

Using Principal Components Analysis and Exploratory Factor Analysis for Refining Survey Modules

Joshua G. Mausolf
jmausolf@uchicago.edu



Advanced Methods in Survey Research
May 25th, 2016

NORC

at the UNIVERSITY *of* CHICAGO

You have designed a survey module with multiple questions hoping to identify a construct, such as “Interview Quality,” “Gentrification,” or “Neighborhood Resilience.” **Do these questions work and which ones should we keep?**

- **Solution, Part 1:** Use correlation, Cronbach’s Alpha, PCA, and EFA to select the best questions.
- **Solution, Part 2:** Test your new components or factors for construct validity.

Aims of this presentation

1. Brief Overview of Primary Methods
2. Outline of Steps to Refine Your Module
3. Fully Worked Factor Analysis Example in Stata
4. Example Test of Our Construct's Validity

Overview of Primary Methods

1. Correlation
2. Cronbach's Alpha
3. Principal Components Analysis (PCA)
4. Elementary Factor Analysis (EFA)

1. Correlation
2. Cronbach's Alpha
3. Principal Components Analysis (PCA)
4. Elementary Factor Analysis (EFA)

A method to look for statistical associations between variables, which in our case are survey questions

Overview of Primary Methods

1. Correlation
2. Cronbach's Alpha
3. Principal Components Analysis (PCA)
4. Elementary Factor Analysis (EFA)

A measure of internal consistency [0, 1]. It indicates how closely related a set of items, such as survey questions, are as a group. Typically, an $\alpha \geq 0.7$ is acceptable. For this exercise, it may be less.

1. Correlation
2. Cronbach's Alpha
3. Principal Components Analysis (PCA)
4. Elementary Factor Analysis (EFA)

A dimensionality reduction technique, which attempts to reduce a large number of variables into a smaller number of variables. A component is a unique combination of variables. An eigenvalue > 1 is significant.

1. Correlation
2. Cronbach's Alpha
3. Principal Components Analysis (PCA)
4. Elementary Factor Analysis (EFA)

An alternate dimensionality reduction technique. A factor is a unique combination of variables. An eigenvalue > 1 is significant.

Outline of Analysis Steps

High Level Overview

1. Preliminary Steps: Data Cleaning
2. First Steps: Analyze Entire Module
3. Next Steps: Determine Factors and Reanalyze
4. Test Final Factor(s) for Construct Validity

Outline of Analysis Steps

Detailed Overview

1. Preliminary Steps: Data Cleaning
 1. Clean Key Punch Data
 2. Reverse Code Relevant Questions in Your Survey Module
2. First Steps: Analyze Entire Module
3. Next Steps: Determine Factors and Reanalyze
4. Test Final Factor(s) for Construct Validity

Outline of Analysis Steps

Detailed Overview

1. Preliminary Steps: Data Cleaning
2. First Steps: Analyze Entire Module
 1. Summarize the Data
 2. Check Variable Correlations
 3. Cronbach's Alpha - First Pass
 4. PCA - First Pass
 5. Factor Analysis (EFA) - First Pass
3. Next Steps: Determine Factors and Reanalyze
4. Test Final Factor(s) for Construct Validity

Outline of Analysis Steps

Detailed Overview

1. Preliminary Steps: Data Cleaning
2. First Steps: Analyze Entire Module
3. Next Steps: Determine Factors and Reanalyze

After examining the results of your first pass of Cronbach's Alpha, PCA, and EFA

- Determine which questions relate as principal components and factors
- Rerun Cronbach's Alpha, PCA, and EFA on each new factor, check results
- "Rotate" factors and save results to analyze for construct validity

4. Test Final Factor(s) for Construct Validity

Outline of Analysis Steps

Detailed Overview

1. Preliminary Steps: Data Cleaning
2. First Steps: Analyze Entire Module
3. Next Steps: Determine Factors and Reanalyze
4. Test Final Factor(s) for Construct Validity

Using your factor as an independent variable,

1. use standard statistical techniques to look for a statistically significant relationship with one or more theoretically relevant dependent variables
2. Examples: F-test, ANOVA, Linear or Logistic Regression

Sample PCA output

```
pca pre_ia p1-p2 p3rev ia_1-ia_3 ia4rev ia_5-ia_8, comp(2)
```

Component	Eigenvalue	Difference	Proportion	Cumulative
Comp1	2.52059	.720161	0.2100	0.2100
Comp2	1.80043	.286807	0.1500	0.3601
...				
Comp12	.195379	.	0.0163	1.0000

Principal components (eigenvectors)

Variable	Comp1	Comp2	Unexplained
pre_ia	-0.0465	0.3995	.7073
p1	0.2529	0.4061	.5419
p2	-0.0484	0.3777	.7372
p3rev	0.1576	-0.0412	.9344
ia_1	0.3401	0.0053	.7084
ia_2	0.5452	0.0573	.2449
ia_3	0.5111	0.2059	.2652
ia4rev	0.2209	-0.2228	.7877
ia_5	0.3247	-0.2656	.6072
ia_6	-0.2529	0.0953	.8224
ia_7	0.0209	0.4419	.6473
ia_8	-0.1179	0.4012	.6752

These measures come from interviewer ratings of respondents in a nationally representative survey (NSHAP). Note that the output has been truncated for display purposes.

Sample PCA output

```
pca pre_ia p1-p2 p3rev ia_1-ia_3 ia4rev ia_5-ia_8, comp(2)
```

Component	Eigenvalue	Difference	Proportion	Cumulative
Comp1	2.52059	.720161	0.2100	0.2100
Comp2	1.80043	.286807	0.1500	0.3601
...				
Comp12	.195379	.	0.0163	1.0000

Principal components (eigenvectors)

Variable	Comp1	Comp2	Unexplained
pre_ia	-0.0465	0.3995	.7073
p1	0.2529	0.4061	.5419
p2	-0.0484	0.3777	.7372
p3rev	0.1576	-0.0412	.9344
ia_1	0.3401	0.0053	.7084
ia_2	0.5452	0.0573	.2449
ia_3	0.5111	0.2059	.2652
ia4rev	0.2209	-0.2228	.7877
ia_5	0.3247	-0.2656	.6072
ia_6	-0.2529	0.0953	.8224
ia_7	0.0209	0.4419	.6473
ia_8	-0.1179	0.4012	.6752

EIGENVALUES: The variance of the component. They add to the sum of the variance in the variables.

Critical eigenvalue: > 1

Sample PCA output

```
pca pre_ia p1-p2 p3rev ia_1-ia_3 ia4rev ia_5-ia_8, comp(2)
```

Component	Eigenvalue	Difference	Proportion	Cumulative
Comp1	2.52059	.720161	0.2100	0.2100
Comp2	1.80043	.286807	0.1500	0.3601
...				
Comp12	.195379	.	0.0163	1.0000

Principal components (eigenvectors)

Variable	Comp1	Comp2	Unexplained
pre_ia	-0.0465	0.3995	.7073
p1	0.2529	0.4061	.5419
p2	-0.0484	0.3777	.7372
p3rev	0.1576	-0.0412	.9344
ia_1	0.3401	0.0053	.7084
ia_2	0.5452	0.0573	.2449
ia_3	0.5111	0.2059	.2652
ia4rev	0.2209	-0.2228	.7877
ia_5	0.3247	-0.2656	.6072
ia_6	-0.2529	0.0953	.8224
ia_7	0.0209	0.4419	.6473
ia_8	-0.1179	0.4012	.6752

The difference between the size of this component's eigenvalue and the *next* component's eigenvalue.

Sample PCA output

```
pca pre_ia p1-p2 p3rev ia_1-ia_3 ia4rev ia_5-ia_8, comp(2)
```

Component	Eigenvalue	Difference	Proportion	Cumulative
Comp1	2.52059	.720161	0.2100	0.2100
Comp2	1.80043	.286807	0.1500	0.3601
...				
Comp12	.195379	.	0.0163	1.0000

Principal components (eigenvectors)

Variable	Comp1	Comp2	Unexplained
pre_ia	-0.0465	0.3995	.7073
p1	0.2529	0.4061	.5419
p2	-0.0484	0.3777	.7372
p3rev	0.1576	-0.0412	.9344
ia_1	0.3401	0.0053	.7084
ia_2	0.5452	0.0573	.2449
ia_3	0.5111	0.2059	.2652
ia4rev	0.2209	-0.2228	.7877
ia_5	0.3247	-0.2656	.6072
ia_6	-0.2529	0.0953	.8224
ia_7	0.0209	0.4419	.6473
ia_8	-0.1179	0.4012	.6752



Proportion of variance explained by each component.

Sample PCA output

```
pca pre_ia p1-p2 p3rev ia_1-ia_3 ia4rev ia_5-ia_8, comp(2)
```

Component	Eigenvalue	Difference	Proportion	Cumulative
Comp1	2.52059	.720161	0.2100	0.2100
Comp2	1.80043	.286807	0.1500	0.3601
...				
Comp12	.195379	.	0.0163	1.0000

Principal components (eigenvectors)

Variable	Comp1	Comp2	Unexplained
pre_ia	-0.0465	0.3995	.7073
p1	0.2529	0.4061	.5419
p2	-0.0484	0.3777	.7372
p3rev	0.1576	-0.0412	.9344
ia_1	0.3401	0.0053	.7084
ia_2	0.5452	0.0573	.2449
ia_3	0.5111	0.2059	.2652
ia4rev	0.2209	-0.2228	.7877
ia_5	0.3247	-0.2656	.6072
ia_6	-0.2529	0.0953	.8224
ia_7	0.0209	0.4419	.6473
ia_8	-0.1179	0.4012	.6752

LOADINGS. Bigger = more associated with this component. Substantively: the correlation between the component and the variable.

Sample PCA output

```
pca pre_ia p1-p2 p3rev ia_1-ia_3 ia4rev ia_5-ia_8, comp(2)
```

Component	Eigenvalue	Difference	Proportion	Cumulative
Comp1	2.52059	.720161	0.2100	0.2100
Comp2	1.80043	.286807	0.1500	0.3601
...				
Comp12	.195379	.	0.0163	1.0000

Principal components (eigenvectors)

Variable	Comp1	Comp2	Unexplained
pre_ia	-0.0465	0.3995	.7073
p1	0.2529	0.4061	.5419
p2	-0.0484	0.3777	.7372
p3rev	0.1576	-0.0412	.9344
ia_1	0.3401	0.0053	.7084
ia_2	0.5452	0.0573	.2449
ia_3	0.5111	0.2059	.2652
ia4rev	0.2209	-0.2228	.7877
ia_5	0.3247	-0.2656	.6072
ia_6	-0.2529	0.0953	.8224
ia_7	0.0209	0.4419	.6473
ia_8	-0.1179	0.4012	.6752

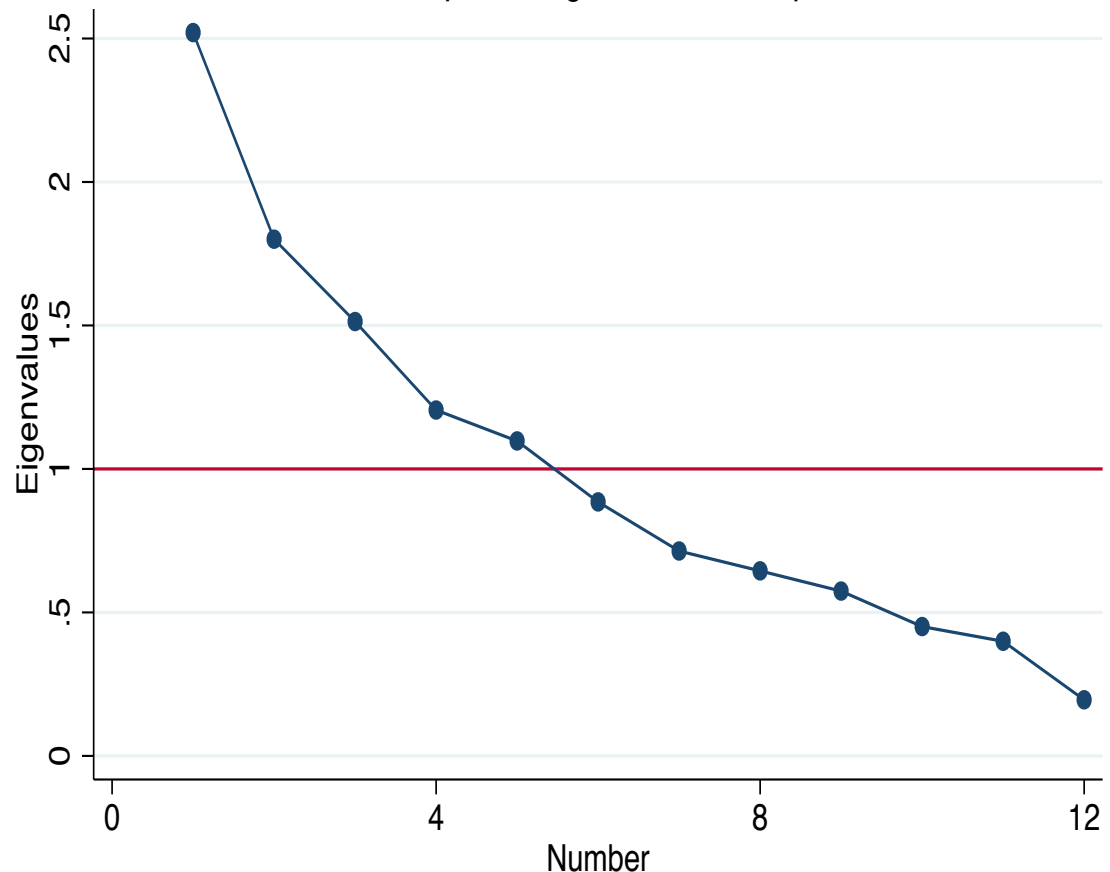
← **Proportion of variance in this measure that is not explained by the displayed components. These would all be '0' if I did not restrict the display to 2 components.**

How many components do I retain?

- Since the aim is data reduction, we need some criteria or we will have the same number of variables after a PCA.
- One rule of thumb is that a component should not be **retained unless it has an eigenvalue greater than or equal to one** (the 'Kaiser' criterion).

Principal Components Analysis Analysis

Scree plot of eigenvalues after pca



Kinds of validity that PCA can assess

- **Convergent:** If our theory predicts that some set of measures should be associated with one another, we should see that they load on the same component.
- **Divergent:** Conversely, if we think two measures really measure different things, they should *not* load on the same component.

Sample EFA output

```
factor pre_ia p1-p2 p3rev ia_1-ia_3 ia4rev ia_5-ia_8, fa(2)
```

Factor	Eigenvalue	Difference	Proportion	Cumulative
Factor1	2.01020	0.90722	0.5268	0.5268
Factor2	1.10298	0.30999	0.2890	0.8158
...				

```
LR test: independent vs. saturated: chi2(66) = 179.46 Prob>chi2 = 0.0000
```

```
Factor loadings (pattern matrix) and unique variances
```

Variable	Factor1	Factor2	Uniqueness
pre_ia	-0.0386	0.4460	0.7996
p1	0.3446	0.4067	0.7159
p2	-0.0519	0.3627	0.8658
p3rev	0.1832	-0.0364	0.9651
ia_1	0.4135	-0.0300	0.8281
ia_2	0.8336	0.0187	0.3048
ia_3	0.7864	0.2402	0.3239
ia4rev	0.2639	-0.2609	0.8623
ia_5	0.4265	-0.3483	0.6968
ia_6	-0.3073	0.1049	0.8946
ia_7	0.0455	0.4274	0.8152
ia_8	-0.1465	0.4047	0.8148

**Same measures as the PCA,
but now in an EFA context.**

Sample EFA output

```
factor pre_ia p1-p2 p3rev ia_1-ia_3 ia4rev ia_5-ia_8, fa(2)
```

Factor	Eigenvalue	Difference	Proportion	Cumulative
Factor1	2.01020	0.90722	0.5268	0.5268
Factor2	1.10298	0.30999	0.2890	0.8158
...				

```
LR test: independent vs. saturated: chi2(66) = 179.46 Prob>chi2 = 0.0000
```

```
Factor loadings (pattern matrix) and unique variances
```

Variable	Factor1	Factor2	Uniqueness
pre_ia	-0.0386	0.4460	0.7996
p1	0.3446	0.4067	0.7159
p2	-0.0519	0.3627	0.8958
p3rev	0.1832	-0.0364	0.9651
ia_1	0.4135	-0.0300	0.8281
ia_2	0.8336	0.0187	0.3048
ia_3	0.7864	0.2402	0.3239
ia4rev	0.2639	-0.2609	0.8623
ia_5	0.4265	-0.3483	0.6968
ia_6	-0.3073	0.1049	0.8946
ia_7	0.0455	0.4274	0.8152
ia_8	-0.1465	0.4047	0.8148

FACTOR LOADINGS:
Correlations between the
measure and the factor, as in
PCA.

Sample EFA output

```
factor pre_ia p1-p2 p3rev ia_1-ia_3 ia4rev ia_5-ia_8, fa(2)
```

Factor	Eigenvalue	Difference	Proportion	Cumulative
Factor1	2.01020	0.90722	0.5268	0.5268
Factor2	1.10298	0.30999	0.2890	0.8158
...				

```
LR test: independent vs. saturated: chi2(66) = 179.46 Prob>chi2 = 0.0000
```

```
Factor loadings (pattern matrix) and unique variances
```

Variable	Factor1	Factor2	Uniqueness
pre_ia	-0.0386	0.4460	0.7996
p1	0.3446	0.4067	0.7159
p2	-0.0519	0.3627	0.8658
p3rev	0.1832	-0.0364	0.9651
ia_1	0.4135	-0.0300	0.8281
ia_2	0.8336	0.0187	0.3048
ia_3	0.7864	0.2402	0.3239
ia4rev	0.2639	-0.2609	0.8623
ia_5	0.4265	-0.3483	0.6968
ia_6	-0.3073	0.1049	0.8946
ia_7	0.0455	0.4274	0.8152
ia_8	-0.1465	0.4047	0.8148

EIGENVALUES: The variance in the variables accounted for by this factor.

Critical eigenvalue: > 1

Sample EFA output

```
factor pre_ia p1-p2 p3rev ia_1-ia_3 ia4rev ia_5-ia_8, fa(2)
```

Factor	Eigenvalue	Difference	Proportion	Cumulative
Factor1	2.01020	0.90722	0.5268	0.5268
Factor2	1.10298	0.30999	0.2890	0.8158
...				

```
LR test: independent vs. saturated: chi2(66) = 179.46 Prob>chi2 = 0.0000
```

```
Factor loadings (pattern matrix) and unique variances
```

Variable	Factor1	Factor2	Uniqueness
pre_ia	-0.0386	0.4460	0.7996
p1	0.3446	0.4067	0.7159
p2	-0.0519	0.3627	0.8658
p3rev	0.1832	-0.0364	0.9651
ia_1	0.4135	-0.0300	0.8281
ia_2	0.8336	0.0187	0.3048
ia_3	0.7864	0.2402	0.3239
ia4rev	0.2639	-0.2609	0.8623
ia_5	0.4265	-0.3483	0.6968
ia_6	-0.3073	0.1049	0.8946
ia_7	0.0455	0.4274	0.8152
ia_8	-0.1465	0.4047	0.8148

Proportion of variance explained by that factor.

Sample EFA output

```
factor pre_ia p1-p2 p3rev ia_1-ia_3 ia4rev ia_5-ia_8, fa(2)
```

Factor	Eigenvalue	Difference	Proportion	Cumulative
Factor1	2.01020	0.90722	0.5268	0.5268
Factor2	1.10298	0.30999	0.2890	0.8158
...				

```
LR test: independent vs. saturated: chi2(66) = 179.46 Prob>chi2 = 0.0000
```

```
Factor loadings (pattern matrix) and unique variances
```

Variable	Factor1	Factor2	Uniqueness
pre_ia	-0.0386	0.4460	0.7996
p1	0.3446	0.4067	0.7159
p2	-0.0519	0.3627	0.8658
p3rev	0.1832	-0.0364	0.9651
ia_1	0.4135	-0.0300	0.8281
ia_2	0.8336	0.0187	0.3048
ia_3	0.7864	0.2402	0.3239
ia4rev	0.2639	-0.2609	0.8623
ia_5	0.4265	-0.3483	0.6968
ia_6	-0.3073	0.1049	0.8946
ia_7	0.0455	0.4274	0.8152
ia_8	-0.1465	0.4047	0.8148

Cumulative variance explained by all factors up to that point.

Sample EFA output

```
factor pre_ia p1-p2 p3rev ia_1-ia_3 ia4rev ia_5-ia_8, fa(2)
```

Factor	Eigenvalue	Difference	Proportion	Cumulative
Factor1	2.01020	0.90722	0.5268	0.5268
Factor2	1.10298	0.30999	0.2890	0.8158
...				

```
LR test: independent vs. saturated: chi2(66) = 179.46 Prob>chi2 = 0.0000
```

```
Factor loadings (pattern matrix) and unique variances
```

Variable	Factor1	Factor2	Uniqueness
pre_ia	-0.0386	0.4460	0.7996
p1	0.3446	0.4067	0.7159
p2	-0.0519	0.3627	0.8658
p3rev	0.1832	-0.0364	0.9651
ia_1	0.4135	-0.0300	0.8281
ia_2	0.8336	0.0187	0.3048
ia_3	0.7864	0.2402	0.3239
ia4rev	0.2639	-0.2609	0.8623
ia_5	0.4265	-0.3483	0.6968
ia_6	-0.3073	0.1049	0.8946
ia_7	0.0455	0.4274	0.8152
ia_8	-0.1465	0.4047	0.8148

UNIQUENESS is the proportion of variance in a variable that is not accounted for in the factor model. These are almost *never* '0,' no matter how many factors you have.

Sample EFA output

```
factor pre_ia p1-p2 p3rev ia_1-ia_3 ia4rev ia_5-ia_8, fa(2)
```

Factor	Eigenvalue	Difference	Proportion	Cumulative
Factor1	2.01020	0.90722	0.5268	0.5268
Factor2	1.10298	0.30999	0.2890	0.8158
...				

LR test: independent vs. saturated: chi2(66) = 179.46 Prob>chi2 = 0.0000

Factor loadings (pattern matrix) and unique variances

Variable	Factor1	Factor2	Uniqueness
pre_ia	-0.0386	0.4460	0.7996
p1	0.3446	0.4067	0.7159
p2	-0.0519	0.3627	0.8658
p3rev	0.1832	-0.0364	0.9651
ia_1	0.4135	-0.0300	0.8281
ia_2	0.8336	0.0187	0.3048
ia_3	0.7864	0.2402	0.3239
ia4rev	0.2639	-0.2609	0.8623
ia_5	0.4265	-0.3483	0.6968
ia_6	-0.3073	0.1049	0.8946
ia_7	0.0455	0.4274	0.8152
ia_8	-0.1465	0.4047	0.8148

LR TEST: Tests for differences between the estimated and actual covariance matrices. You want this to be insignificant – but it almost never is in large samples.

Aims of this presentation

1. Brief Overview of Primary Methods
2. Outline of Steps to Refine Your Module
3. Fully Worked Factor Analysis Example in Stata
4. Example Test of Our Construct's Validity

→ We will work on (3) and (4) in Stata.

Special thanks to:

Kathleen A. Cagney

Louise Hawkley

Brent W. Roberts

L. Phillip Schumm

Linda J. Waite

James Iveniuk

Joshua G. Mausolf
Department of Sociology
The University of Chicago
jmausolf@uchicago.edu

Thank You!



NORC
at the UNIVERSITY of CHICAGO

Slides, Data, and Code Available for Download:
<https://uchicago.box.com/v/factor>

Please Feel Free to Email with Any Questions

 insight for informed decisions™