# Moderated Non-Linear Factor Analysis

Kevin M. Donovan

August 27, 2020

## 1 Motivation

A fundamental step in scientific research entails the selection of specific concepts to investigate and the corresponding methods to measure differing levels of these concepts. Generally, these concepts are considered *latent* constructs, or unobservable directly in a population, with the measurements of these constructs being observable quantifications which are associated with these constructs. The quality of these measurements needs to be analyzed to ensure they are accurate surrogates of the unobservable latent constructs. Otherwise, conclusions from the observed data may not be reflective of the true characteristics of these latent constructs.

One component of the accuracy of the items used for measurement is *measurement invariance* (MI). Measurement invariance exists for an item when the item provides the same result/scale across different sample populations for a given level of the construct. In contrast, *differential item functioning* (DIF) exists for an item when two individuals of similar level for the construct do not have the same probability of obtaining a given value for the item. For example, a test measuring mathematical ability for two children of similar underlying skill results in different test scores between the children (or on average for two populations) due traits of the children (gender, race, etc.). Thus, the construct is not being accurately measured in the population due to these systematic item differences.

For a set of observed items, these latent constructs are often modeled using factor analysis. However, if DIF exists for some of these items, then the results from the factor analysis will not accurate reflect the relationship between the items and the construct due to this confounding. Not only does this DIF need to be detected, but it may be of interest to quantify the level of DIF within the items along with controlling for it when conducting inference on the underlying constructs (*factors*). Common methods to accomplish this are multiple groups analysis (MG), multiple-indicator multiple-cause (MIMIC) and moderated non-linear factor analysis (MNLFA). MNLFA is the most general of the three (the others are forms of MNLFA, MNLFA allows continuous and categorical items and covariates). The method is detailed in this document.

# 2 Methodology

## 2.1 Measurement Invariance

Following are adapted from [1]

Let $p$ by 1 $\boldsymbol{y_i}$ denote the item responses for subject $i$.
Let $r$ by 1 $\boldsymbol{\eta_i}$ denote the factor values for subject $i$
Let $q$ by 1 $\boldsymbol{x_i}$ denote the covariate vector for subject $i$

Let $f(\boldsymbol{y_i}|\boldsymbol{\eta_i}, \boldsymbol{x_i})$ denote the conditional density (assuming continuous items, "density" language used in general).

$$\text{MI} \leftrightarrow \boldsymbol{y_i}|\boldsymbol{\eta_i} \perp \boldsymbol{x_i} \leftrightarrow f(\boldsymbol{y_i}|\boldsymbol{\eta_i}, \boldsymbol{x_i}) = f(\boldsymbol{y_i}|\boldsymbol{\eta_i})$$

Weaker form of MI is *first-order MI* (FMI)
$$\text{FI} \leftrightarrow \text{E}(\boldsymbol{y_i}|\boldsymbol{\eta_i}) \perp \boldsymbol{x_i} \leftrightarrow \text{E}(\boldsymbol{y_i}|\boldsymbol{x_i}, \boldsymbol{\eta_i}) = \text{E}(\boldsymbol{y_i}|\boldsymbol{\eta_i})$$

## 2.2 Factor Analysis

We begin with a review of the standard factor analysis model for continuous items to motivate MNLFA.

Let $\boldsymbol{\delta_i}$ denote the residuals where

$y_i = \Lambda_y \eta_i + \delta_i$ s.t.
$E(\delta_i) = 0$, $Cov(\eta_{i,j}, \delta_{i,k}) = 0$ for $j = 1, \ldots, r$ and $k = 1, \ldots, p$ with
$Cov(\delta_i) = \Theta$
Independent factors, i.e. $\Theta = \text{diag}(\sigma^2{}_{11}, \ldots, \sigma^2{}_{rr})$ sometimes assumed.
The model is fit, i.e. estimates and standard errors of parameters are computed, using maximum likelihood assuming a multivariate normal distribution for the residuals.

## 2.3 MNLFA

Following are adapted from [1]

Let $E(y_i | \eta_i, x_i) = v_i + \Lambda_i \eta_i$
$Cov(y_i | \eta_i, x_i) = \Sigma_i$
$E(\eta_i | x_i) = \alpha_i$
$Cov(\eta_i | x_i) = \Psi_i$ s.t.

### 2.3.1 Modeling the means

$v_i = f_v(x_i)$ where $f_v(.)$ is some chosen functional form for the intercept term.

Linearity assumed for intercepts and factor means $\rightarrow$
$v_i = v_0 + K x_i$ and $\alpha_i = \alpha_0 + \Gamma_i x_i$

where $K$ and $\Gamma$ capture linear associations with $x_i$ for the factor intercepts and factor means respectively.

For the factor loadings, assuming linearity for the functional form
$\lambda_{a,i} = \lambda_{a,0} + \Omega_z x_i$ s.t.
$\lambda_{a,0}$ is $p$ by 1 vector of "baseline" loadings for factor $a$ when $x_i = 0$
and
$\Omega_a$ is $p$ by $q$ matrix of coefficients producing linear changes in loading $a$

### 2.3.2 Modeling the covariance

Restrictions on variance, correlation domains $\rightarrow$ functional form for covariance components more difficult (variance$>0$, $-1 <$correlation$< 1$)

**Variance**
One possibility:

$$\text{Var}(\eta_{i,aa}|\boldsymbol{x_i}) = \Psi_{i,aa} = \Psi_{i(aa)0} \exp(\boldsymbol{\beta'_{aa}} \boldsymbol{x_i}) \text{ s.t. } \Psi_{i(aa)0} > 0$$

$\Psi_{i(aa)0}$ captures baseline variance in factors when $\boldsymbol{x_i} = \boldsymbol{0}$.
$\boldsymbol{\beta_{aa}}$ is $q$ by 1 vector of moderation effects of $\boldsymbol{x_i}$ on factor $a$ variance

**Covariance**
One possibility:

$$\text{Z}[\text{Cor}(\eta_{a,i}, \eta_{b,i}|\boldsymbol{x_i})] = \xi_{i,ab} = \xi_{(ab),0} + \boldsymbol{u'_{ab}} \boldsymbol{x_i}$$

with Z[.] denoting Fisher's Z transformation. Correlation found using inverse Z transfom. Covariance calculated using variance and correlation formulas.

**NOTE:** Models for residual variance and covariance for $f(\boldsymbol{y_i}|\boldsymbol{\eta_i}, \boldsymbol{x_i})$ model can be specified similarly.

## 2.4   Fitting the model

Following are adapted from [1, 2] using depression-related example from [3]

**FITTING OF MODEL DONE USING MAXIMUM LIKELIHOOD (SEE PAPER)**

**Iterative Process** (various authors slightly differ on specifics):
**Step 1:** Determine factor structure

1. Select *calibration sample* = subsample of data randomly such that observations are independent (one per cluster for clustered data)

2. Run exploratory factor analysis (EFA) on calibration sample to determine factor structure (# of factors, which items to retain, etc.)

Example:
*cite paper* Consider survey of items assessing psychological wellness (level

of depression, stress, etc.)

EFA on calibration sample resulted in 2 factor solution being "optimal" based on visual inspection of eigenvalues/scree plot:

17 items $\rightarrow$ elevated loadings in factor 1 $\leftrightarrow$ "depression" markers
15 items $\rightarrow$ elevated loadings in factor 2 $\leftrightarrow$ "stress" markers

Based on this, decide which items to retain (dimension reduction step)
Possible criteria: keep items which "uniquely" load on one factor by comparing magnitude of loadings.
To derive less noisy estimation of "depression" factor, re-run EFA only with the 17 items retained.

**Step 2:** Test for DIF and develop conditional factor model

1. First, consider single-dimension models for each factor (if more then 1 factor was retained).

2. Single-dimensional model consists of the following

   (a) $y_i|\eta_i, x_i$ model: no DIF or covariates (simple FA model)

   (b) $\eta_i|x_i$ model: mean and variance chosen functions of covariates (ex. linear for mean, exponential for variance, no covariance due to single-dimension/factor)

   (c) Set scale of factor (usual step in FA) by setting baseline/intercept of conditional factor mean to 0, baseline/intercept of conditional factor variance to 1

   (d) Return estimates of model to be used in next step

3. Using results in above step, fit new model which considers DIF item-by-item. Process is:

   (a) Select a candidate set of items which a priori may be subject to DIF (could be all items)

   (b) For first item in set, say $y_{i,1}$ for simplicity, run MNLFA model with is the same as in the above step but with general DIF allowed for $y_{i,1}$ (done by specifying corresponding parameters/relationship

5

in software). For starting values in estimation procedure, use 0 (or more informed set) for parameters introduced per DIF and use estimates from previous step for parameters in both this and previous model. Could also simply fix parameters from previous model to their estimates.

(c) Conduct hypothesis test (likelihood ration test with previous model as null model) to assess degree of DIF. Could also do parameter specific selection for item to represent type of DIF using parameter-specific tests of significance (ex. Wald tests)

4. Repeat for all items in candidate set. May want to correct the DIF hypothesis tests for multiple comparisons. Choose final set of items which show evidence of DIF (informed by "significance" of tests).

**Step 3:** Run final model

1. Combining results from above models, have the "final" model

   (a) $\boldsymbol{\eta_i}|\boldsymbol{x_i}$ model: mean and variance chosen functions of covariates (ex. linear for mean, exponential for variance, no covariance due to single-dimension/factor)

   (b) Set scale of factor (usual step in FA) by setting baseline/intercept of conditional factor mean to 0, baseline/intercept of conditional factor variance to 1

   (c) $\boldsymbol{y_i}|\boldsymbol{\eta_i},\boldsymbol{x_i}$ model: MNLFA specification with DIF incorporated per results in previous model (by parameter for each item or for all parameters for each item)

**Step 4:** Predict factor scores, evaluate results

1. From final model estimates and SEs, calculate predicted factor scores (with measure of uncertainty if possible, prediction intervals for example)

2. Provide summaries of factor analysis results

   (a) Look at parameter estimates and SEs. These will provide information on degree and type of DIF for items, association between covariates and factor mean and variance, associations between factor scores and items, etc.
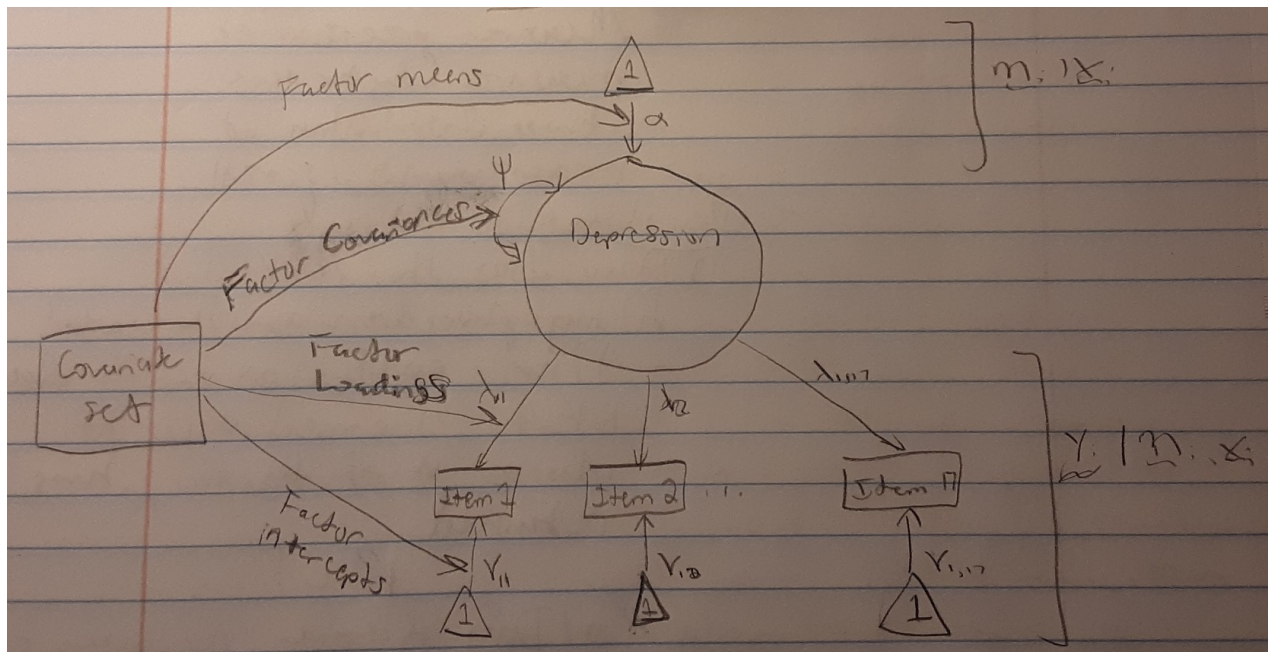
6

(b) Look at factor score distributions (visually, numerically, etc.) and how the distributions differ by covariates (age, gender, etc.)

3. Evaluate "accuracy/reliability" of factor scores. Possible methods:

   (a) Re-compute all steps in another calibration sample (or separate sample from same population). Compare factor scores within common factor between two samples (ex. using correlation). **May be ambiguous with multiple factor analysis due to lack of unique factor ordering. Also for a single sample, calibration samples likely to be highly correlated**.

   (b) Measure precision of factor scores using *Total Information Curve*.

## 2.5  MNLFA Path Diagram: Depression Example

Following adapted from [3]

Path diagram for depression example used in 2.4.

Item 1, ..., Item 17 are set of behavioral questions.



7

# References

[1] Daniel J Bauer. A more general model for testing measurement invariance and differential item functioning. *Psychological methods*, 22(3):507, 2017.

[2] Daniel J Bauer. Supplement to "A more general model for testing measurement invariance and differential item functioning". *Psychological methods*, 22(3):507, 2017.

[3] Patrick J Curran, James S McGinley, Daniel J Bauer, Andrea M Hussong, Alison Burns, Laurie Chassin, Kenneth Sher, and Robert Zucker. A moderated nonlinear factor model for the development of commensurate measures in integrative data analysis. *Multivariate behavioral research*, 49(3):214–231, 2014.