# Covariance structure and factor models

PSYC/EPSY/SOC/STAT 588

Spring 2011

Instructor: Sungjin Hong Office: Psych 429 (244-8296, hongsj AT illinois DOT edu) Class meetings: TR, 3:00-4:50, room 29 (219A/289 to access AMOS; 289 also reserved for Wed 3-5pm) Office hours: W, 2:00-3:00 or by appointment Course website: https://compass.illinois.edu

- 1. Path analysis
- 2. Factor analysis (measurement model in SEM terms)
- 3. General estimation procedure

## 1. Path analysis

 To parsimoniously describe relationships between a set of observed variables with directed paths by which we wish to represent causal influences of one variable to another



 $GPA = a IQ + E_G, \qquad AI = b GPA + E_A$ 

 $\operatorname{cov}(\operatorname{GPA},\operatorname{AI}) = a^2 b \operatorname{var}(\operatorname{IQ}) + b \operatorname{var}(\operatorname{E}_{\operatorname{G}})$ 

### 2. Factor analysis

- Scientifically useful variables are often not directly observable (measurable) --- e.g., IQ
- The FA model defines the unobservable (latent) variables such that they **linearly** relate to (influence or determine) a set of correlated, yet distinctive, observable variables (indicators)

$$X_{1} = a_{1} + b_{1} \operatorname{IQ} + E_{1}$$
  

$$\vdots \qquad \vdots$$
  

$$X_{q} = a_{q} + b_{q} \operatorname{IQ} + E_{q}$$

$$cf. \quad IQ = c_1 X_1 + c_2 X_2 + \dots + c_q X_q + E_{IQ}$$

### 3. General estimation procedure

- Relationships among observed and unobserved variables are quantitatively defined ---- "parameterized", "modeled"
- A number is derived to indicate discrepancy between the modeled structure and the data with some side conditions (assumptions) --- model fit

$$F = f$$
 (data, parameters, combining rule)  
 $F = tr [(\mathbf{S} - \mathbf{\Sigma}(\mathbf{\theta}))^2]$ : ULS

• Distributional characteristics of F (as a random variable) are established according to a set of parametric assumptions (e.g., multivariate normal data), which allows for inference on statistical reliability of some hypotheses, e.g., F = 0 or  $F_1 - F_2 = 0$ 

### SEM related myths

- Confirmatory vs. exploratory:
  - a priori model specified
  - statistical test available (parametric or nonparametric)
  - what software used (e.g., AMOS or SPSS)
- An SEM analysis proves something:

*"I got a dataset from somewhere and model fit from an AMOS output is acceptable. Now, I have proved my hypothesis (the model) and so I am entitled a Ph.D."* 

• Hopefully, we will unveil many more in this course

• Data matrix tall or wide?

The principal component model written for tall data matrix ---- "mathematical" notation



The principal component model written for a set of random variables (wide data matrix) --- "statistical" notation



- Variables in the data vector x are considered as random quantities; not interested in their realized values themselves but their distributions --- equivalent to the data matrix represented horizontally
- The "LISREL" notation follows this convention as is most common for SEM; and we'll do so in this course

• Given data matrix  $\mathbf{X}$  ( $N \times q$ ), deviation matrix (column mean centered) is written as:

 $\mathbf{Z} = \mathbf{L}\mathbf{X}, \quad \mathbf{L} = \mathbf{I} - \frac{1}{N}\mathbf{1}\mathbf{1'}$ : idempotent centering matrix

And unbiased sample covariance matrix:

$$\mathbf{S} = \frac{1}{N-1} \mathbf{Z}' \mathbf{Z}$$
 or  $\frac{1}{N-1} \mathbf{z} \mathbf{z}'$  with the horizontal notation

