

Bias Reduction PML

2023-11-20

Introduction

Let $\mathbf{Y} = (Y_1, \dots, Y_p)^\top \in \{0, 1\}^p$ be a vector of Bernoulli random variables. Consider a response pattern $\mathbf{y} = (y_1, \dots, y_p)^\top$, where each $y_i \in \{0, 1\}$. The probability of observing such a response pattern is given by the joint distribution

$$\pi = \Pr(\mathbf{Y} = \mathbf{y}) = \Pr(Y_1 = y_1, \dots, Y_p = y_p). \quad (1)$$

Note that there are a total of $R = 2^p$ possible joint probabilities π_r corresponding to all possible two-way response patterns \mathbf{y}_r .

When we consider a parametric model with parameter vector $\boldsymbol{\theta} \in \mathbb{R}^m$, we write $\pi_r(\boldsymbol{\theta})$ to indicate each joint probability, and

$$\boldsymbol{\pi}(\boldsymbol{\theta}) = \begin{pmatrix} \pi_1(\boldsymbol{\theta}) \\ \vdots \\ \pi_R(\boldsymbol{\theta}) \end{pmatrix} \in [0, 1]^R \quad (2)$$

for the vector of joint probabilities, with $\sum_{r=1}^R \pi_r(\boldsymbol{\theta}) = 1$.

Binary factor models

The model of interest is a factor model, commonly used in social statistics. Using an underlying variable (UV) approach, the observed binary responses y_i are manifestations of some latent, continuous variables Y_i^* , $i = 1, \dots, p$. The connection is made as follows:

$$Y_i = \begin{cases} 1 & Y_i^* > \tau_i \\ 0 & Y_i^* \leq \tau_i, \end{cases}$$

where τ_i is the threshold associated with the variable Y_i^* . For convenience, Y_i^* is taken to be standard normal random variables¹. The factor model takes the form

$$\mathbf{Y}^* = \mathbf{\Lambda}\boldsymbol{\eta} + \boldsymbol{\epsilon},$$

where each component is explained below:

- $\mathbf{Y}^* = (Y_1^*, \dots, Y_p^*)^\top \in \mathbf{R}^p$ are the underlying variables;
- $\mathbf{\Lambda} \in \mathbf{R}^{p \times q}$ is the matrix of loadings;
- $\boldsymbol{\eta} = (\eta_1, \dots, \eta_q)^\top \in \mathbf{R}^q$ is the vector of latent factors;
- $\boldsymbol{\epsilon} \in \mathbf{R}^p$ are the error terms associated with the items (aka unique variables).

We also make some distributional assumptions, namely

¹For parameter identifiability, the location and scale of the normal distribution have to be fixed if the thresholds are to be estimated.

1. $\boldsymbol{\eta} \sim N_q(\mathbf{0}, \boldsymbol{\Psi})$, where $\boldsymbol{\Psi}$ is a correlation matrix, i.e. for $k, l \in \{1, \dots, q\}$,

$$\boldsymbol{\Psi}_{kl} = \begin{cases} 1 & \text{if } k = l \\ \rho(\eta_k, \eta_l) & \text{if } k \neq l. \end{cases}$$

2. $\boldsymbol{\epsilon} \sim N_p(\mathbf{0}, \boldsymbol{\Theta}_\epsilon)$, with $\boldsymbol{\Theta}_\epsilon = \mathbf{I} - \text{diag}(\boldsymbol{\Lambda} \boldsymbol{\Psi} \boldsymbol{\Lambda}^\top)$.

These two assumptions, together with $\text{Cov}(\boldsymbol{\eta}, \boldsymbol{\epsilon}) = \mathbf{0}$, implies that $\mathbf{Y}^* \sim N_p(\mathbf{0}, \boldsymbol{\Sigma}_{\mathbf{Y}^*})$, where

$$\boldsymbol{\Sigma}_{\mathbf{Y}^*} = \text{Var}(\mathbf{Y}^*) = \boldsymbol{\Lambda} \boldsymbol{\Phi} \boldsymbol{\Lambda}^\top + \boldsymbol{\Theta}_\epsilon. \quad (3)$$

The parameter vector for this factor model is denoted $\boldsymbol{\theta}^\top = (\boldsymbol{\lambda}, \boldsymbol{\psi}, \boldsymbol{\tau}) \in \mathbb{R}^m$, where it contains the vectors of the free non-redundant parameters in $\boldsymbol{\Lambda}$ and $\boldsymbol{\Psi}$ respectively, as well as the vector of all free thresholds.

Under this factor model, the probability of response pattern \mathbf{y}_r is

$$\pi_r(\boldsymbol{\theta}) = \Pr(\mathbf{Y} = \mathbf{y}_r \mid \boldsymbol{\theta}) \quad (4)$$

$$= \int \cdots \int_A \phi_p(\mathbf{y}^* \mid \mathbf{0}, \boldsymbol{\Sigma}_{\mathbf{Y}^*}) d\mathbf{y}^* \quad (5)$$

where $\phi_p(\cdot \mid \boldsymbol{\mu}, \boldsymbol{\Sigma})$ is the density function of the p -dimensional normal distribution with mean $\boldsymbol{\mu}$ and variance $\boldsymbol{\Sigma}$. This integral is evaluated on the set

$$A = \{\mathbf{Y}^* \in \mathbb{R}^p \mid Y_1 = y_1, \dots, Y_p = y_p\}.$$

Pairwise likelihood estimation

In order to define the pairwise likelihood, let $\pi_{y_i y_j}^{(ij)}(\boldsymbol{\theta})$ be the probability under the model that $Y_i = y_i \in \{0, 1\}$ and $Y_j = y_j \in \{0, 1\}$ for a pair of variables Y_i and Y_j , $i, j = 1, \dots, p$ and $i < j$. The pairwise log-likelihood takes the form

$$\ell_P(\boldsymbol{\theta}) = \sum_{i < j} \sum_{y_i} \sum_{y_j} \hat{n}_{y_i y_j}^{(ij)} \log \pi_{y_i y_j}^{(ij)}(\boldsymbol{\theta}), \quad (6)$$

where $\hat{n}_{y_i y_j}^{(ij)}$ is the observed (weighted) frequency of sample units with $Y_i = y_i$ and $Y_j = y_j$,

$$\hat{n}_{y_i y_j}^{(ij)} = \sum_h w_h [\mathbf{y}_i^{(h)} = y_i, \mathbf{y}_j^{(h)} = y_j].$$

Here the w_h refers to the design weight for any individual h in the sample. For simplicity, we may assume that these weights are normalised such that $\sum w_h = N$. In such a case, a simple random sampling design would imply all weights are equal to one, and the weighted pairwise likelihood reduces to the usual pairwise likelihood function.

Let us also define the corresponding observed pairwise proportions $p_{y_i y_j}^{(ij)} = \hat{n}_{y_i y_j}^{(ij)} / n$. There are a total of $\tilde{R} = 4 \times \binom{p}{2}$ summands, where the ‘4’ is representative of the total number of pairwise combinations of binary choices ‘00’, ‘10’, ‘01’, and ‘11’.

The *pairwise maximum likelihood estimator* $\hat{\boldsymbol{\theta}}_{\text{PL}}$ satisfies $\hat{\boldsymbol{\theta}}_{\text{PL}} = \arg\max_{\boldsymbol{\theta}} \ell_P(\boldsymbol{\theta})$. Under certain regularity conditions,

$$\sqrt{N}(\hat{\boldsymbol{\theta}}_{\text{PL}} - \boldsymbol{\theta}) \xrightarrow{D} N_m(\mathbf{0}, \mathcal{H}(\boldsymbol{\theta}) \mathcal{J}(\boldsymbol{\theta})^{-1} \mathcal{H}(\boldsymbol{\theta})), \quad (7)$$

where

- $\mathcal{H}(\boldsymbol{\theta}) = -\mathbb{E} \nabla^2 \ell_P(\boldsymbol{\theta}; \mathbf{y}^{(h)})$ is the *sensitivity matrix*; and
- $\mathcal{J}(\boldsymbol{\theta}) = \text{Var}(\nabla \ell_P(\boldsymbol{\theta}; \mathbf{y}^{(h)}))$ is the *variability matrix*.

In practice, we may estimate these matrices using the following estimators:

$$\hat{\mathbf{H}} := \mathbf{H}(\hat{\boldsymbol{\theta}}) = -\frac{1}{\sum_h w_h} \sum_{h=1}^N w_h \nabla^2 \ell_P(\boldsymbol{\theta}; \mathbf{y}^{(h)}) \Big|_{\boldsymbol{\theta}=\hat{\boldsymbol{\theta}}}$$

$$\hat{\mathbf{J}} := \mathbf{J}(\hat{\boldsymbol{\theta}}) = \frac{1}{\sum_h w_h} \sum_{h=1}^N w_h^2 \nabla \ell_P(\boldsymbol{\theta}; \mathbf{y}^{(h)}) \nabla \ell_P(\boldsymbol{\theta}; \mathbf{y}^{(h)})^\top \Big|_{\boldsymbol{\theta}=\hat{\boldsymbol{\theta}}}$$

That is, $\hat{\mathbf{H}}$ is the Hessian resulting from the optimisation of the pairwise likelihood function, while $\hat{\mathbf{J}}$ is the cross product of the gradient of the pairwise likelihood function—each evaluated at the maximum PLE. Note that both are considered “unit” information matrices, as they are normalised by the sum of the weights (sample size).

Bias reduction

Define

$$A(\hat{\boldsymbol{\theta}}) = -\frac{1}{2} \nabla \text{tr} (\mathbf{H}(\boldsymbol{\theta})^{-1} \mathbf{J}(\boldsymbol{\theta})^{-1}) \Big|_{\boldsymbol{\theta}=\hat{\boldsymbol{\theta}}}.$$

Then, an improved estimator $\tilde{\boldsymbol{\theta}}$ is given by

$$\tilde{\boldsymbol{\theta}} = \hat{\boldsymbol{\theta}} + \mathbf{H}(\hat{\boldsymbol{\theta}})^{-1} A(\hat{\boldsymbol{\theta}}).$$

Some computational notes:

- The Hessian $\mathbf{H}(\boldsymbol{\theta})$ matrix is obtained as a byproduct of the optimisation routine in `{lavaan}` (or using my manually coded `pml` function and `optim`). There is no explicit code for it.
- The variability matrix $\mathbf{J}(\boldsymbol{\theta})$ is obtained from `{lavaan}`, by tricking it into accepting starting values $\boldsymbol{\theta}$ as converged parameter values and extracting the $\mathbf{J}(\boldsymbol{\theta})$ accordingly.
- Then form the $\mathbf{A}(\boldsymbol{\theta})$ matrix and obtain the gradient using the `numDeriv` package.

Example

Consider $p = 5$ binary items generated using the UV approach as above with the following true parameter values:

- $\boldsymbol{\lambda} = 0.8, 0.7, 0.47, 0.38, 0.34$
- $\boldsymbol{\tau} = -1.43, -0.55, -0.13, -0.72, -1.13$

NB: This is called Model no 1 from a previous work (Jamil, Moustaki, and Skinner 2023).

```
set.seed(123)
model_no <- 1
mod <- txt_mod(model_no)
(dat <- gen_data_bin(model_no))
```

```
## # A tibble: 1,000 x 5
##   y1    y2    y3    y4    y5
##   <ord> <ord> <ord> <ord> <ord>
## 1 1      1      0      1      1
## 2 1      1      1      1      1
## 3 1      0      1      0      1
## 4 1      0      1      1      1
## 5 1      0      1      0      1
```

```
## 6 1      0      0      1      0
## 7 1      1      1      0      1
## 8 1      0      0      1      1
## 9 1      0      1      1      1
## 10 1     1      1      0      1
## # i 990 more rows
```

The fit from the PL routine is obtained from {lavaan}:

```
fit_lav <- sem(mod, dat, std.lv = TRUE, estimator = "PML")
summary(fit_lav)
```

```
## lavaan 0.6.17.1946 ended normally after 17 iterations
##
##      Estimator                      PML
##      Optimization method            NLMINB
##      Number of model parameters      10
##
##      Number of observations          1000
##
## Model Test User Model:
##
##              Standard      Scaled
##      Test Statistic        2.478    3.571
##      Degrees of freedom         5    5.320
##      P-value (Unknown)         NA    0.655
##      Scaling correction factor    0.694
##      mean+var adjusted correction (PML)
##
## Parameter Estimates:
##
##      Standard errors              Sandwich
##      Information bread            Observed
##      Observed information based on    Hessian
##
## Latent Variables:
##
##              Estimate  Std.Err  z-value  P(>|z|)
##      eta1 =~
##      y1          0.696   0.077   9.092   0.000
##      y2          0.700   0.062  11.302   0.000
##      y3          0.447   0.060   7.496   0.000
##      y4          0.493   0.061   8.138   0.000
##      y5          0.417   0.072   5.798   0.000
##
## Intercepts:
##
##              Estimate  Std.Err  z-value  P(>|z|)
##      .y1          0.000
##      .y2          0.000
##      .y3          0.000
##      .y4          0.000
##      .y5          0.000
##      eta1        0.000
##
## Thresholds:
##
##              Estimate  Std.Err  z-value  P(>|z|)
##      y1|t1       -1.522   0.062 -24.640   0.000
```

```
##      y2|t1      -0.550    0.042  -13.135    0.000
##      y3|t1      -0.171    0.040   -4.297    0.000
##      y4|t1      -0.678    0.043  -15.714    0.000
##      y5|t1      -1.190    0.052  -23.011    0.000
##
## Variances:
##              Estimate Std.Err  z-value  P(>|z|)
##      .y1             0.516
##      .y2             0.510
##      .y3             0.800
##      .y4             0.757
##      .y5             0.826
##      eta1            1.000
##
## Scales y*:
##              Estimate Std.Err  z-value  P(>|z|)
##      y1              1.000
##      y2              1.000
##      y3              1.000
##      y4              1.000
##      y5              1.000
```

We compare the fit of the coefficients and the true value:

```
theta_hat <- coef(fit_lav)
theta_true <- c(lavaan.bingof:::get_Lambda(model_no),
               lavaan.bingof:::get_tau(model_no))
tibble(
  coef = names(theta_hat),
  theta_hat = round(theta_hat, 2),
  theta_true = theta_true
) |>
  kbl(booktabs = TRUE)
```

coef	theta_hat	theta_true
eta1=~y1	0.70	0.80
eta1=~y2	0.70	0.70
eta1=~y3	0.45	0.47
eta1=~y4	0.49	0.38
eta1=~y5	0.42	0.34
y1 t1	-1.52	-1.43
y2 t1	-0.55	-0.55
y3 t1	-0.17	-0.13
y4 t1	-0.68	-0.72
y5 t1	-1.19	-1.13

Now apply the bias reduction method

```
# Assume HHH() and JJJ() are the functions to obtain the H and J matrices at a
# given theta
AAA <- function(.theta) {
  tmp <- function(x) {
    Hinv <- HHH(x) |> MASS::ginv()
    J <- JJJ(x)
```

```

    -0.5 * sum(diag(Hinv %*% J))
  }
  numDeriv::grad(tmp, .theta)
}

# Bias reduction
(A <- AAA(theta_hat))

## [1] -0.1296364813  0.1811323957 -0.5299060523  0.0124908790  0.2713109474
## [6]  3.0055899345  0.7016799736 -0.0006567022  1.0727559368  2.2152931318

Hinv <- solve(HHH(theta_hat))
theta_tilde <- theta_hat + Hinv %*% A

```

coef	theta_hat	theta_true	theta_tilde
eta1=~y1	0.70	0.80	0.10
eta1=~y2	0.70	0.70	2.15
eta1=~y3	0.45	0.47	-1.20
eta1=~y4	0.49	0.38	0.57
eta1=~y5	0.42	0.34	1.44
y1 t1	-1.52	-1.43	1.38
y2 t1	-0.55	-0.55	-0.06
y3 t1	-0.17	-0.13	-0.04
y4 t1	-0.68	-0.72	-0.05
y5 t1	-1.19	-1.13	0.41

References

Jamil, Haziq, Irini Moustaki, and Chris Skinner. 2023. “Pairwise Likelihood Estimation and Limited Information Goodness-of-Fit Test Statistics for Binary Factor Analysis Models Under Complex Survey Sampling.” *arXiv Preprint arXiv:2311.02543*.