Rethinking Complexity in Item Response Theory Models

Wes Bonifay, Ph.D.







Excessive flexibility makes a model "so weak that there is no way to find evidence either for or against it." (Wexler, 1978, p. 346)

Complexity

Definition: The ability of a model to fit a wide range of data patterns (*Myung*, *Pitt*, & *Kim*, 2005)

• Affected by number of parameters & functional form (e.g., Collyer, 1985)

• y = x + b vs. $y = e^{xb}$

In latent variable modeling, complexity is typically gauged by counting parameters

- Model evaluation metrics like AIC, BIC, DIC penalize for complexity
- Minimum description length (*Rissanen, 1978*) also accounts for the functional form of the model

Fitting Propensity

A baseline for model fit can be established by fitting models to random data (*Cutting et al., 1992*)

Fitting propensity: A model's ability to fit diverse patterns of data, all else being equal

• Models with same number of parameters but different structures may exhibit different fitting propensities (*Preacher, 2006*)

Herein, we investigate the fitting propensities of 5 popular item response theory (IRT) models

Measurement Models

Hypothesis **1**: The EFA model will exhibit, on average, the highest fitting propensity

Hypothesis 2: The bifactor model will display higher fitting propensity than the DINA and DINO models

Unidimensional 3PL model?

Method

Data generation:

Sampling from a unit simplex (Smith & Tromble, 2004)

- Generate all possible item response patterns for a small number of binary response items
 - \rightarrow 2⁷ = 128 possible patterns
- Assign to each a random weight representing the number of simulees (of *N* = 10,000) who supplied that particular pattern

1,000 random data sets to represent the complete data space

Example data set:

0	0	0	0	0	0	0	20.348
0	0	0	0	0	0	1	60.367
0	0	0	0	0	1	0	27.676
0	0	0	0	0	1	1	7.983
0	0	0	0	1	0	0	3.714
0	0	0	0	1	0	1	10.517
0	0	0	0	1	1	0	47.091
					:		
1	1	1	1	0	0	1	4.961
1	1	1	1	0	1	0	30.693
1	1	1	1	0	1	1	14.663
1	1	1	1	1	0	0	1.993
1	1	1	1	1	0	1	2.673
1	1	1	1	1	1	0	67.551
1	1	1	1	1	1	1	40.919

Method

Evaluation measure: Y2/N statistic (Bartholomew & Leung, 2002; Cai et al., 2006):

Y2/N

Complete	Region	$Y2/N \le .01$
Data Space	EFA	1.4
	Bifactor	0.9
	DINA	0
	DINO	0
	Uni	0
	А	0.4
В	В	1.0
	С	0.5
	Unoccupied	98.1

Local dependence LD X² (Chen & Thissen, 1997)

Discussion

Hypotheses: confirmed

The importance of *functional form*:

- Arrangement of the variables in a model affects ability to fit well
- More complex models (EFA & bifactor) displayed high propensity to fit *any* data
- DINA & DINO models had low fitting propensity; theoretical difference → models fit different patterns
- Unidimensional model had an additional free parameter, but much lower fitting propensity!
- Strong implications re: model evaluation via goodness of fit

References

- Bartholomew, D. J., & Leung, S. O. (2002). A goodness of fit test for sparse 2^p contingency tables. *British Journal of Mathematical and Statistical Psychology*, 55:1-15.
- Cai, L., Maydeu-Olivares, A., Coffman, D., & Thissen, D. (2006). Limited-information goodness-of-fit testing of item response theory models for sparse 2^p tables. *British Journal of Mathematical and Statistical Psychology*, *59*, 173-194.
- Collyer, C. E. (1985). Comparing strong and weak models by fitting them to computer-generated data. *Perception & Psychophysics*, *38*(5), 476-481.
- Cutting, J. E., Bruno, N., Brady, N. P., & Moore, C. (1992). Selectivity, scope, and simplicity of models: A lesson from fitting judgments of perceived depth. *Journal of Experimental Psychology: General*, 121(3), 364-381.
- Myung, I. J., Pitt, M. A., & Kim, W. (2005). Model evaluation, testing and selection. In Lamberts, K. & Goldstone, R., (Eds.), *Handbook of Cognition*. London, UK: Sage Publications Ltd.
- Preacher, K. J. (2006). Quantifying parsimony in structural equation modeling. *Multivariate Behavioral Research*, 41(3), 227-259.

Rissanen, J. (1978). Modeling by shortest data description. Automatica, 14(5), 465-471.

Smith, N. A., & Tromble, R. W. (2004). Sampling uniformly from the unit simplex. *Johns Hopkins University, Tech. Rep.* 29.

Wexler, K. (1978). A review of John R. Anderson's *Language, Memory, and Thought. Cognition, 6,* 327-351.

Thanks to Dr. Li Cai for his assistance with this work.

Contact: bonifayw@missouri.edu

$$EFA = \begin{cases} 7 c \text{ (intercept) parameters} \\ 7 a \text{ (discrimination) parameters for Factor 1} \\ + 6 a \text{ (discrimination) parameters for Factor 2} \\ 20 \text{ free parameters} \\ \end{cases}$$

$$Bifactor = \begin{cases} 7 c \text{ (intercept) parameters} \\ 7 a \text{ (discrimination) parameters for the General Factor} \\ 5 a \text{ (discrimination) parameters for Specific Factor 1} \\ + 1 a \text{ (discrimination) parameters for Specific Factor 2} \\ 20 \text{ free parameters} \\ \end{cases}$$

$$\begin{cases} 7 \lambda_{I,0} \text{ (intercept) parameters} \\ 2 \lambda_{I,1,(2)} \text{ (main effect) parameters for Attribute 1} \\ 1 \lambda_{I,2,(2)} \text{ (main effect) parameters for Attribute 2} \\ 2 \lambda_{I,1,(2)} \text{ (main effect) parameters for Attribute 3} \\ 1 \lambda_{I,2,(2,3)} \text{ (interaction effect) parameters} \\ 2 \text{ (attribute intercept) parameters} \\ 2 \text{ (attribute intercept) parameters} \\ 2 \text{ (attribute discrimination) parameters} \\ 2 \text{ (attribute discrimination) parameters} \\ 2 \text{ (item discrimination) parameters} \\ 2 \text{ (item discrimination) parameters} \\ 2 \text{ (item discrimination) parameters} \\ 1 \text{ free parameters} \\ \end{cases}$$

