

Essence of the Latent Variable Models

Shin-ichi Mayekawa

950629,960823,011102,030618,050714

1 Models

1.1 Latent Variable Models with Continuous Latent Variables

1. Factor Analysis Model : continuous observed variable

$\Lambda = \{\boldsymbol{\mu}, \boldsymbol{\Lambda}, \boldsymbol{\Psi}\}$ (grand mean, factor pattern and error variances)

$\boldsymbol{\rho} = \{\boldsymbol{\mu}_x, \boldsymbol{\Phi}\}$ (factor mean and factor dispersion)

$$\mathbf{y} \sim N(\boldsymbol{\mu} + \boldsymbol{\Lambda}\mathbf{x}, \boldsymbol{\Psi}) \quad (1)$$

$$f(\mathbf{y}|\mathbf{x}, \boldsymbol{\Lambda}) = \frac{1}{\sqrt{2\pi|\boldsymbol{\Psi}|}} \exp\left(-\frac{1}{2}(\mathbf{y} - \boldsymbol{\mu} + \boldsymbol{\Lambda}\mathbf{x})' \boldsymbol{\Psi}^{-1}(\mathbf{y} - \boldsymbol{\mu} + \boldsymbol{\Lambda}\mathbf{x})\right) \quad (2)$$

$$\mathbf{x} \sim N(\boldsymbol{\mu}_x, \boldsymbol{\Phi}) \quad (\boldsymbol{\mu}_x = \mathbf{0}, \boldsymbol{\Phi} = \mathbf{I}) \quad (3)$$

$$h(\mathbf{x}|\boldsymbol{\rho}) = \frac{1}{\sqrt{2\pi|\boldsymbol{\Phi}|}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_x)' \boldsymbol{\Phi}^{-1}(\mathbf{x} - \boldsymbol{\mu}_x)\right). \quad (4)$$

2. IRT Model : binary observed variable

$\boldsymbol{\Lambda} = \{\boldsymbol{\lambda}_j, j=1, 2, \dots, p\}$ (item parameters), $\boldsymbol{\rho} = \{\mu_x, \sigma_x^2\}$ (ability mean and variance)

$$\mathbf{y} \sim \text{independent Bernoulli with prob} = P_j = \text{ICC of item } j \text{ at } x \quad (5)$$

$$f(\mathbf{y}|\mathbf{x}, \boldsymbol{\Lambda}) = \prod_{j=1}^p P_j(x|\boldsymbol{\lambda}_j)^{y_j} (1 - P_j(x|\boldsymbol{\lambda}_j))^{1-y_j} \quad (6)$$

$$\mathbf{x} \sim N(\mu_x, \sigma_x^2) \quad (\mu_x = 0, \sigma_x^2 = 1) \quad (7)$$

$$h(\mathbf{x}|\boldsymbol{\rho}) = \frac{1}{\sqrt{2\pi\sigma_x^2}} \exp\left(-\frac{1}{2}(x - \mu_x)^2/\sigma_x^2\right). \quad (8)$$

1.2 Latent Variable Models with Categorical Latent Variables

1. Classification Model : continuous observed variable

$\boldsymbol{\Lambda} = \{\boldsymbol{\mu}_q, \boldsymbol{\Sigma}_q, q=1, 2, \dots, Q\}$ (class means and dispersions),

$\boldsymbol{\rho} = \{\rho_q, q=1, 2, \dots, Q\}$ (prior probabilities)

$$\mathbf{y} \sim N(\boldsymbol{\mu}_q, \boldsymbol{\Sigma}_q), \text{ if } \mathbf{x} = q, \quad (9)$$

$$f(\mathbf{y}|\mathbf{x}^*, \boldsymbol{\Lambda}) = \prod_{q=1}^Q \left[\frac{1}{\sqrt{2\pi|\boldsymbol{\Sigma}_q|}} \exp\left(-\frac{1}{2}(\mathbf{y} - \boldsymbol{\mu}_q)' \boldsymbol{\Sigma}_q^{-1}(\mathbf{y} - \boldsymbol{\mu}_q)\right) \right]^{x_q^*} \quad (10)$$

2. Latent Class Analysis Model : binary observed variable

$\boldsymbol{\Lambda} = \{\boldsymbol{\lambda}_{jq}, j=1, 2, \dots, p, q=1, 2, \dots, Q\}$ (class response probabilities),

$\boldsymbol{\rho} = \{\rho_q, q=1, 2, \dots, Q\}$ (prior probabilities)

$$\mathbf{y} \sim \text{independent Bernoulli with prob} = \boldsymbol{\lambda}_{jq}, \text{ if } \mathbf{x} = q, \quad (11)$$

$$f(\mathbf{y}|\mathbf{x}^*, \boldsymbol{\Lambda}) = \prod_{q=1}^Q \left[\prod_{j=1}^p \boldsymbol{\lambda}_{jq}^{y_j} (1 - \boldsymbol{\lambda}_{jq})^{1-y_j} \right]^{x_q^*} \quad (12)$$

2 Parameter Estimation using EM Algorithm

2.1 EM Algorithm

EM algorithm can be used to maximize the likelihood of the form

$$mL = \prod_{i=1}^N f(\mathbf{y}_i | \Lambda, \rho) \quad (13)$$

$$= \prod_{i=1}^N \int_{\mathbf{X}} f(\mathbf{y}_i, \mathbf{x} | \Lambda, \rho) d\mathbf{x} \quad (14)$$

$$= \prod_{i=1}^N \int_{\mathbf{X}} f(\mathbf{y}_i | \mathbf{x}, \Lambda) h(\mathbf{x} | \rho) d\mathbf{x} \quad (15)$$

or its logarithm

$$\ln mL = \sum_{i=1}^N \ln f(\mathbf{y}_i | \Lambda, \rho) \quad (16)$$

$$= \sum_{i=1}^N \ln \left(\int_{\mathbf{X}} f(\mathbf{y}_i, \mathbf{x} | \Lambda, \rho) d\mathbf{x} \right) \quad (17)$$

Outline

1. Find the initial value of the parameters Λ and ρ
 Λ_0, ρ_0
2. Calculate the log complete data likelihood
 $\ln L(\Lambda, \rho)$
3. Calculate the conditional distribution (posterior distribution) of \mathbf{x} given \mathbf{y} .
 $h(\mathbf{x} | \mathbf{y}, \Lambda_0, \rho_0)$
4. Calculate the expectation of the log complete data likelihood with respect to \mathbf{x} given \mathbf{y} and the current value of the parameters
 $\mathcal{E} \ln L(\Lambda, \rho)$
5. Maximize $\mathcal{E} \ln L(\Lambda, \rho)$ with respect to the parameters
6. Evaluate $\mathcal{E} \ln L(\Lambda, \rho)$ with the updated parameters
7. If the increase of $\mathcal{E} \ln L(\Lambda, \rho)$ is negligible, stop iteration,
else set $\Lambda_0 = \Lambda, \rho_0 = \rho$ and go back to step 2.

2.2 Continuous Latent Variables

Notation:

- $\mathbf{y} = p \times 1$ observed variable
- $\mathbf{x} = r \times 1$ continuous latent variable
- $\mathbf{y}_i, i=1 \dots N$, observed data
- $\mathbf{x}_i, i=1 \dots N$, missing data
- $\Lambda, \rho =$ parameters to be estimated

1. Model

$$f(\mathbf{y} | \mathbf{x}, \Lambda) = \prod_{j=1}^p f_j(y_j | \mathbf{x}, \Lambda) \quad \text{if locally independent.} \quad (18)$$

2. Prior

$$h(\mathbf{x}|\boldsymbol{\rho}) \tag{19}$$

3. Joint

$$f(\mathbf{y}, \mathbf{x}|\boldsymbol{\Lambda}, \boldsymbol{\rho}) = f(\mathbf{y}|\mathbf{x}, \boldsymbol{\Lambda})h(\mathbf{x}|\boldsymbol{\rho}) \tag{20}$$

4. Marginal

$$f(\mathbf{y}|\boldsymbol{\Lambda}, \boldsymbol{\rho}) = \int_{\mathbf{X}} f(\mathbf{y}, \mathbf{x}|\boldsymbol{\Lambda}, \boldsymbol{\rho})d\mathbf{x} \tag{21}$$

5. Posterior

$$h(\mathbf{x}|\mathbf{y}, \boldsymbol{\Lambda}, \boldsymbol{\rho}) = \frac{f(\mathbf{y}, \mathbf{x}|\boldsymbol{\Lambda}, \boldsymbol{\rho})}{f(\mathbf{y}|\boldsymbol{\Lambda}, \boldsymbol{\rho})} \tag{22}$$

$$\propto f(\mathbf{y}, \mathbf{x}|\boldsymbol{\Lambda}, \boldsymbol{\rho}) = f(\mathbf{y}|\mathbf{x}, \boldsymbol{\Lambda})h(\mathbf{x}|\boldsymbol{\rho}) \tag{23}$$

6. Marginal Likelihood

This is the likelihood on the basis of the marginal density of \mathbf{y} .

$$mL(\boldsymbol{\Lambda}, \boldsymbol{\rho}) = \prod_{i=1}^N f(\mathbf{y}_i|\boldsymbol{\Lambda}, \boldsymbol{\rho}) \tag{24}$$

$$= \prod_{i=1}^N \int_{\mathbf{X}} f(\mathbf{y}_i, \mathbf{x}|\boldsymbol{\Lambda}, \boldsymbol{\rho})d\mathbf{x} \tag{25}$$

7. Complete Data Likelihood

This is the likelihood on the basis of the joint density of \mathbf{y} and \mathbf{x} .

$$L(\boldsymbol{\Lambda}, \boldsymbol{\rho}) = \prod_{i=1}^N f(\mathbf{y}_i, \mathbf{x}_i|\boldsymbol{\Lambda}, \boldsymbol{\rho}) \tag{26}$$

$$= \prod_{i=1}^N f(\mathbf{y}_i|\mathbf{x}_i, \boldsymbol{\Lambda})h(\mathbf{x}_i|\boldsymbol{\rho}) \tag{27}$$

8. Log Complete Data Likelihood

$$\ln L(\boldsymbol{\Lambda}, \boldsymbol{\rho}) = \sum_{i=1}^N \ln f(\mathbf{y}_i, \mathbf{x}_i|\boldsymbol{\Lambda}, \boldsymbol{\rho}) \tag{28}$$

$$= \sum_{i=1}^N \ln f(\mathbf{y}_i|\mathbf{x}_i, \boldsymbol{\Lambda}) + \sum_{i=1}^N \ln h(\mathbf{x}_i|\boldsymbol{\rho}) \tag{29}$$

9. Expected Log Complete Data Likelihood

This is the expectation of $\ln L$ with respect to the posterior distribution of \mathbf{X} given \mathbf{Y} . The posterior distribution is treated as known with $\boldsymbol{\Lambda}_0$ and $\boldsymbol{\rho}_0$ being the provisional values of the parameters.

$$\mathcal{E} \ln L(\boldsymbol{\Lambda}, \boldsymbol{\rho}) = \mathcal{E}(\ln L(\boldsymbol{\Lambda}, \boldsymbol{\rho}) | \mathbf{Y}) \tag{30}$$

$$= \sum_{i=1}^N \mathcal{E}(\ln f(\mathbf{y}_i, \mathbf{x}|\boldsymbol{\Lambda}, \boldsymbol{\rho}) | \mathbf{y}_i) \tag{31}$$

$$= \sum_{i=1}^N \mathcal{E}(\ln f(\mathbf{y}_i|\mathbf{x}, \boldsymbol{\Lambda}) | \mathbf{y}_i) + \sum_{i=1}^N \mathcal{E}(\ln h(\mathbf{x}|\boldsymbol{\rho}) | \mathbf{y}_i) \tag{32}$$

$$\begin{aligned}
&= \sum_{i=1}^N \int_{\mathbf{X}} \ln f(\mathbf{y}_i | \mathbf{x}, \boldsymbol{\Lambda}) h(\mathbf{x} | \mathbf{y}_i, \boldsymbol{\Lambda}_0, \boldsymbol{\rho}_0) d\mathbf{x} \\
&\quad + \sum_{i=1}^N \int_{\mathbf{X}} \ln h(\mathbf{x} | \boldsymbol{\rho}) h(\mathbf{x} | \mathbf{y}_i, \boldsymbol{\Lambda}_0, \boldsymbol{\rho}_0) d\mathbf{x}
\end{aligned} \tag{33}$$

In the M-step, this is maximized with respect to $\boldsymbol{\Lambda}$ and $\boldsymbol{\rho}$. $\boldsymbol{\Lambda}_0$ and $\boldsymbol{\rho}_0$ are treated as known.

10. Numerical Integration

The general numerical integration formula is:

$$I = \int_{\mathbf{X}} G(x) H(x) dx \approx \sum_{q=1}^Q G(x_q) w(x_q) \tag{34}$$

where x_q and $w(x_q)$ are, respectively, the quadrature points and the associated weights which depend on the form of H and the range of integration.

- Gauss-Legendre Quadrature

$$G(x) := \ln f(\mathbf{y}_i | x, \boldsymbol{\Lambda}) h(x | \mathbf{y}_i, \boldsymbol{\Lambda}_0, \boldsymbol{\rho}_0), \quad H(x) := 1 \tag{35}$$

- Gauss-Hermite Quadrature

$$G(x) := \ln f(\mathbf{y}_i | x, \boldsymbol{\Lambda}) \frac{h(x | \mathbf{y}_i, \boldsymbol{\Lambda}_0, \boldsymbol{\rho}_0)}{\text{density of } N(\mu, \sigma^2)}, \quad H(x) := \text{density of } N(\mu, \sigma^2) \tag{36}$$

11. Fixed Point Approximation

When

$$\int_{\mathbf{X}} H(\mathbf{x}) d\mathbf{x} = 1, \tag{37}$$

we may use the following approximation:

$$I = \int_{\mathbf{X}} G(\mathbf{x}) H(\mathbf{x}) d\mathbf{x} \approx \sum_{q=1}^Q G(\mathbf{x}_q) H^*(\mathbf{x}_q) \tag{38}$$

where

$$H^*(\mathbf{x}_q) = \frac{H(\mathbf{x}_q)}{\sum_{q=1}^Q H(\mathbf{x}_q)}. \tag{39}$$

Therefore,

$$\begin{aligned}
\mathcal{E} \ln L(\boldsymbol{\Lambda}, \boldsymbol{\rho}) &\approx \sum_{i=1}^N \sum_{q=1}^Q \ln f(\mathbf{y}_i | \mathbf{x}_q, \boldsymbol{\Lambda}) h^*(\mathbf{x}_q | \mathbf{y}_i, \boldsymbol{\Lambda}_0, \boldsymbol{\rho}_0) + \sum_{i=1}^N \sum_{q=1}^Q \ln h(\mathbf{x}_q | \boldsymbol{\rho}) h^*(\mathbf{x}_q | \mathbf{y}_i, \boldsymbol{\Lambda}_0, \boldsymbol{\rho}_0) \tag{40} \\
&= \sum_{q=1}^Q \sum_{i=1}^N h^*(\mathbf{x}_q | \mathbf{y}_i, \boldsymbol{\Lambda}_0, \boldsymbol{\rho}_0) \ln f(\mathbf{y}_i | \mathbf{x}_q, \boldsymbol{\Lambda}) + \sum_{q=1}^Q \left(\sum_{i=1}^N h^*(\mathbf{x}_q | \mathbf{y}_i, \boldsymbol{\Lambda}_0, \boldsymbol{\rho}_0) \right) \ln h(\mathbf{x}_q | \boldsymbol{\rho})
\end{aligned}$$

2.3 Categorical Latent Variables

Notation:

- $\mathbf{y} = p \times 1$ observed variable
- $\mathbf{x} = Q \times 1$ discrete latent variable
- $\mathbf{x}^* = q \times 1$ dummy expanded \mathbf{x} variable (indicator)
- $\mathbf{y}_i, i=1 \dots N$, observed data
- $\mathbf{x}_i, i=1 \dots N$, missing data
- $\Lambda, \rho =$ parameters to be estimated

1. Model

$$f(\mathbf{y}|\mathbf{x}^*, \Lambda) = \prod_{q=1}^Q (f_q(\mathbf{y}|\lambda_q))^{x_q^*} \quad (42)$$

$$= \prod_{q=1}^Q \left(\prod_{j=1}^p f_{jq}(y_j|\lambda_q) \right)^{x_q^*} \quad (43)$$

$$\text{if locally independent.} \quad (44)$$

The subscript q for f is not necessary if all the class shares the same form of pdf with different values of parameters.

2. Prior

$$h(\mathbf{x}^*|\rho) = \prod_{q=1}^Q \rho_q^{x_q^*}, \quad \text{where} \quad \sum_{q=1}^Q \rho_q = 1 \quad (45)$$

3. Joint

$$f(\mathbf{y}, \mathbf{x}^*|\Lambda, \rho) = f(\mathbf{y}|\mathbf{x}^*, \Lambda)h(\mathbf{x}^*|\rho) \quad (46)$$

$$= \prod_{q=1}^Q (f_q(\mathbf{y}|\lambda_q) \rho_q)^{x_q^*} \quad (47)$$

4. Marginal

$$f(\mathbf{y}|\Lambda, \rho) = \sum_{\mathbf{x}^*} f(\mathbf{y}, \mathbf{x}^*|\Lambda, \rho) \quad (48)$$

$$= \sum_{\mathbf{x}^*} \prod_{q=1}^Q (f_q(\mathbf{y}|\lambda_q) \rho_q)^{x_q^*} \quad (49)$$

$$= \sum_{q=1}^Q f_q(\mathbf{y}|\lambda_q) \rho_q \quad (50)$$

5. Posterior

$$h(\mathbf{x}^*|\mathbf{y}, \Lambda, \rho) = \frac{f(\mathbf{y}, \mathbf{x}^*|\Lambda, \rho)}{f(\mathbf{y}|\Lambda)} \quad (51)$$

$$= \frac{\prod_{q=1}^Q (f_q(\mathbf{y}|\lambda_q) \rho_q)^{x_q^*}}{\sum_{q=1}^Q f_q(\mathbf{y}|\lambda_q) \rho_q} = \prod_{q=1}^Q \left(\frac{f_q(\mathbf{y}|\lambda_q) \rho_q}{\sum_{q=1}^Q f_q(\mathbf{y}|\lambda_q) \rho_q} \right)^{x_q^*} \quad (52)$$

$$= \prod_{q=1}^Q h_q(\mathbf{y}, \Lambda, \rho)^{x_q^*} \quad (53)$$

6. Marginal Likelihood

This is the likelihood on the basis of the marginal density of \mathbf{y} .

$$mL(\mathbf{\Lambda}, \boldsymbol{\rho}) = \prod_{i=1}^N f(\mathbf{y}_i | \mathbf{\Lambda}, \boldsymbol{\rho}) \quad (54)$$

$$= \prod_{i=1}^N \left(\sum_{q=1}^Q f_q(\mathbf{y}_i | \boldsymbol{\lambda}_q) \rho_q \right) \quad (55)$$

7. Complete Data Likelihood

This is the likelihood on the basis of the joint density of \mathbf{y} and \mathbf{x} .

$$L(\mathbf{\Lambda}, \boldsymbol{\rho}) = \prod_{i=1}^N f(\mathbf{y}_i, \mathbf{x}_i^* | \mathbf{\Lambda}, \boldsymbol{\rho}) \quad (56)$$

$$= \prod_{i=1}^N \prod_{q=1}^Q (f_q(\mathbf{y}_i | \boldsymbol{\lambda}_q) \rho_q)^{x_{iq}^*} \quad (57)$$

8. Log Complete Data Likelihood

$$\ln L(\mathbf{\Lambda}, \boldsymbol{\rho}) = \sum_{i=1}^N \ln f(\mathbf{y}_i, \mathbf{x}_i^* | \mathbf{\Lambda}, \boldsymbol{\rho}) \quad (58)$$

$$= \sum_{i=1}^N \sum_{q=1}^Q x_{iq}^* \ln f_q(\mathbf{y}_i | \boldsymbol{\lambda}_q) + \sum_{i=1}^N \sum_{q=1}^Q x_{iq}^* \ln \rho_q \quad (59)$$

9. Expected Log Complete Data Likelihood

This is the expectation of $\ln L$ with respect to the posterior distribution of \mathbf{X} given \mathbf{Y} . The posterior distribution is treated as known with $\mathbf{\Lambda}_0$ and $\boldsymbol{\rho}_0$ being the provisional values of the parameters. Note that the posterior expectation of x_{iq}^* is $h_q(\mathbf{y}_i, \mathbf{\Lambda}_0, \boldsymbol{\rho}_0)$.

$$\mathcal{E} \ln L(\mathbf{\Lambda}, \boldsymbol{\rho}) = \mathcal{E}(\ln L | \mathbf{Y}) \quad (60)$$

$$= \mathcal{E} \left(\sum_{i=1}^N \sum_{q=1}^Q x_{iq}^* \ln f_q(\mathbf{y}_i | \boldsymbol{\lambda}_q) + \sum_{i=1}^N \sum_{q=1}^Q x_{iq}^* \ln \rho_q \mid \mathbf{Y} \right) \quad (61)$$

$$= \sum_{i=1}^N \sum_{q=1}^Q h_q(\mathbf{y}_i, \mathbf{\Lambda}_0, \boldsymbol{\rho}_0) \ln f_q(\mathbf{y}_i | \boldsymbol{\lambda}_q) + \sum_{i=1}^N \sum_{q=1}^Q h_q(\mathbf{y}_i, \mathbf{\Lambda}_0, \boldsymbol{\rho}_0) \ln \rho_q \quad (62)$$

$$= \sum_{q=1}^Q \sum_{i=1}^N h_q(\mathbf{y}_i, \mathbf{\Lambda}_0, \boldsymbol{\rho}_0) \ln f_q(\mathbf{y}_i | \boldsymbol{\lambda}_q) + \sum_{q=1}^Q \left(\sum_{i=1}^N h_q(\mathbf{y}_i, \mathbf{\Lambda}_0, \boldsymbol{\rho}_0) \right) \ln \rho_q \quad (63)$$

In the M-step, this is maximized with respect to $\mathbf{\Lambda}$ and $\boldsymbol{\rho}$. $\mathbf{\Lambda}_0$ and $\boldsymbol{\rho}_0$ are treated as known.

$$\text{the first term of } \mathcal{E} \ln L(\mathbf{\Lambda}, \boldsymbol{\rho}) = \sum_{q=1}^Q \sum_{i=1}^N h_q(\mathbf{y}_i, \mathbf{\Lambda}_0, \boldsymbol{\rho}_0) \ln f_q(\mathbf{y}_i | \boldsymbol{\lambda}_q) \quad (64)$$

$$= \sum_{q=1}^Q \sum_{i=1}^N h_q(\mathbf{y}_i, \mathbf{\Lambda}_0, \boldsymbol{\rho}_0) \sum_{j=1}^p \ln f_{jq}(\mathbf{y}_{ij} | \boldsymbol{\lambda}_q) \quad (65)$$