

# Principal Components and Factor Analysis

PSY 5102: Advanced Statistics for Psychological and Behavioral Research 2

---

---

---

---

---

---

---

---

## Goals

- What is factor analysis?
- What are factors?
- Representing factors
  - Graphs and equations
- Extracting factors
  - Methods and criteria
- Interpreting factor structures
  - Factor rotation
- Reliability
  - Cronbach's alpha

---

---

---

---

---

---

---

---

## When and Why Do We Use Factor Analysis?

- Take many variables and explain them with a few “factors” or “components”
- To see whether different measures are tapping aspects of a common dimension

---

---

---

---

---

---

---

---

### General Steps in Factor Analysis

- ⦿ Step 1: Select and measure a set of variables in a given domain
- ⦿ Step 2: Screen data in order to prepare the correlation matrix
- ⦿ Step 3: Factor extraction
- ⦿ Step 4: Factor rotation to increase interpretability
- ⦿ Step 5: Interpretation of factors
- ⦿ Further Steps: Validate and determine reliability of the scales

---

---

---

---

---

---

---

---

### Correlation Matrix

	Talk 1	Social Skills	Interest	Talk 2	Selfish	Liar
Talk 1	1.000					
Social Skills	.772	1.000				
Interest	.646	.879	1.000			
Talk 2	.074	-.120	.054	1.000		
Selfish	-.131	.031	-.101	.441	1.000	
Liar	.068	.012	.110	.261	.277	1.000

- ⦿ In factor analysis, we look to reduce the correlation matrix into smaller sets of *uncorrelated* dimensions

---

---

---

---

---

---

---

---

### What is a Factor?

- ⦿ If several variables correlate highly, they might measure aspects of a common underlying dimension
  - These dimensions are called factors
- ⦿ Factors are classification axes along which the measures can be plotted
  - The greater the loading of variables on a factor, the more that factor explains relationships between those variables

---

---

---

---

---

---

---

---

### What is a "Good" Factor?

- A good factor should...
  - Make sense
  - Be easy to interpret
  - Have a simple structure
  - Lack complex loadings

---

---

---

---

---

---

---

---

### Problems with Factor Analysis

- Unlike many of the other analyses we have covered, there is no statistical criterion to serve as a comparison for the linear combination
- It is more art than science
  - There are a number of extraction methods (e.g., principal-components analysis)
  - There are a number of rotation methods (e.g., Orthogonal, Oblique)
  - Choice of the number of factors to extract
  - Communality estimates

---

---

---

---

---

---

---

---

### Types of Factor Analysis

- Exploratory factor analysis
  - Summarizing data by grouping correlated variables
  - Investigating sets of measured variables related to theoretical constructs
  - Usually done near the onset of research
  - This is the type of factor analysis that we will address

---

---

---

---

---

---

---

---

## Types of Factor Analysis

- ◎ **Confirmatory Factor Analysis**
  - More advanced technique
  - When factor structure is known...or at least theorized
  - This basically involves testing the generalization of a factor structure to new data
  - This is tested through Structural Equation Modeling which we will discuss later

---

---

---

---

---

---

---

---

## Basic Terminology

- ◎ **Orthogonal Rotation**
  - Loading Matrix: correlation between each variable and the factor
- ◎ **Oblique Rotation**
- ◎ **Factor Correlation Matrix: correlations between the factors**
  - Structure Matrix: correlation between factors and variables
  - Pattern Matrix: unique relationship between each factor and variable uncontaminated by overlap between the factors
- ◎ **Factor Coefficient matrix: coefficients used to calculate factor scores (like regression coefficients)**

---

---

---

---

---

---

---

---

## Factor Analysis vs. Principal-Components Analysis

- ◎ **FA produces factors**
  - PCA produces components
- ◎ **Factors cause variables**
  - Components are aggregates of the variables
- ◎ **FA analyzes only the variance shared among the variables (common variance without error or unique variance)**
  - PCA analyzes all of the variance
- ◎ **FA: "What are the underlying processes that could produce these correlations?"**
  - PCA: Just summarize empirical associations, very data driven

---

---

---

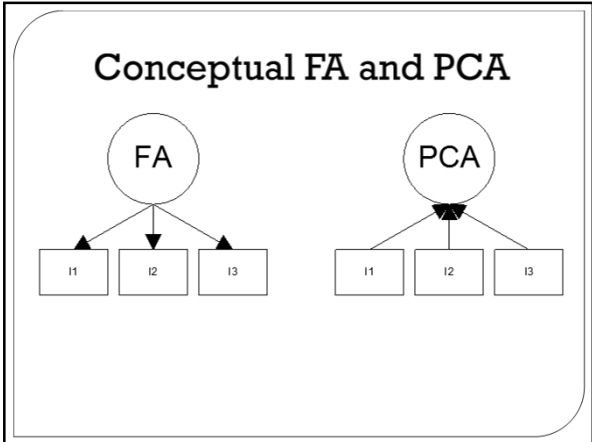
---

---

---

---

---



---

---

---

---

---

---

---

---

### Principal-Components Analysis

- ◉ The goal of principal-components analysis is to identify a new set of a few variables that explain all (or nearly all) of the total variance
  - The goal is parsimony
- ◉ These principal-components are a linear function of the original variables
  - The first principal-component maximizes the amount of variance that is explained
- ◉ Eigenvector: the linear function
  - Their goal is to explain as much variability as possible
  - The number varies between analyses
- ◉ Eigenvalue: the total amount of variance that is explained by an eigenvector

---

---

---

---

---

---

---

---

### How Many Factors to Extract?

- ◉ There are a number of strategies for deciding when to stop extracting factors or components
  - Percentage of variance that is explained
  - A specific number of factors may be extracted
  - Kaiser's stopping rule: only extracts (and retains) those factors with eigenvalues of at least 1
  - Scree test is a graphical procedure that depicts the eigenvalues

---

---

---

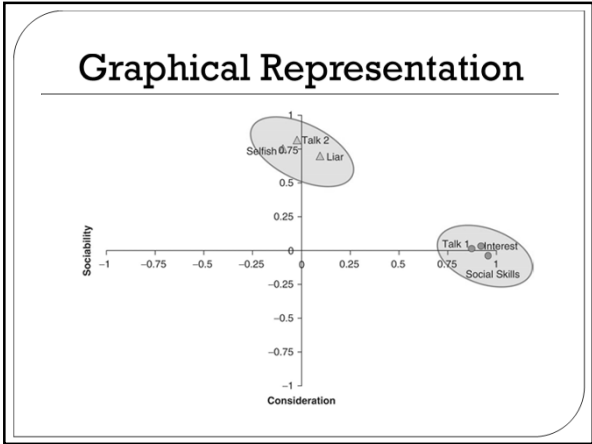
---

---

---

---

---




---

---

---

---

---

---

---

---

### Mathematical Representation

$$Y = b_1X_1 + b_2X_2 \dots b_nX_n$$

$$Factor_i = b_1Variable_1 + b_2Variable_2 \dots b_nVariable_n$$

$$Y = b_1X_1 + b_2X_2 \dots b_nX_n$$

Sociability =  $b_1$ Talk1 +  $b_2$ Social Skills +  $b_3$ Interest  
 +  $b_4$ Talk2 +  $b_5$ Selfish +  $b_6$ Liar

Consideration =  $b_1$ Talk1 +  $b_2$ Social Skills +  $b_3$ Interest  
 +  $b_4$ Talk2 +  $b_5$ Selfish +  $b_6$ Liar

---

---

---

---

---

---

---

---

### Factor Loadings

- ⦿ The  $b$  values in the equation represent the weights of a variable on a factor
- ⦿ These values are the same as the coordinates on a factor plot
- ⦿ They are called Factor Loadings
- ⦿ These values are stored in a *Factor pattern matrix* ( $A$ )

$$A = \begin{pmatrix} 0.87 & 0.01 \\ 0.96 & -0.03 \\ 0.92 & 0.04 \\ 0.00 & 0.82 \\ -0.10 & 0.75 \\ 0.09 & 0.70 \end{pmatrix}$$

---

---

---

---

---

---

---

---

## Practical Concerns

- How many participants?
  - Subjects-to-variables ratio should be at least 5-to-1 (5 participants for every variable)

---

---

---

---

---

---

---

---

### The SPSS Anxiety Questionnaire

23 items that are intended to capture aspects of anxiety concerning statistical analyses

	SD	D	N	A	SA
1					
2					
3					
4					
5					
6					
7					
8					
9					
10					
11					
12					
13					
14					
15					
16					
17					
18					
19					
20					
21					
22					
23					

---

---

---

---

---

---

---

---

## Initial Considerations

- The quality of analysis depends upon the quality of the data (“Garbage In” → “Garbage Out”)
- Test variables should correlate quite well
  - $r > .3$
- Avoid Multicollinearity:
  - several variables highly correlated,  $r > .80$
- Avoid Singularity:
  - some variables perfectly correlated,  $r = 1$
- Screen the correlation matrix and eliminate any variables that obviously cause concern

---

---

---

---

---

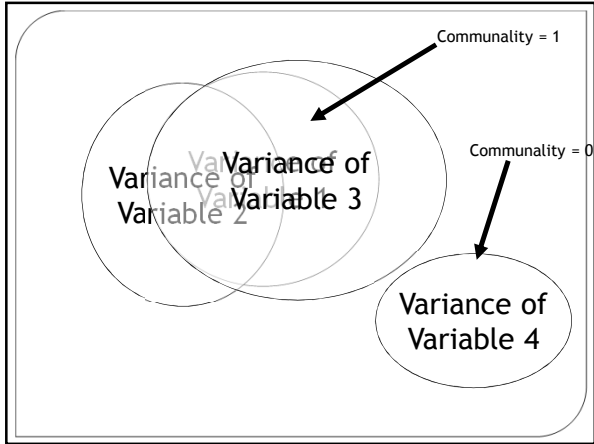
---

---

---








---

---

---

---

---

---

---

---

### Finding Factors

- ◉ We find factors by calculating the amount of common variance
  - Circularity
- ◉ Principal Components Analysis:
  - Assume all variance is shared
  - All communalities = 1
- ◉ Factor Analysis
  - Estimate communality
  - Use squared multiple correlation (SMC)
- ◉ Principal Components and Factor Analysis will identify similar factors when there are a large number of variables (i.e., more than 30) and the communalities are high (i.e., greater than .7)

---

---

---

---

---

---

---

---

Communalities		
	Initial	Extraction
Q01	1.000	.435
Q02	1.000	.414
Q03	1.000	.530
Q04	1.000	.469
Q05	1.000	.343
Q06	1.000	.654
Q07	1.000	.545
Q08	1.000	.739
Q09	1.000	.484
Q10	1.000	.335
Q11	1.000	.690
Q12	1.000	.513
Q13	1.000	.536
Q14	1.000	.488
Q15	1.000	.378
Q16	1.000	.487
Q17	1.000	.683
Q18	1.000	.597
Q19	1.000	.343
Q20	1.000	.484
Q21	1.000	.550
Q22	1.000	.464
Q23	1.000	.412

Extraction Method: Principal Component

---

---

---

---

---

---

---

---

## Factor Extraction

- **Kaiser's Extraction**
  - Kaiser (1960): retain factors with Eigenvalues > 1
- **Scree Plot**
  - Cattell (1966): use 'point of inflexion' of the scree plot
- **Which Rule?**
  - Use Kaiser's Extraction when
    - less than 30 variables, communalities after extraction > 0.7
    - sample size > 250 and mean communality ≥ 0.6
  - Scree plot is good if sample size is > 200

---

---

---

---

---

---

---

---

---

---

---

---

Total Variance Explained

Component	Initial Eigenvalues			Extraction Sums of Squared Loadings			Rotation Sums of Squared Loadings		
	Total	% of Variance	Cumulative %	Total	% of Variance	Cumulative %	Total	% of Variance	Cumulative %
1	7.290	31.696	31.696	7.290	31.696	31.696	3.730	16.219	16.219
2	1.739	7.560	39.256	1.739	7.560	39.256	3.340	14.523	30.742
3	1.317	5.725	44.981	1.317	5.725	44.981	2.553	11.099	41.842
4	1.227	5.330	50.317	1.227	5.330	50.317	1.949	8.475	50.317
5	.988	4.295	54.612						
6	.895	3.893	58.504						
7	.806	3.502	62.007						
8	.753	3.404	65.410						
9	.751	3.265	68.676						
10	.717	3.117	71.793						
11	.684	2.972	74.765						
12	.670	2.911	77.676						
13	.612	2.661	80.337						
14	.578	2.512	82.849						
15	.549	2.388	85.236						
16	.523	2.275	87.511						
17	.508	2.210	89.721						
18	.456	1.982	91.704						
19	.424	1.843	93.546						
20	.408	1.773	95.319						
21	.379	1.650	96.969						
22	.364	1.583	98.552						
23	.333	1.448	100.000						

Extraction Method: Principal Component Analysis.

---

---

---

---

---

---

---

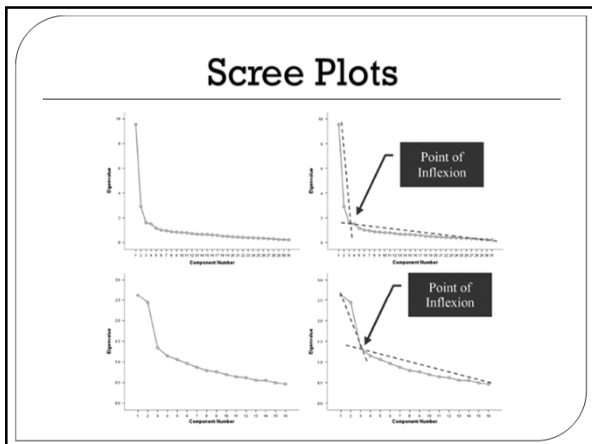
---

---

---

---

---




---

---

---

---

---

---

---

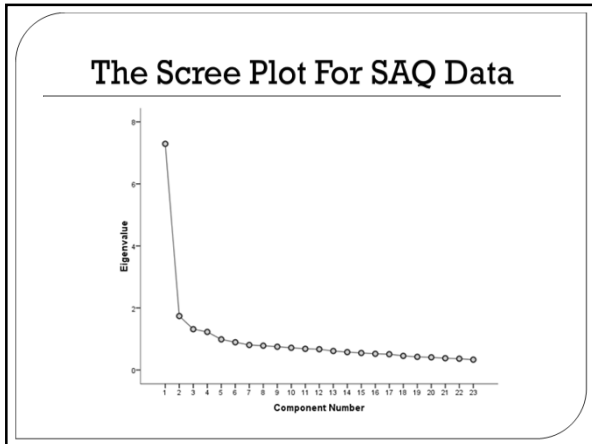
---

---

---

---

---



---

---

---

---

---

---

---

---

### Rotation

- To aid interpretation it is possible to maximize the loading of a variable on one factor while minimizing its loading on all other factors
- This is known as Factor Rotation
- There are two types
  - Orthogonal (factors are uncorrelated)
  - Oblique (factors intercorrelate)

---

---

---

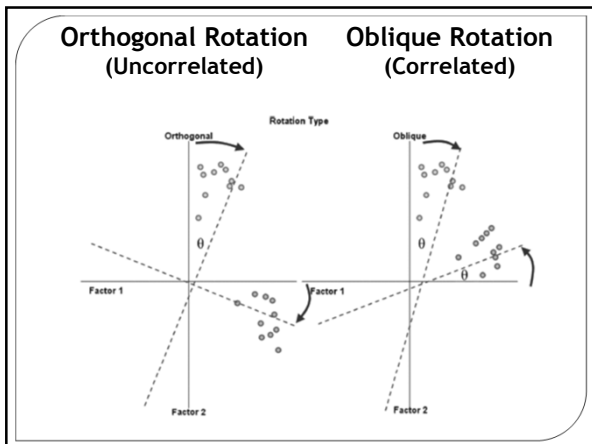
---

---

---

---

---



---

---

---

---

---

---

---

---

### Before Rotation

Component Matrix

	Component			
	1	2	3	4
Q18	.701			
Q07	.656			
Q16	.679			
Q13	.673			
Q12	.659			
Q21	.658			
Q14	.656			
Q11	.662			-.400
Q17	.643			
Q04	.634			
Q03	-.029			
Q15	.593			
Q01	.596			
Q05	.558			
Q08	.549	.401		-.417
Q10	.437			
Q20	.408			-.404
Q19	-.427			
Q09	.327			
Q02	.548			
Q22	-.455			
Q06	.562		.571	
Q23				.507

Extraction Method: Principal Component Analysis.  
a. 4 components extracted.

---

---

---

---

---

---

---

---

---

---

---

---

### Orthogonal Rotation

Rotated Component Matrix\*

	Component			
	1	2	3	4
I have little experience of computers	.800			
SPSS always crashes when I try to use it	.654			
I worry that I will cause irreparable damage because of my incompetence with computers	.647			
All computers hate me	.630			
Computers have minds of their own and deliberately go wrong whenever I use them	.576			
Computers are useful only for playing games	.550			
Computers are out to get me	.489			
I can't sleep for thoughts of eigen vectors		.677		
I sleep up under my duvet thinking that I am trapped under a normal distribution		.661		
Standard deviations excite me		-.587		
People try to tell you that SPSS makes statistics easier to understand but I object	.473	.523		
I dream that Pearson is attacking me with correlation coefficients		.516		
I sleep openly at the mention of central tendency		.514		
Statistics makes me cry		.486		
I don't understand statistics		.429		
I have never been good at mathematics			.633	
I slip into a coma whenever I see an equation			.747	
I'd rather do statistics at school			.747	.646
My friends are better at statistics than me				.646
My friends are better at SPSS than I am				.588
If I'm good at statistics my friends will think I'm a nerd				.543
My friends will think I'm stupid for not being able to cope with SPSS				.443
Everybody looks at me when I use SPSS				.427

Extraction Method: Principal Component Analysis.  
Rotation Method: Varimax with Kaiser Normalization.  
\*. Rotation converged in 9 iterations.

---

---

---

---

---

---

---

---

---

---

---

---

### Oblique Rotation

Pattern Matrix\*

	Component			
	1	2	3	4
I can't sleep for thoughts of eigen vectors	.708			
I sleep up under my duvet thinking that I am trapped under a normal distribution	.691			
Standard deviations excite me	-.511			
I dream that Pearson is attacking me with correlation coefficients	.403			
I sleep openly at the mention of central tendency	.400			
Statistics makes me cry				
I don't understand statistics				
My friends are better at SPSS than I am		.643		
My friends are better at statistics than me		.621		
If I'm good at statistics my friends will think I'm a nerd		.618		
My friends will think I'm stupid for not being able to cope with SPSS		.567		
Everybody looks at me when I use SPSS			.688	
I have little experience of computers			.713	
SPSS always crashes when I try to use it			.683	
All computers hate me			.660	
I worry that I will cause irreparable damage because of my incompetence with computers			.630	
Computers have minds of their own and deliberately go wrong whenever I use them			.588	
Computers are useful only for playing games			.550	
Computers are out to get me			.489	
People try to tell you that SPSS makes statistics easier to understand but I object	.412	.462		
Computers are out to get me			.411	
I have never been good at mathematics				-.902
I slip into a coma whenever I see an equation				-.774
I'd rather do statistics at school				-.774

Extraction Method: Principal Component Analysis.  
Rotation Method: Oblimin with Kaiser Normalization.  
\*. Rotation converged in 23 iterations.

---

---

---

---

---

---

---

---

---

---

---

---

### What Do the Factors Represent?

- ◉ We assume that algebraic factors represent psychological constructs
  - Factor 1: Fear of statistics (e.g., "I can't sleep for thoughts of eigenvectors")
  - Factor 2: Fear of peer evaluation (e.g., "My friends are better at statistics than me")
  - Factor 3: Fear of computers(e.g., "All computers hate me")
  - Factor 4: Fear of mathematics (e.g., "I have never been good at mathematics")
- ◉ The nature of these psychological dimensions is 'guessed at' by looking at the loadings for a factor
- ◉ There is no way to "know" what the factor represents...rather we have to decide ourselves
  - The same set of items may be referred to as "social dominance" by one researcher but "aggression" by another

---

---

---

---

---

---

---

---

### Reliability

- ◉ Test-Retest Method
  - Complete the same measure on two occasions and calculates the correlation
- ◉ Alternate Form Method
  - Complete two slightly different forms of the same measure and calculates the correlation
- ◉ Split-Half Method
  - Splits the questionnaire into two random halves and calculates the correlation
- ◉ Cronbach's Alpha
  - Basically splits the questionnaire into all possible halves, calculates the scores, correlates them, and averages the correlation for all splits
  - Ranges from 0 (no reliability) to 1 (complete reliability)

---

---

---

---

---

---

---

---

### Interpreting Cronbach's Alpha

- ◉Kline (1999)
  - Reliable if  $\alpha > .7$
- ◉Depends on the number of items
  - More questions = bigger  $\alpha$
- ◉Treat Subscales separately
- ◉Remember to reverse score reverse phrased items!
  - If not,  $\alpha$  is reduced and can even be negative

---

---

---

---

---

---

---

---



## Reliability for Fear of Mathematics Subscale

**Item-Total Statistics**

	Scale Mean if Item Deleted	Scale Variance if Item Deleted	Corrected Item-Total Correlation	Squared Multiple Correlation	Cronbach's Alpha if Item Deleted
I have never been good at mathematics	4.72	2.470	.684	.470	.740
I did badly at mathematics at school	4.70	2.453	.682	.467	.742
I slip into a coma whenever I see an equation	4.49	2.504	.652	.425	.772

**Reliability Statistics**

Cronbach's Alpha	Cronbach's Alpha Based on Standardized Items	N of Items
.819	.819	3

---



---



---



---



---



---



---



---