

Simulation demo

Example from Muthen & Muthen, 2014

Adam Garber

1/28/2020

References:

Muthén, L. K., & Muthén, B. O. (2002). How to use a Monte Carlo study to decide on sample size and determine power. *Structural equation modeling*, 9(4), 599-620.

Hallquist, M. N., & Wiley, J. F. (2018). MplusAutomation: An R Package for Facilitating Large-Scale Latent Variable Analyses in Mplus. *Structural equation modeling: a multidisciplinary journal*, 25(4), 621-638.

Getting started:

1. Create an R-Project
2. Install packages
3. Load packages

R-Project instructions:

1. click “NEW PROJECT” (upper right corner of window)
2. choose option “NEW DIRECTORY”
3. choose location of project (avoid long filepaths)

Within R-studio under the files pane (bottom right):

1. click “New Folder” and name folder “data”
2. click “New Folder” and name folder “efa_mplus”
3. click “New Folder” and name folder “figures”

```
install.packages(c("MplusAutomation",  
                  "tidyverse",  
                  "here",  
                  "glue",  
                  "psych"))
```

IF package rhdf5 does not load then run:

```
if (!requireNamespace("BiocManager", quietly = TRUE))
  install.packages("BiocManager")

BiocManager::install("rhdf5")
```

Simulation demonstration

Factors that influence minimum sample size requirements:

- size of the model (number of parameters)
- distribution of the variables (skew, kurtosis, multi-modal..)
- amount of missing data & pattern of missingness
- reliability of the variables
- strength of the relations among the variables

“This article focuses on parameter estimates, standard errors, coverage, and power assuming correctly specified models. Misspecified models can also be studied in the Mplus Monte Carlo framework but are not included here.” - Muthen & Muthen, 2014

CFA Model example:

- 2 factors
- 10 indicators (5 per factor)
- 31 free parameters & 24 *df*
- factor loadings = 0.8 (freely estimated in the model)
- residual variances = 0.36 (error)
- factor variances = 1 (fixed)
- reliability of factor = $0.64 = .8^2(1) / .8^2(1 + 0.36)$

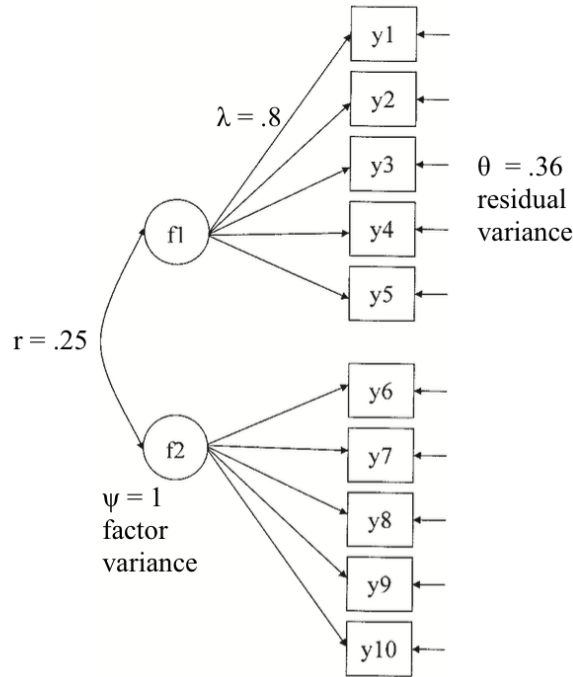


Figure 1. CFA model example. Picture adapted from, Muthen & Muthen 2104

Monte Carlo conditions

1. normally distributed continuous factor indicators without missing data
2. normally distributed continuous factor indicators with missing data
3. non-normal continuous factor indicators without missing data
4. non-normal continuous factor indicators with missing data

Missing data:

“all participants have data on $y_1, y_2, y_3, y_4,$ and y_5 , and 50% of the participants have data on $y_6, y_7, y_8, y_9,$ and y_{10} ”

Non-normal data (how to engineer skew & kurtosis):

- “non-normal data are created using a mixture of two normal subpopulations or classes of individuals. Normal data are generated for two classes that have different means and variances for the factor indicators. The combined data are analyzed as though they come from a single population.”
- “The first step is to generate data for two classes such that the combination of the data from the two classes **has the desired skewness and kurtosis.**”
- “For the CFA model with nonnormal data, Class 1, the outlier class, contains 12% of the participants and Class 2 contains the remaining 88%. Only the factor indicators for the second factor are non-normal. Therefore, the Class 1 mean for the second factor is chosen to be 15 and the variance 5, as compared to the Class 2 mean and variance of zero and 1. The resulting population univariate skewness for variables y_6 through y_{10} is 1.2. The resulting population univariate kurtosis for variables y_6 through y_{10} ranges from 1.5 to 1.6.”

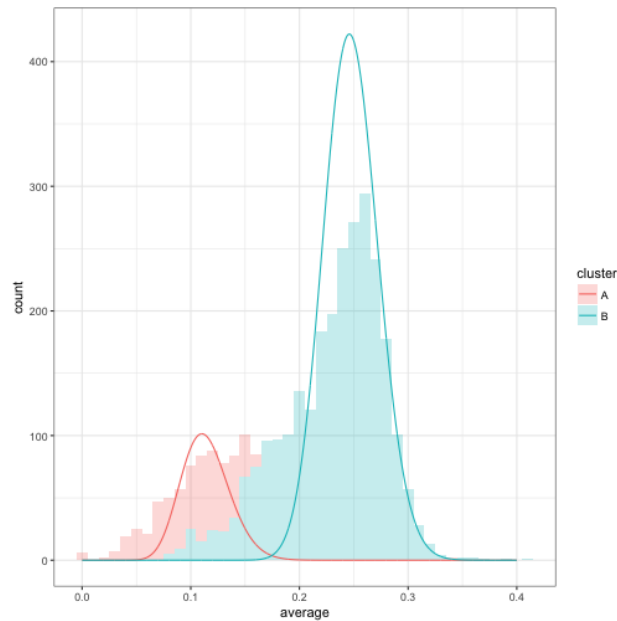


Figure 2. Two class model. What are the data generating processes that might result in skew?

- How many data generating processes can we think of that will result in non-normality?
- Will a single homogenous group always look normally distributed due to noise?

loading packages...

```
library(tidyverse)
library(MplusAutomation)
library(rhdf5)
library(here)
library(glue)
library(psych)
```

Simulation 0 - tuning skew exercise

- “The second step is to run the analysis with one replication and a large sample to obtain approximate population values for the one class model. (e.g., achieve factor indicator reliabilities of 0.64)”
- We will just do 1 replication with an N size of 100,000 (exploit law of large numbers!)

```
# testing & tuning

cfa_tune <- lapply(1:5, function(k) {
  cfa_0 <- mplusObject(

    TITLE = "CFA 1 - non-normal, no missing",

    MONTECARLO =
      sprintf("NAMES ARE y1-y10;
```

```

NOBSERVATIONS = 100000;
NREPS = 1;
SEED = 53487;
CLASSES = C(1);
GENCLASSES = C(2);
SAVE = cfa0_%d.sav;",k),

ANALYSIS =
"TYPE = MIXTURE;
ESTIMATOR = MLR;",

MODELPOPULATION =
glue("%OVERALL%
f1 BY y1-y5*.8;
f2 BY y6-y10*.8;
f1@1 f2@1;
y1-y5*.36 y6-y10*9;
f1 WITH f2*.95;
[C#1@-{k*.5}]; ! parameter we will tune to adjust the size of the outliers

%C#1%          ! outlier class

[f1@0 f2@15]; ! means (factor 2 set to 15 to tune skewness & kurtosis)
f1@1 f2@5;     ! variances (factor 2 set to 5 to tune skewness & kurtosis)

%C#2%          ! majority class

[f1@0 f2@0];
f1@1 f2@1;"),

MODEL =
"%OVERALL%
f1 BY y1-y5*.8;
f2 BY y6-y10*4;
f1@1 f2@1;
y1-y5*.36 y6-y10*9;
f1 WITH f2*.20;

[y6-y10*1.42];" ,

OUTPUT = " SAMPSTAT TECH9;")

cfa_0_fit <- mplusModeler(cfa_0,
                        dataout=here("mplus_tune", "cfa0_demo_sim.dat"),
                        modelout=sprintf(here("mplus_tune", "%d_cfa0_demo_sim.inp"), k),
                        check=TRUE, run = TRUE, hashfilename = FALSE)
})

data0_1 <- read.delim(here("mplus_tune", "cfa0_1.sav"), sep = "",
                      header = FALSE,
                      na.strings = "999.000000")

data0_2 <- read.delim(here("mplus_tune", "cfa0_5.sav"), sep = "",

```

```

        header = FALSE,
        na.strings = "999.000000")

describe(data0_1)

data0_1 %>% ggplot(aes(x=V10)) +
  geom_density()

data0_2 %>% ggplot(aes(x=V10)) +
  geom_density()

```

Criteria to assess level of bias

- “The first criterion is that parameter and standard error biases do not exceed 10% for any parameter in the model.”
- “The second criterion is that the standard error bias for the parameter for which power is being assessed does not exceed 5%”
- “The third criterion is that coverage remains between 0.91 and 0.98.”
- “Once these three conditions are satisfied, the sample size is chosen to keep power close to 0.80.”

Simulation 1

- Normally distributed
- No missing
- Sample size (n) = 150
- Number of repetition = 10,000

```

cfa_1 <- mplusObject(

  TITLE = "CFA 1 - normal, no missing",

  MONTECARLO =
    "NAMES ARE y1-y10;
    NOOBSERVATIONS = 150;
    NREPS = 10000;
    SEED = 53487;
    CLASSES = C(1);
    GENCLASSES = C(1);
    SAVE = cfa1.sav;",

  ANALYSIS =
    "TYPE = MIXTURE;
    ESTIMATOR = ML; ! when normal MLR simplifies to ML",

  MODELPOPULATION =
    "%OVERALL%           !
    f1 BY y1-y5*.8;      ! factor loadings = .8 (average?)
    f2 BY y6-y10*.8;     !
    f1@1 f2@1;           ! factor variances = 1 (fixed)
    y1-y10*.36;          ! residual variances = .36
    f1 WITH f2*.25;      ! factor correlation = .25

```

```

",

MODEL =
  "%OVERALL%
  f1 BY y1-y5*.8;
  f2 BY y6-y10*.8;
  f1@1 f2@1;
  y1-y10*.36;
  f1 WITH f2*.25;" ,

OUTPUT = "TECH9;")

cfa_1_fit <- mplusModeler(cfa_1,
                          dataout=here("mplus_files", "cfa1_demo_sim.dat"),
                          modelout=here("mplus_files", "cfa1_demo_sim.inp"),
                          check=TRUE, run = TRUE, hashfilename = FALSE)

```

Examples from simulation 1 (Y1):

- Parameter bias: columns 1 & 2 $(.8-.7979)/.8 = .0026$
- Standard error bias: columns 2 & 3 $(.0706-.0699)/.0706 = 0.0099$
- Coverage: column 6 = .947
- Power: column 7 = 1.0

Simulation 2

- Normally distributed
- **missing set to 50% for y6 - y10**

```

cfa_2 <- mplusObject(

  TITLE = "CFA 2 - normal, missing (50%)",

  MONTECARLO =
    "NAMES ARE y1-y10;
    NOBSERVATIONS = 175;
    NREPS = 10000;
    SEED = 53487;
    CLASSES = C(1);
    GENCLASSES = C(1);
    PATMISS = y6 (.5) y7 (.5) y8 (.5) y9 (.5) y10 (.5);
    PATPROB = 1;
    SAVE = cfa2.sav;",

  ANALYSIS =
    "TYPE = MIXTURE MISSING;
    ESTIMATOR = ML; ! when normal MLR simplifies to ML",

  MODELPOPULATION =

```

```

"%OVERALL%
f1 BY y1-y5*.8;
f2 BY y6-y10*.8;
f1@1 f2@1;
y1-y10*.36;
f1 WITH f2*.25;";

MODEL =
"%OVERALL%
f1 BY y1-y5*.8;
f2 BY y6-y10*.8;
f1@1 f2@1;
y1-y10*.36;
f1 WITH f2*.25;" ,

OUTPUT = "PATTERNS TECH9;")

cfa_2_fit <- mplusModeler(cfa_2,
                        dataout=here("mplus_files", "cfa2_demo_sim.dat"),
                        modelout=here("mplus_files", "cfa2_demo_sim.inp"),
                        check=TRUE, run = TRUE, hashfilename = FALSE)

```

Simulation 3

- Non-normally distributed
- No missing

```

cfa_3 <- mplusObject(

  TITLE = "CFA 3 - non-normal, no missing",

  MONTECARLO =
    "NAMES ARE y1-y10;
    NOOBSERVATIONS = 265;
    NREPS = 10000;
    SEED = 53487;
    CLASSES = C(1);
    GENCLASSES = C(2);
    SAVE = cfa3.sav;";

  ANALYSIS =
    "TYPE = MIXTURE;
    ESTIMATOR = MLR;";

  MODELPOPULATION =
    "%OVERALL%
    f1 BY y1-y5*.8;
    f2 BY y6-y10*.8;
    f1@1 f2@1;

```



```

y1-y5*.36 y6-y10*9;
f1 WITH f2*.95;
[C#1@-2];

%C#1%          ! outlier class (size = 12%)

[f1@0 f2@15]; ! means (facotr 2 set to 15 to tune skewness & kurtosis)
f1@1 f2@5;    ! variances (facotr 2 set to 5 to tune skewness & kurtosis)

%C#2%          ! majority class (size = 88%)

[f1@0 f2@0];
f1@1 f2@1;";

MODEL =
"%OVERALL%
f1 BY y1-y5*.8;
f2 BY y6-y10*4;
f1@1 f2@1;
y1-y5*.36 y6-y10*9;
f1 WITH f2*.20;

[y6-y10*1.42];" ,

OUTPUT = "TECH9;")

cfa_3_fit <- mplusModeler(cfa_3,
                          dataout=here("mplus_files", "cfa3_demo_sim.dat"),
                          modelout=here("mplus_files", "cfa3_demo_sim.inp"),
                          check=TRUE, run = TRUE, hashfilename = FALSE)

```

Simulation 4

- Non-normally distributed
- missing set to 50% for y6 - y10

```

cfa_4 <- mplusObject(

  TITLE = "CFA 4 - non-normal, missing (50%)",

  MONTECARLO =
  "NAMES ARE y1-y10;
  NOBSEVATIONS = 315;
  NREPS = 10000;
  SEED = 53487;
  CLASSES = C(1);
  GENCLASSES = C(2);
  PATMISS = y6 (.5) y7 (.5) y8 (.5) y9(.5) y10 (.5);
  PATPROB = 1;
  SAVE = cfa4.sav;";

```

```

ANALYSIS =
  "TYPE = MIXTURE;
  ESTIMATOR = MLR;",

MODELPOPULATION =
  "%OVERALL%
  f1 BY y1-y5*.8;
  f2 BY y6-y10*.8;
  f1@1 f2@1;
  y1-y5*.36 y6-y10*9;
  f1 WITH f2*.95;
  [C#1@-2];
  %C#1%
  [f1@0 f2@15];
  f1@1 f2@5;
  %C#2%
  [f1@0 f2@0];
  f1@1 f2@1;";

MODEL =
  "%OVERALL%
  f1 BY y1-y5*.8;
  f2 BY y6-y10*4;
  f1@1 f2@1;
  y1-y5*.36 y6-y10*9;
  f1 WITH f2*.20;
  [y6-y10*1.42];" ,

OUTPUT = "PATTERNS TECH9;")

cfa_4_fit <- mplusModeler(cfa_4,
  dataout=here("mplus_files", "cfa4_demo_sim.dat"),
  modelout=here("mplus_files", "cfa4_demo_sim.inp"),
  check=TRUE, run = TRUE, hashfilename = FALSE)

```

view characteristics of simulated data

```

cfa4 <- read.delim(here("mplus_files", "cfa4.sav"), sep = "",
  header = FALSE,
  na.strings = "999.000000")

describe(cfa4)

```

Simulation 004

- explore biased outputs
- vary sample size to see changes in bias

testing & tuning

```
cfa_bias <- lapply(1:5, function(k) {
```

```

cfa_004 <- mplusObject(

TITLE = "CFA 1 - non-normal, no missing",

MONTECARLO =
  glue("NAMES ARE y1-y10;
    NOBSERVATIONS = {315-k*25}; ! vary sample size
    NREPS = 10000;
    SEED = 53487;
    CLASSES = C(1);
    GENCLASSES = C(2);
    SAVE = cfa004_{k}.sav;"),

ANALYSIS =
  "TYPE = MIXTURE;
    ESTIMATOR = MLR;",

MODELPOPULATION =
  "%OVERALL%
    f1 BY y1-y5*.8;
    f2 BY y6-y10*.8;
    f1@1 f2@1;
    y1-y5*.36 y6-y10*9;
    f1 WITH f2*.95;
    [C#1@-2];

    %C#1%          ! outlier class

    [f1@0 f2@15]; ! means (factor 2 set to 15 to tune skewness & kurtosis)
    f1@1 f2@5;    ! variances (factor 2 set to 5 to tune skewness & kurtosis)

    %C#2%          ! majority class

    [f1@0 f2@0];
    f1@1 f2@1;";

MODEL =
  "%OVERALL%
    f1 BY y1-y5*.8;
    f2 BY y6-y10*4;
    f1@1 f2@1;
    y1-y5*.36 y6-y10*9;
    f1 WITH f2*.20;

    [y6-y10*1.42];" ,

OUTPUT = " SAMPSTAT TECH9;")

cfa_004_fit <- mplusModeler(cfa_004,
  dataout=here("mplus_bias", "cfa004_demo_sim.dat"),
  modelout=sprintf(here("mplus_bias", "%d_cfa004_demo_sim.inp"), k),
  check=TRUE, run = TRUE, hashfilename = FALSE)
})

```

End of running CFA simulation
