoulli}(\rho_i), \quad \rho_i = \Pr(Z < \mathbf{x}'_i \boldsymbol{\beta}),$
 $Z \sim \mathcal{N}(0, 1).$
- ▶ **Complete Data:** $z_i | \mathbf{x}_i \stackrel{\text{ind}}{\sim} \mathcal{N}(\mathbf{x}'_i \boldsymbol{\beta}), \quad y_i = \mathbb{1}\{z_i > 0\}.$
- ▶ **Complete Data Likelihood:** $\ell(\boldsymbol{\beta} | \mathbf{z}, \mathbf{y}, \mathbf{X}) = -\frac{1}{2} \sum_{i=1}^n (z_i - \mathbf{x}'_i \boldsymbol{\beta})^2$
- ▶ **E-Step:** $Q_t(\boldsymbol{\beta}) = -\frac{1}{2} \sum_{i=1}^n \mathcal{G}(\mathbf{x}'_i(\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}^{(t)}), (2 \times \mathbb{1}\{y_i = 1\} - 1) \cdot \mathbf{x}'_i \hat{\boldsymbol{\beta}}^{(t)}),$

where $\mathcal{G}(a, b) = E[(Z - a)^2 | Z < b]$. To calculate $\mathcal{G}(a, b)$

Probit Regression: EM Algorithm

► **Probit Regression:** $y_i | \mathbf{x}_i \stackrel{\text{ind}}{\sim} \text{Bernoulli}(\rho_i),$ $\rho_i = \Pr(Z < \mathbf{x}'_i \boldsymbol{\beta}),$
 $Z \sim \mathcal{N}(0, 1).$

► **Complete Data:** $z_i | \mathbf{x}_i \stackrel{\text{ind}}{\sim} \mathcal{N}(\mathbf{x}'_i \boldsymbol{\beta}),$ $y_i = \mathbb{1}\{z_i > 0\}.$

► **Complete Data Likelihood:**

$$\ell(\boldsymbol{\beta} | \mathbf{z}, \mathbf{y}, \mathbf{X}) = -\frac{1}{2} \sum_{i=1}^n (z_i - \mathbf{x}'_i \boldsymbol{\beta})^2 = -\frac{1}{2} \sum_{i=1}^n z_i^2 - 2(\mathbf{x}'_i \boldsymbol{\beta}) \cdot z_i + (\mathbf{x}'_i \boldsymbol{\beta})^2$$

► **E-Step:**

$$\begin{aligned} Q_t(\boldsymbol{\beta}) &= E[\ell(\boldsymbol{\beta} | \mathbf{z}, \mathbf{y}, \mathbf{X}) | \mathbf{y}, \mathbf{X}, \hat{\boldsymbol{\beta}}^{(t)}] \\ &= -\frac{1}{2} \sum_{i=1}^n \left\{ E[z_i^2 | y_i, \mathbf{x}_i, \hat{\boldsymbol{\beta}}^{(t)}] - 2(\mathbf{x}'_i \boldsymbol{\beta}) \cdot E[z_i | y_i, \mathbf{x}_i, \hat{\boldsymbol{\beta}}^{(t)}] + (\mathbf{x}'_i \boldsymbol{\beta})^2 \right\}. \end{aligned}$$

Probit Regression: EM Algorithm

► **Probit Regression:** $y_i | \mathbf{x}_i \stackrel{\text{ind}}{\sim} \text{Bernoulli}(\rho_i),$ $\rho_i = \Pr(Z < \mathbf{x}'_i \boldsymbol{\beta}),$
 $Z \sim \mathcal{N}(0, 1).$

► **Complete Data:** $z_i | \mathbf{x}_i \stackrel{\text{ind}}{\sim} \mathcal{N}(\mathbf{x}'_i \boldsymbol{\beta}),$ $y_i = \mathbb{1}\{z_i > 0\}.$

► **E-Step:**

$$Q_t(\boldsymbol{\beta}) = -\frac{1}{2} \sum_{i=1}^n \left\{ E[z_i^2 | y_i, \mathbf{x}_i, \hat{\boldsymbol{\beta}}^{(t)}] - 2(\mathbf{x}'_i \boldsymbol{\beta}) \cdot E[z_i | y_i, \mathbf{x}_i, \hat{\boldsymbol{\beta}}^{(t)}] + (\mathbf{x}'_i \boldsymbol{\beta})^2 \right\}.$$

To calculate the expectations, note that $\pm(z_i - \mathbf{x}'_i \hat{\boldsymbol{\beta}}^{(t)}) \sim \mathcal{N}(0, 1)$, such that

$$\begin{aligned} E[z_i | y_i, \mathbf{x}_i, \hat{\boldsymbol{\beta}}^{(t)}] &= \mathbf{x}'_i \hat{\boldsymbol{\beta}}^{(t)} + E[z_i - \mathbf{x}'_i \hat{\boldsymbol{\beta}}^{(t)} | y_i, \mathbf{x}_i, \hat{\boldsymbol{\beta}}^{(t)}] \\ &= \mathbf{x}'_i \hat{\boldsymbol{\beta}}^{(t)} + \begin{cases} E[Z | Z < \mathbf{x}'_i \hat{\boldsymbol{\beta}}^{(t)}] & y_i = 1 \\ E[Z | Z < -\mathbf{x}'_i \hat{\boldsymbol{\beta}}^{(t)}] & y_i = 0, \end{cases} \end{aligned}$$

where $Z \sim \mathcal{N}(0, 1)$, and similarly for $E[z_i^2 | y_i, \mathbf{x}_i, \hat{\boldsymbol{\beta}}^{(t)}]$.

⇒ Need to calculate $g(a) = E[Z | Z < a]$ and $h(a) = E[Z^2 | Z < a]$.

Probit Regression: EM Algorithm

► **Probit Regression:** $y_i | \mathbf{x}_i \stackrel{\text{ind}}{\sim} \text{Bernoulli}(\rho_i)$, $\rho_i = \Pr(Z < \mathbf{x}'_i \boldsymbol{\beta})$,
 $Z \sim \mathcal{N}(0, 1)$.

► **Complete Data:** $z_i | \mathbf{x}_i \stackrel{\text{ind}}{\sim} \mathcal{N}(\mathbf{x}'_i \boldsymbol{\beta})$, $y_i = \mathbb{1}\{z_i > 0\}$.

► **E-Step:** $Q_t(\boldsymbol{\beta}) = -\frac{1}{2} \sum_{i=1}^n \left\{ E[z_i^2 | y_i, \mathbf{x}_i, \hat{\boldsymbol{\beta}}^{(t)}] - 2(\mathbf{x}'_i \boldsymbol{\beta}) \cdot E[z_i | y_i, \mathbf{x}_i, \hat{\boldsymbol{\beta}}^{(t)}] + (\mathbf{x}'_i \boldsymbol{\beta})^2 \right\}$.

► Requires $g(a) = E[Z | Z < a]$ and $h(a) = E[Z^2 | Z < a]$, where $Z \sim \mathcal{N}(0, 1)$.

► Moment-generating function (MGF) of a truncated normal:

$$M(t) = E[e^{Zt} | Z < a] = \frac{\int_{-\infty}^a e^{tz} \cdot e^{-z^2/2} dz}{\int_{-\infty}^a e^{-z^2/2} dz} = \frac{e^{t^2/2} \Phi(a-t)}{\Phi(a)}$$

$$\implies g(a) = \frac{dM(0)}{dt} = -1 \times \frac{\phi(a)}{\Phi(a)}, \quad h(a) = \frac{d^2 M(0)}{dt^2} = 1 - a \times \frac{\phi(a)}{\Phi(a)},$$

where $\phi(z)$ and $\Phi(z)$ are the PDF and CDF of $Z \sim \mathcal{N}(0, 1)$.

Probit Regression: EM Algorithm

► **Probit Regression:** $y_i | \mathbf{x}_i \stackrel{\text{ind}}{\sim} \text{Bernoulli}(\rho_i)$, $\rho_i = \Pr(Z < \mathbf{x}'_i \boldsymbol{\beta})$,
 $Z \sim \mathcal{N}(0, 1)$.

► **Complete Data:** $z_i | \mathbf{x}_i \stackrel{\text{ind}}{\sim} \mathcal{N}(\mathbf{x}'_i \boldsymbol{\beta})$, $y_i = \mathbb{1}\{z_i > 0\}$.

► **E-Step:** After some algebra, get

$$Q_t(\boldsymbol{\beta}) = -\frac{1}{2} \sum_{i=1}^n (\hat{z}_i^{(t)} - \mathbf{x}'_i \boldsymbol{\beta})^2,$$

► **M-Step:** Equivalent to maximizing the likelihood of $\hat{z}_i^{(t)} \stackrel{\text{ind}}{\sim} \mathcal{N}(\mathbf{x}'_i \boldsymbol{\beta}, 1)$

$$\implies \hat{\boldsymbol{\beta}}^{(t+1)} = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}' \hat{\mathbf{z}}^{(t)}.$$

Example: Multivariate t-Distribution

- **Definition:** Let $\mathbf{z} = (z_1, \dots, z_d) \sim \mathcal{N}(\mathbf{0}, \Sigma)$ $\perp\!\!\!\perp$ $x \sim \chi^2_{(\nu)}$. Then

$$\mathbf{y} = \frac{\mathbf{z}}{\sqrt{x/\nu}} + \boldsymbol{\mu}$$

has a multivariate Student-t distribution, denoted $\mathbf{y} \sim t_{(\nu)}(\boldsymbol{\mu}, \Sigma)$.

- **EM Setup:** To simulate observations $\mathbf{y}_i \stackrel{\text{iid}}{\sim} t_{(\nu)}(\boldsymbol{\mu}, \Sigma)$, do

$$x_i \stackrel{\text{iid}}{\sim} \chi^2_{(\nu)}$$
$$\mathbf{y}_i | x_i \stackrel{\text{ind}}{\sim} \mathcal{N}(\boldsymbol{\mu}, \nu \Sigma / x_i).$$

This suggests the setup for EM is

- $\mathbf{y}_{\text{obs}} = \mathbf{Y} = (\mathbf{y}_1, \dots, \mathbf{y}_n)$.
- $\mathbf{y}_{\text{comp}} = (\mathbf{Y}, \mathbf{x})$, where $\mathbf{x} = (x_1, \dots, x_n)$.
- $\mathbf{y}_{\text{miss}} = \mathbf{x}$.

Multivariate t: EM Algorithm

► **Model:** $y_i \stackrel{\text{iid}}{\sim} t_{(\nu)}(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \cong \mathcal{N}(0, \boldsymbol{\Sigma}) / \sqrt{\chi_{(\nu)}^2 / \nu} + \boldsymbol{\mu}$.

► **Complete Data:** $x_i \stackrel{\text{iid}}{\sim} \chi_{(\nu)}^2$, $\mathbf{y}_i | x_i \stackrel{\text{iid}}{\sim} \mathcal{N}(\boldsymbol{\mu}, \nu \boldsymbol{\Sigma} / x_i)$.

► **Complete Data Likelihood:** With $\boldsymbol{\Omega} = \nu \boldsymbol{\Sigma}$ and $\boldsymbol{\theta} = (\boldsymbol{\mu}, \boldsymbol{\Omega}, \nu)$,

$$\begin{aligned} \ell(\boldsymbol{\theta} | \mathbf{Y}, \mathbf{x}) = & -\frac{1}{2} \left[n \log |\boldsymbol{\Omega}| + \sum_{i=1}^n x_i \cdot (\mathbf{y}_i - \boldsymbol{\mu})' \boldsymbol{\Omega}^{-1} (\mathbf{y}_i - \boldsymbol{\mu}) \right] \\ & - \frac{1}{2} \left[n\nu \log(2) + 2n \log \Gamma(\nu/2) - \nu \sum_{i=1}^n \log(x_i) \right]. \end{aligned}$$

► **E-Step:**

$$\begin{aligned} Q_t(\boldsymbol{\theta}) = & -\frac{1}{2} \left[n \log |\boldsymbol{\Omega}| + \sum_{i=1}^n E[x_i | \mathbf{y}_i, \hat{\boldsymbol{\theta}}^{(t)}] \cdot (\mathbf{y}_i - \boldsymbol{\mu})' \boldsymbol{\Omega}^{-1} (\mathbf{y}_i - \boldsymbol{\mu}) \right] \\ & - \frac{1}{2} \left[n\nu \log(2) + 2n \log \Gamma(\nu/2) - \nu \sum_{i=1}^n E[\log(x_i) | \mathbf{y}_i, \hat{\boldsymbol{\theta}}^{(t)}] \right]. \end{aligned}$$

Multivariate t: EM Algorithm

- ▶ **Model:** $\mathbf{y}_i \stackrel{\text{iid}}{\sim} t_{(\nu)}(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \cong \mathcal{N}(0, \boldsymbol{\Sigma}) / \sqrt{\chi_{(\nu)}^2 / \nu} + \boldsymbol{\mu}$.
- ▶ **Complete Data:** $x_i \stackrel{\text{iid}}{\sim} \chi_{(\nu)}^2$, $\mathbf{y}_i | x_i \stackrel{\text{iid}}{\sim} \mathcal{N}(\boldsymbol{\mu}, \nu \boldsymbol{\Sigma} / x_i)$.
- ▶ **E-Step:** Requires $E[x | \mathbf{y}, \boldsymbol{\theta}]$ and $E[\log(x) | \mathbf{y}, \boldsymbol{\theta}]$.
- ▶ Conditional distribution of x :

$$\begin{aligned} p(x | \mathbf{y}, \boldsymbol{\theta}) &\propto p(\mathbf{y} | x, \boldsymbol{\theta}) \cdot p(x | \boldsymbol{\theta}) \\ &\propto \exp \left\{ \frac{\nu - 2}{2} \log(x) - \frac{1}{2} x \cdot (\mathbf{y} - \boldsymbol{\mu})' \boldsymbol{\Omega}^{-1} (\mathbf{y} - \boldsymbol{\mu}) + \frac{d}{2} \log(x) \right\} \\ &= \exp \{ (\alpha - 1) \log(x) - \beta \cdot x \}, \end{aligned}$$

where $\alpha = \alpha(\boldsymbol{\theta}) = \frac{1}{2}(\nu + d)$

$$\beta = \beta(\boldsymbol{\theta}) = \frac{1}{2} [(\mathbf{y} - \boldsymbol{\mu})' \boldsymbol{\Omega}^{-1} (\mathbf{y} - \boldsymbol{\mu}) + 1].$$

$\implies x | \mathbf{y} \sim \text{Gamma}(\alpha, \beta)$.

Multivariate t: EM Algorithm

- ▶ **Model:** $\mathbf{y}_i \stackrel{\text{iid}}{\sim} t_{(\nu)}(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \cong \mathcal{N}(0, \boldsymbol{\Sigma}) / \sqrt{\chi_{(\nu)}^2 / \nu} + \boldsymbol{\mu}$.
- ▶ **Complete Data:** $x_i \stackrel{\text{iid}}{\sim} \chi_{(\nu)}^2$, $\mathbf{y}_i | x_i \stackrel{\text{iid}}{\sim} \mathcal{N}(\boldsymbol{\mu}, \nu \boldsymbol{\Sigma} / x_i)$.
- ▶ **E-Step:** Requires $E[x | \mathbf{y}, \boldsymbol{\theta}]$ and $E[\log(x) | \mathbf{y}, \boldsymbol{\theta}]$.
 - ▶ $x | \mathbf{y}, \boldsymbol{\theta} \sim \text{Gamma}(\alpha, \beta)$, where $\alpha = \frac{1}{2}(\nu + d)$
$$\beta = \frac{1}{2}[(\mathbf{y} - \boldsymbol{\mu})' \boldsymbol{\Omega}^{-1} (\mathbf{y} - \boldsymbol{\mu}) + 1].$$
 - ▶ Gamma distribution is an **exponential family**

$$p(x | \mathbf{y}, \boldsymbol{\theta}) = \exp\{\alpha \log(x) - \beta \cdot x - \Psi(\alpha, -\beta)\} \cdot h(x),$$

where $\Psi(\alpha, -\beta) = -\alpha \log(\beta) + \log \Gamma(\alpha)$.

\implies Sufficient statistics are $\mathbf{T} = (\log(x), x)$, such that

$$E[\mathbf{T} | \mathbf{y}] = \nabla \Psi(\alpha, -\beta) = \left(-\log(\beta) + \frac{\Gamma'(\alpha)}{\Gamma(\alpha)}, \frac{\alpha}{\beta} \right).$$

Multivariate t: EM Algorithm

► **Model:** $\mathbf{y}_i \stackrel{\text{iid}}{\sim} t_{(\nu)}(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \cong \mathcal{N}(0, \boldsymbol{\Sigma}) / \sqrt{\chi_{(\nu)}^2 / \nu} + \boldsymbol{\mu}$.

► **Complete Data:** $x_i \stackrel{\text{iid}}{\sim} \chi_{(\nu)}^2$, $\mathbf{y}_i | x_i \stackrel{\text{iid}}{\sim} \mathcal{N}(\boldsymbol{\mu}, \nu \boldsymbol{\Sigma} / x_i)$.

► **E-Step:**

$$Q_t(\boldsymbol{\theta}) = -\frac{1}{2} \left[n \log |\boldsymbol{\Omega}| + \sum_{i=1}^n \hat{x}_i^{(t)} \cdot (\mathbf{y}_i - \boldsymbol{\mu})' \boldsymbol{\Omega}^{-1} (\mathbf{y}_i - \boldsymbol{\mu}) \right] \\ - \frac{1}{2} \left[n\nu \log(2) + 2n \log \Gamma(\nu/2) - \nu \sum_{i=1}^n \hat{w}_i^{(t)} \right],$$

where

$$\hat{x}_i^{(t)} = \frac{\hat{\alpha}^{(t)}}{\hat{\beta}_i^{(t)}}$$

$$\hat{w}_i^{(t)} = -\log \hat{\beta}_i^{(t)} + \frac{\Gamma'(\hat{\alpha}^{(t)})}{\Gamma(\hat{\alpha}^{(t)})}$$

$$\hat{\alpha}^{(t)} = \frac{1}{2}(\hat{\nu}^{(t)} + d)$$

$$\hat{\beta}_i^{(t)} = \frac{1}{2}[(\mathbf{y}_i - \hat{\boldsymbol{\mu}}^{(t)})' \boldsymbol{\Omega}^{-1} (\mathbf{y}_i - \hat{\boldsymbol{\mu}}^{(t)}) + 1].$$

Multivariate t: EM Algorithm

► **Model:** $y_i \stackrel{\text{iid}}{\sim} t(\nu)(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \cong \mathcal{N}(0, \boldsymbol{\Sigma}) / \sqrt{\chi^2(\nu) / \nu} + \boldsymbol{\mu}$.

► **Complete Data:** $x_i \stackrel{\text{iid}}{\sim} \chi^2(\nu)$, $\mathbf{y}_i | x_i \stackrel{\text{iid}}{\sim} \mathcal{N}(\boldsymbol{\mu}, \nu \boldsymbol{\Sigma} / x_i)$.

► **E-Step:**

$$Q_t(\boldsymbol{\theta}) = -\frac{1}{2} \left[n \log |\boldsymbol{\Omega}| + \sum_{i=1}^n \hat{x}_i^{(t)} \cdot (\mathbf{y}_i - \boldsymbol{\mu})' \boldsymbol{\Omega}^{-1} (\mathbf{y}_i - \boldsymbol{\mu}) \right] \\ - \frac{1}{2} \left[n\nu \log(2) + 2n \log \Gamma(\nu/2) - \nu \sum_{i=1}^n \hat{w}_i^{(t)} \right].$$

► **M-Step:**

► $\hat{\boldsymbol{\mu}}^{(t+1)} = \frac{\sum_{i=1}^n \hat{x}_i^{(t)} \mathbf{y}_i}{\sum_{i=1}^n \hat{x}_i^{(t)}}, \quad \hat{\boldsymbol{\Omega}}^{(t+1)} = \frac{\sum_{i=1}^n \hat{x}_i^{(t)} (\mathbf{y}_i - \hat{\boldsymbol{\mu}}^{(t+1)}) (\mathbf{y}_i - \hat{\boldsymbol{\mu}}^{(t+1)})'}{\sum_{i=1}^n \hat{x}_i^{(t)}}.$

► $\hat{\nu}^{(t+1)} = \arg \min_{\nu} \left\{ n\nu \log(2) + 2n \log \Gamma(\nu/2) - \nu \sum_{i=1}^n \hat{w}_i^{(t)} \right\}$, a convex optimization problem.