

PCA & Exploratory Factor Analysis Quick Reference

PCA: Principal Component Analysis

Goal: Reduce dimensions, maximize variance

Steps:

1. Standardize data (mean=0, sd=1)
2. Compute correlation matrix **R**
3. Eigendecomposition: $\mathbf{R}\mathbf{v} = \lambda\mathbf{v}$
4. Project data onto eigenvectors

Key formula:

$$PC_j = v_{1j}X_1 + v_{2j}X_2 + \dots + v_{pj}X_p$$

Variance explained by PC_j :

$$\frac{\lambda_j}{\sum_i \lambda_i} \times 100\%$$

How many components?

Kaiser: keep if $\lambda > 1$

Scree: look for "elbow"

Cumulative variance $\geq 80\%$

Loadings vs. Scores:

Loading: correlation of variable with PC

Score: each observation's value on PC

EFA: Exploratory Factor Analysis

Goal: Find latent factors causing correlations

Model: $\mathbf{X} = \mathbf{\Lambda}\mathbf{F} + \boldsymbol{\varepsilon}$

X = observed variables

$\mathbf{\Lambda}$ = loading matrix

F = latent factors

$\boldsymbol{\varepsilon}$ = unique variance

Pre-checks:

KMO > 0.6 (sampling adequacy)

Bartlett $p < 0.05$ (not identity matrix)

How many factors?

Parallel analysis (recommended)

Scree plot

Theory / content knowledge

Loading interpretation:

$|\lambda| > 0.4$: meaningful

$|\lambda| > 0.7$: strong

Communality:

$$h_i^2 = \sum_j \lambda_{ij}^2 \quad (\text{shared variance})$$

PCA vs. EFA: When to Use

PCA:

Data compression / reduction

No model assumptions

Components = linear combos of *all* variables

Use when: preprocessing, visualization

EFA:

Discover latent structure

Model-based (factors *cause* correlations)

Factors explain correlations, not total variance

Use when: questionnaire validation, scale construction

Decision:

"Reduce variables?" \rightarrow PCA

"Find hidden constructs?" \rightarrow EFA

"Factors expected to correlate?" \rightarrow promax

"Factors independent?" \rightarrow varimax

Key distinction:

PCA: total variance (including unique)

EFA: shared variance (communality) only

Rotation Methods

Varimax (orthogonal):

Factors stay uncorrelated

Maximizes simple structure

Most common default

Promax (oblique):

Factors may correlate

Check Φ matrix for correlations

More realistic for psychology / social science

Rule of thumb:

If Φ correlations < 0.3 \rightarrow varimax sufficient

If Φ correlations ≥ 0.3 \rightarrow use promax

Simple structure (Thurstone):

Each variable loads high on one factor

Each factor has some high, some near-zero loadings

R Commands

PCA:

`pca <- prcomp(data, scale.=TRUE)`

`summary(pca)`

`pca$rotation`

`pca$x`

`biplot(pca)`

`screeplot(pca)`

EFA:

`library(psych)`

`KMO(cor(data))`

`fa.parallel(data)`

`fa(data, nfactors=k, rotate="varimax")`

`fa(..., rotate="promax")`

`$loadings`

`$communalities`

`$Phi`

variance explained

loadings

scores

visualization

scree plot

sampling adequacy

number of factors

factor loadings

shared variance

factor correlations (promax)

Common Mistakes

PCA:

Forgetting `scale.=TRUE`

Interpreting PCA loadings as factor loadings (they differ)

Keeping too many components (use scree + variance)

EFA:

Choosing # factors by Kaiser alone

\rightarrow use parallel analysis instead

Using varimax when factors should correlate

Reporting unrotated loadings

Ignoring communalities < 0.3 (poorly explained items)

Not checking KMO / Bartlett before EFA

Both:

Using PCA when you need EFA (or vice versa)

Skipping standardization

Over-interpreting borderline loadings (0.3-0.4)

Always report: KMO/Bartlett, # factors + criterion, rotation method, rotated loadings, variance explained, communalities