

# Estimation of the covariance matrix based on two types of the forward search algorithm

**Aleš Toman**

23<sup>rd</sup> International Workshop on Matrices and Statistics  
Ljubljana, June 9, 2014

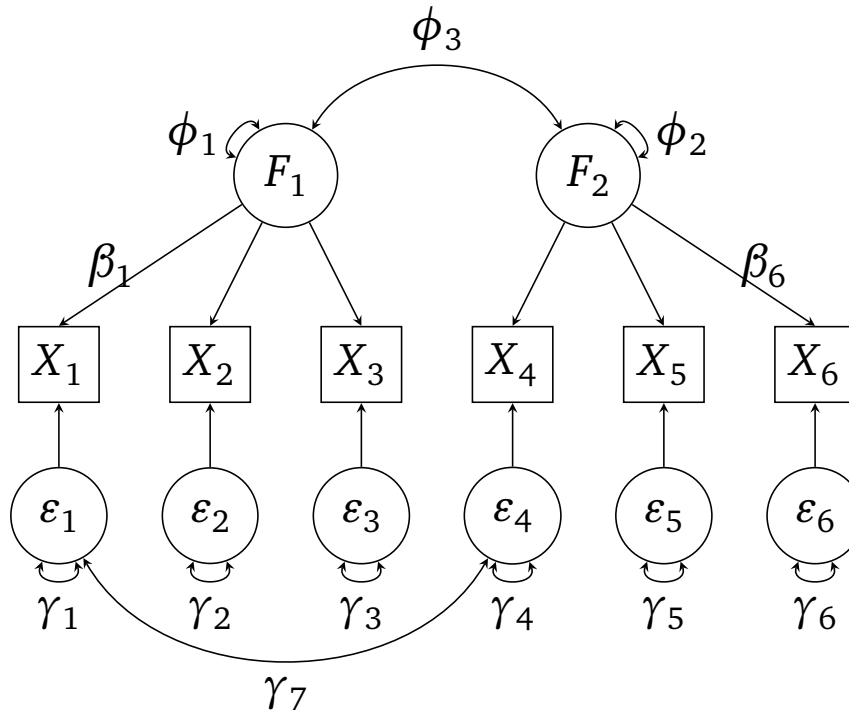
# Outline

- ▷ Confirmatory factor analysis
- ▶ Forward search algorithm
  - ▷ Outlier identification
  - ▶ Robust covariance matrix estimate
  - ▶ Robust confirmatory factor analysis
- ▶ Conclusions



*The research was partially supported by Slovenian Research Agency and IMFM.*

# Confirmatory factor analysis [2]



Linear dependence:  $X = \mu + \Lambda F + \varepsilon$ .

Does the restricted model make *reasonable fit* to the data?

# Confirmatory factor analysis [2]

- ▶  $\text{var}(X) = \Sigma \quad p \times p$
- ▶  $\text{var}(F) = \Phi \quad q \times q$
- ▶  $\text{var}(\varepsilon) = \Psi \quad p \times p$

Model implied covariance matrix  $\tilde{\Sigma} = \Lambda\Phi\Lambda^T + \Psi$ .

## *Maximum likelihood estimates*

- ▶ Multivariate normal distribution.
- ▶  $\hat{\Sigma}$  maximum likelihood estimate of the covariance matrix.
- ▶ Minimize  $F_{\text{ML}} = \text{trace}(\tilde{\Sigma}^{-1}\hat{\Sigma}) - \log(\det(\tilde{\Sigma}^{-1}\hat{\Sigma})) - p$ .

Model estimation and fit evaluation are *based on the matrix*  $\hat{\Sigma}$ !

# Forward search algorithm [1]

The forward search algorithm is an *iterative method*, that orders the data according to their *distances from the proposed model*. It helps us to identify observations with disproportionately *high influence* on statistical inference.

# Forward search algorithm [1]

The forward search algorithm is an *iterative method*, that orders the data according to their *distances from the proposed model*. It helps us to identify observations with disproportionately *high influence* on statistical inference.

1. Split the sample into 2 subsets ► outlier free *basic set*,  
► non-basic set.
2. Add observations to the basic set.
3. Use *forward plots* to show the dynamics of estimates.

The algorithm enables ► data exploration,  
► robust parameter estimation.

# Data

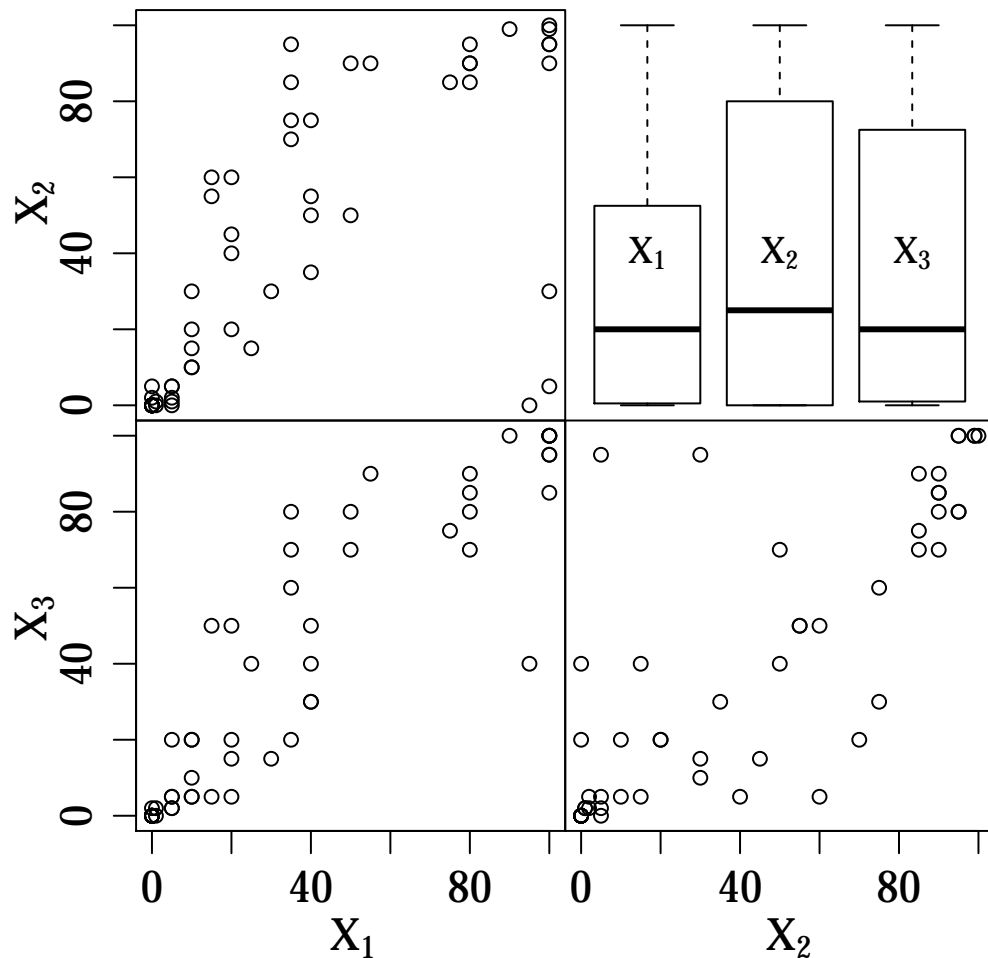
**G. W. Cermak, K. A. Bollen:** *Observer consistency in judging extent of cloud cover*, *Atmospheric Environment*, **17** (1983) 2109–2121.

- ▶  $n = 60$  slides (July 1980).
- ▶  $p = 3$  judges.
- ▶ Percent of the sky containing clouds.



# Data

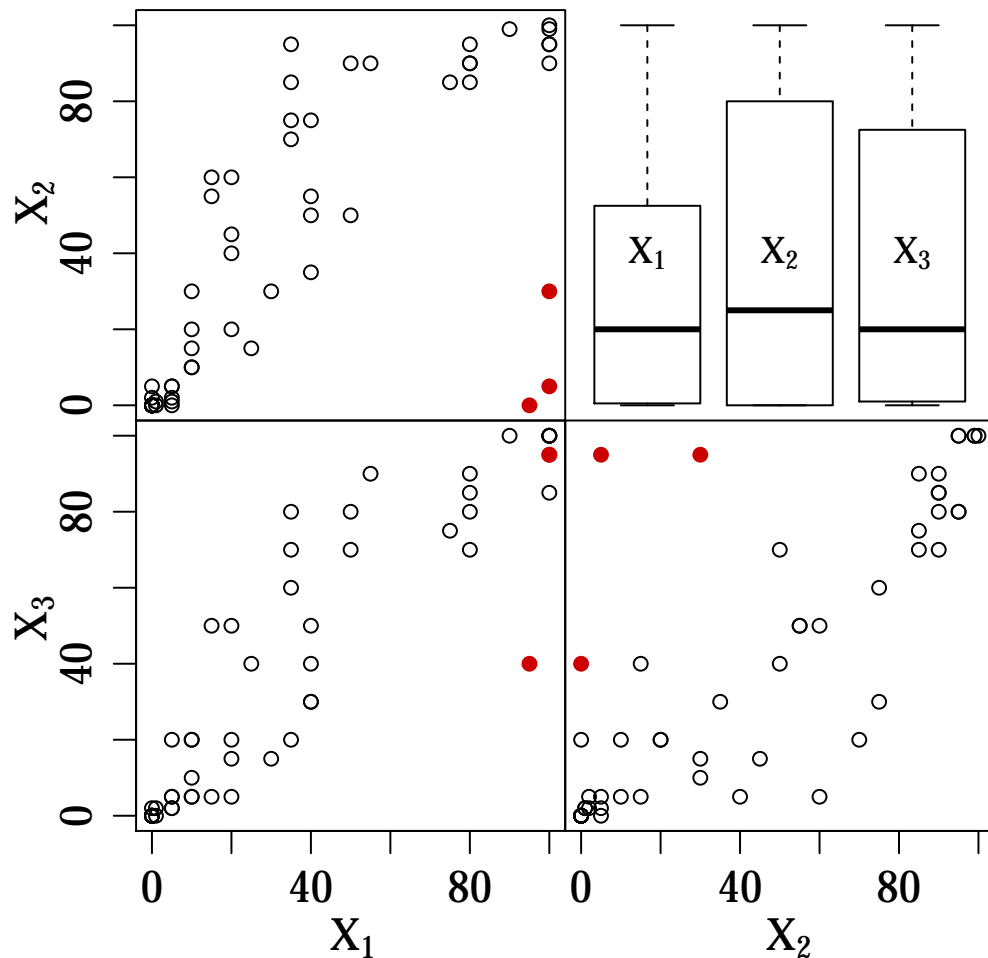
- ▶ Small sample?
- ▶ Asymmetric and non-normal distribution.
- ▶ Frequently used.





## Data

- ▶ Small sample?
- ▶ Asymmetric and non-normal distribution.
- ▶ Frequently used.
- ▶ Observations 40, 51, 52.



# Full sample estimates

$$\bar{x} = \begin{bmatrix} 32.95 \\ 37.65 \\ 35.55 \end{bmatrix}$$

$$S = \begin{bmatrix} 1301 & 1020 & 1237 \\ 1020 & 1463 & 1200 \\ 1237 & 1200 & 1404 \end{bmatrix}$$

# Robust covariance matrix estimate [3]

1. Split the sample into 2 subsets
  - ▶ outlier free basic set,
  - ▶ non-basic set.

$S$  sample covariance matrix

$S_{(-i)}$  sample covariance matrix *with observation  $i$  excluded*

$$S = \begin{bmatrix} \triangle_1 & \cdot & \cdot \\ \square_1 & \square_2 & \cdot \\ \pentagon_1 & \pentagon_2 & \pentagon_3 \end{bmatrix} \Rightarrow \text{vecs}(S) = \begin{bmatrix} \triangle_1 \\ \square_1 \\ \square_2 \\ \pentagon_1 \\ \pentagon_2 \\ \pentagon_3 \end{bmatrix}$$

▶  $s = \text{vecs}(S)$

▶  $s_{(-i)} = \text{vecs}(S_{(-i)})$

# Robust covariance matrix estimate [3]

1. Split the sample into 2 subsets ► outlier free basic set,  
► non-basic set.

(Squared) *Cook's distance* of observation  $i$

$$CD_i^2 = (s_{(-i)} - s)^T (\widehat{\text{cov}} s)^{-1} (s_{(-i)} - s)$$

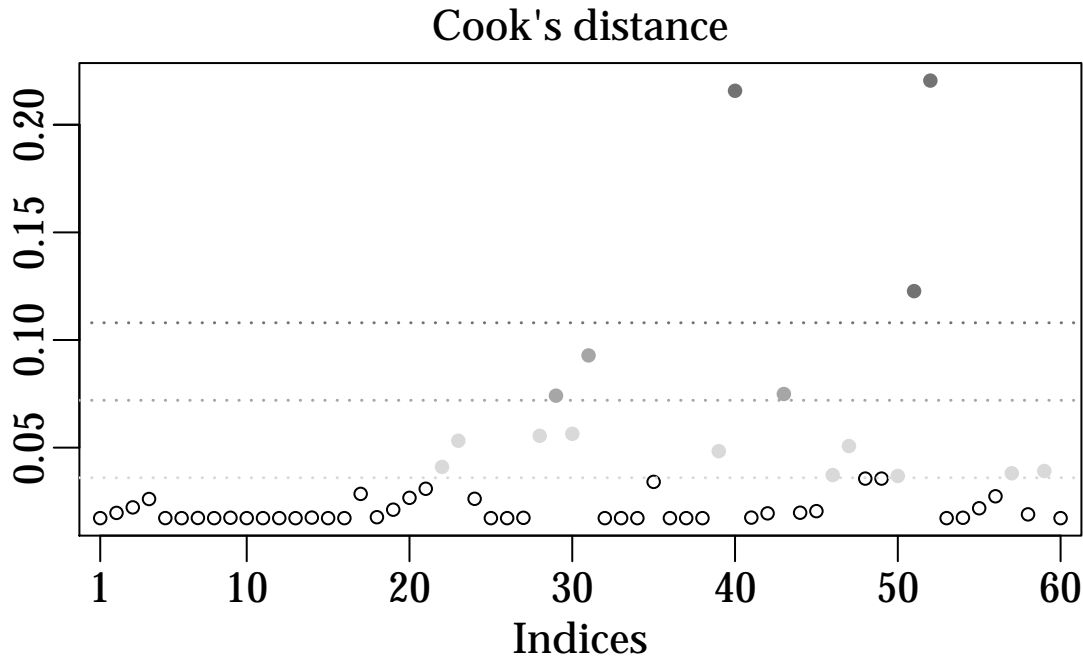
Multivariate normal distribution of  $X$ :

- $\text{cov}(s_{gh}, s_{jk}) \propto \sigma_{gj}\sigma_{hk} + \sigma_{gk}\sigma_{hj}$
- $\widehat{\text{cov}}(s_{gh}, s_{jk}) \propto s_{gj}s_{hk} + s_{gk}s_{hj}$

Take  $m = 30$  observations with lowest Cook's distances.

# Robust covariance matrix estimate [3]

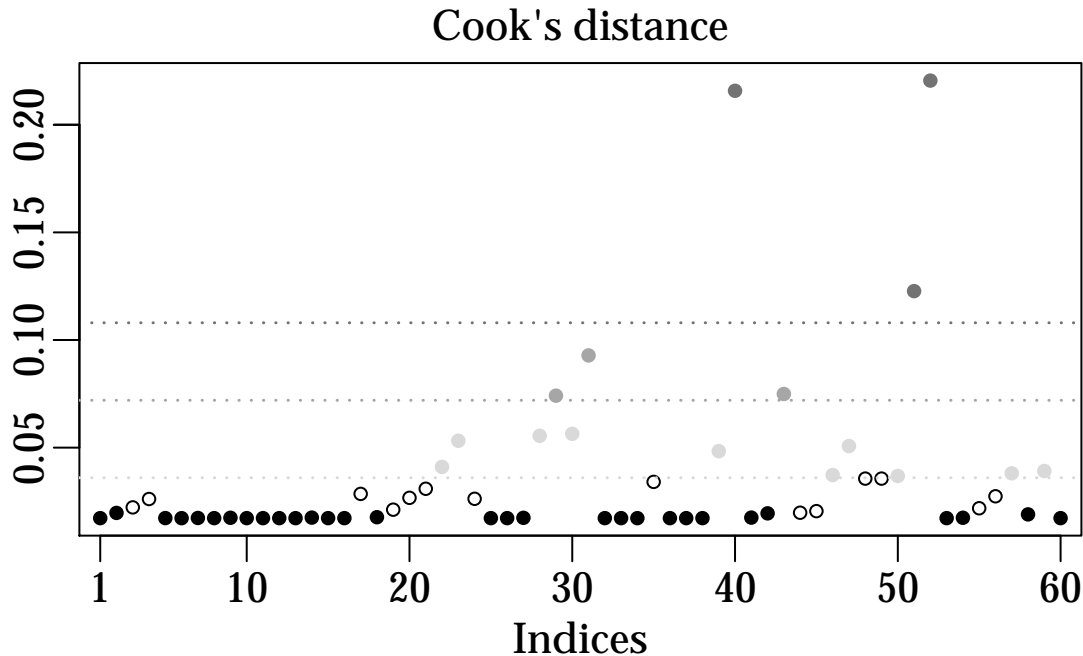
1. Split the sample into 2 subsets ► outlier free basic set,  
► non-basic set.



$$CD = \frac{1}{n} \sum_{i=1}^n CD_i \text{ and levels } 1 \cdot CD, 2 \cdot CD \text{ in } 3 \cdot CD$$

# Robust covariance matrix estimate [3]

1. Split the sample into 2 subsets ► outlier free basic set,  
► non-basic set.



$$CD = \frac{1}{n} \sum_{i=1}^n CD_i \text{ and levels } 1 \cdot CD, 2 \cdot CD \text{ in } 3 \cdot CD$$

# Robust covariance matrix estimate [3]

## 2. Add observations to the basic set.

$X^{(\ell)}$  the basic set with  $\ell$  observations: we wish to include one more.

$S^{(\ell)}$  sample covariance matrix of the basic set.

- ▶  $S_{(-i)}^{(\ell)}$  covariance matrix *with observation  $i$  excluded*,
- ▶  $S_{(+i)}^{(\ell)}$  covariance matrix *with observation  $i$  added*.

$$s^{(\ell)} = \text{vecs}(S^{(\ell)}) \quad s_{(-i)}^{(\ell)} = \text{vecs}(S_{(-i)}^{(\ell)}) \quad s_{(+i)}^{(\ell)} = \text{vecs}(S_{(+i)}^{(\ell)})$$

# Robust covariance matrix estimate [3]

## 2. Add observations to the basic set.

(Squared) *Cook's distance* of observation  $i$

$$CD_i^{2(\ell)} = \begin{cases} (s_{(-i)}^{(\ell)} - s^{(\ell)})^T (\widehat{\text{cov}} s^{(\ell)})^{-1} (s_{(-i)}^{(\ell)} - s^{(\ell)}); & i \in X^{(\ell)} \\ (s_{(+i)}^{(\ell)} - s^{(\ell)})^T (\widehat{\text{cov}} s^{(\ell)})^{-1} (s_{(+i)}^{(\ell)} - s^{(\ell)}); & i \notin X^{(\ell)} \end{cases}$$

$$\blacktriangleright \widehat{\text{cov}}(s_{gh}^{(\ell)}, s_{jk}^{(\ell)}) \propto s_{gj}^{(\ell)} s_{hk}^{(\ell)} + s_{gk}^{(\ell)} s_{hj}^{(\ell)}$$

$CD_i^{(\ell)}$  measures the *the influence of observation  $i$  on  $S^\ell$* .

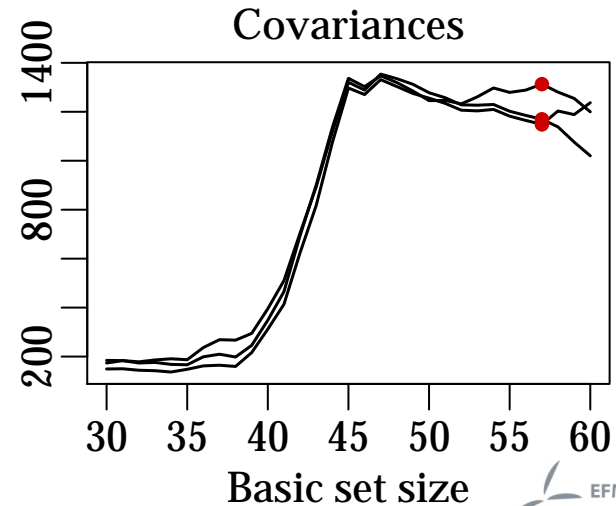
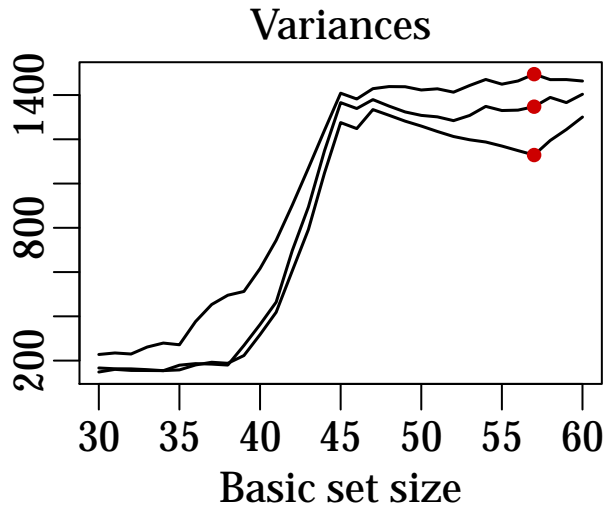
Take  $\ell + 1$  observations with lowest Cook's distances.



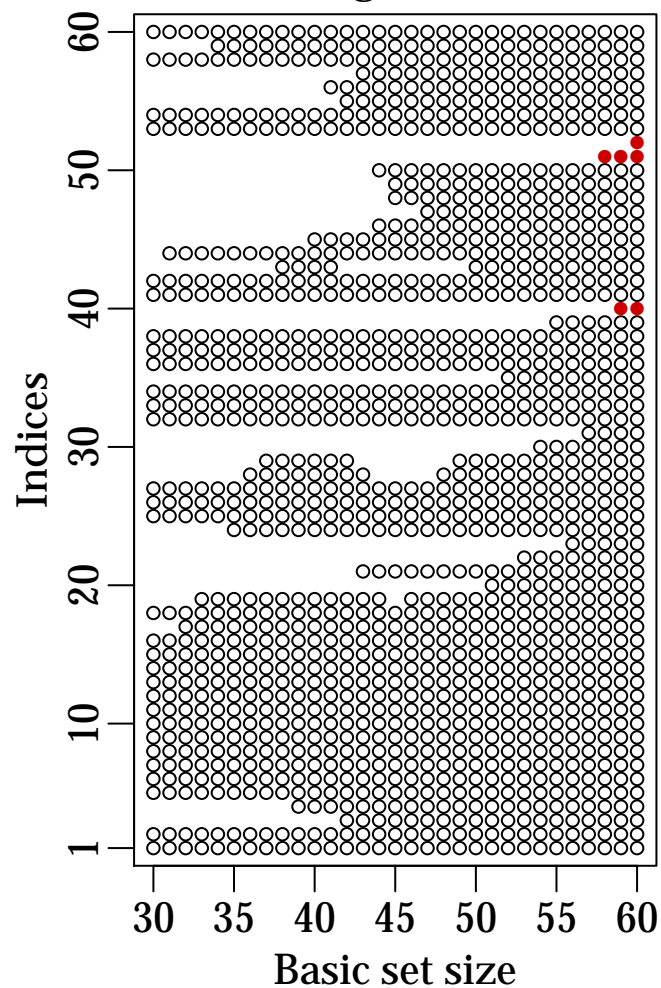
# Robust covariance matrix estimate [3]

3. Use forward plots to show the dynamics of estimates.

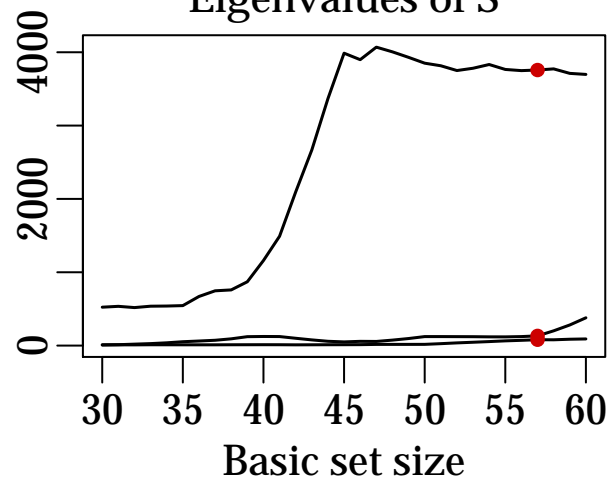
- ▶ Ordering of the data.
- ▶ Variances, covariances, and their functions.
- ▶ Cook's distances.



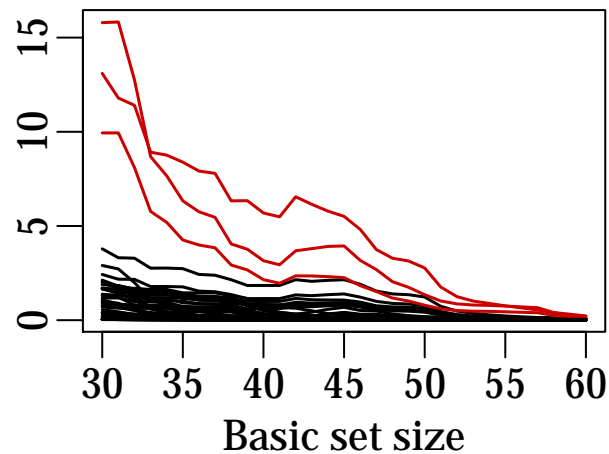
### Ordering of the data



### Eigenvalues of S



### Cook's distances

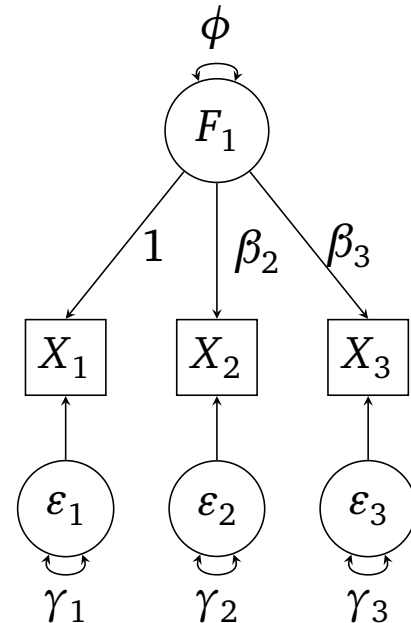


# Data and the confirmatory factor model

$$\bar{x} = \begin{bmatrix} 32.95 \\ 37.65 \\ 35.55 \end{bmatrix}$$

$$\hat{\Sigma} = \begin{bmatrix} 1279 & 1003 & 1216 \\ 1003 & 1439 & 1180 \\ 1216 & 1180 & 1380 \end{bmatrix}$$

- ▶ One-factor model.
- ▶ Independent errors.

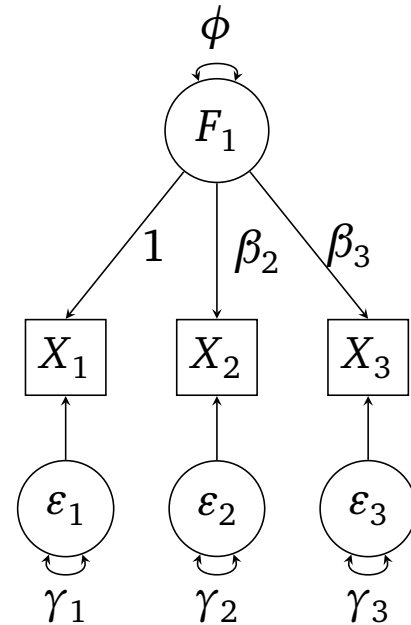


# Data and the confirmatory factor model

$$\bar{x} = \begin{bmatrix} 32.95 \\ 37.65 \\ 35.55 \end{bmatrix}$$

$$\hat{\Sigma} = \begin{bmatrix} 1279 & 1003 & 1216 \\ 1003 & 1439 & 1180 \\ 1216 & 1180 & 1380 \end{bmatrix}$$

- ▶ One-factor model.
- ▶ Independent errors.
- ▶ *Heywood case*:  $\hat{\gamma}_3 = -50.584$



# Robust confirmatory factor analysis [4]

1. Split the sample into 2 subsets ► outlier free basic set,  
► non-basic set.

Distance from the model is measured with *observational residuals*

$$e_i = x_i - \hat{x} = x_i - \bar{x} - \hat{\Lambda} \hat{f}_i.$$

Regression method of *factor scores estimation*

$$\hat{f}_i = \hat{\Phi} \hat{\Lambda}^T \hat{\Sigma}^{-1} (x_i - \bar{x}), \quad \text{where} \quad \hat{\Sigma} = \hat{\Lambda} \hat{\Phi} \hat{\Lambda}^T + \hat{\Psi}.$$

$$e_i = (I - \hat{\Lambda} \hat{\Phi} \hat{\Lambda}^T \hat{\Sigma}^{-1}) (x_i - \bar{x})$$

*Summarized observational residuals* are

$$(e_i^S)^2 = e_i^T (\hat{\Psi} \hat{\Sigma}^{-1} \hat{\Psi})^{-1} e_i.$$

# Robust confirmatory factor analysis [4]

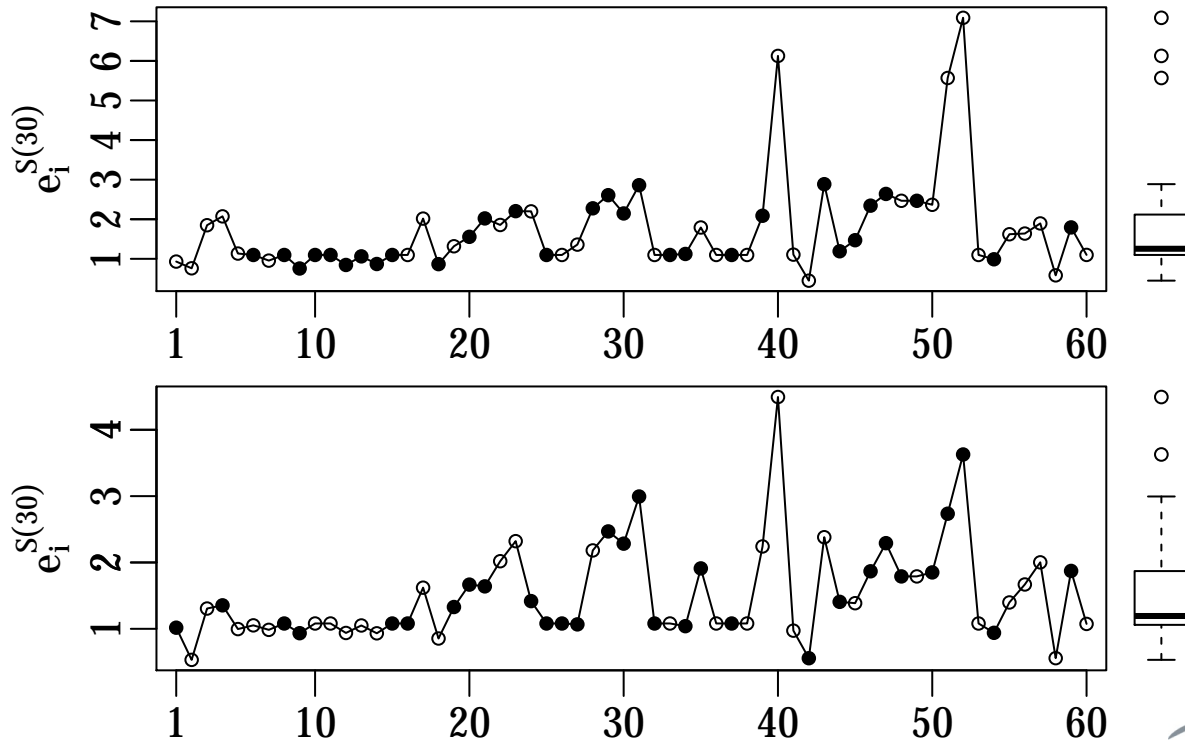
1. Split the sample into 2 subsets ► outlier free basic set,  
► non-basic set.

1. Take a *random subet*  $X^{(m)}$  of size  $m = 30$ .
2. Estimate the confirmatory factor model on the subset.
3. Compute  $e_i^{(m)} = \left( I - \widehat{\Lambda}^{(m)} \widehat{\Phi}^{(m)} (\widehat{\Lambda}^{(m)})^T (\widehat{\Sigma}^{(m)})^{-1} \right) (x_i - \bar{x}^{(m)})$ .
4. Compute  $(e_i^{S(m)})^2 = (e_i^{(m)})^T \left( \widehat{\Psi}^{(m)} (\widehat{\Sigma}^{(m)})^{-1} \widehat{\Psi}^{(m)} \right)^{-1} e_i^{(m)}$ .
5. Find the median of summarized observational residuals.

**Repeat 1000-times** and take the subset with the smallest median.

# Robust confirmatory factor analysis [4]

1. Split the sample into 2 subsets ► outlier free basic set,  
► non-basic set.



# Robust confirmatory factor analysis [4]

## 2. Add observations to the basic set.

$X^{(\ell)}$  the basic set with  $\ell$  observations: we wish to include one more.

1. Estimate the confirmatory factor model on the basic set.
2. Compute  $e_i^{(\ell)} = \left( I - \widehat{\Lambda}^{(\ell)} \widehat{\Phi}^{(\ell)} (\widehat{\Lambda}^{(\ell)})^T (\widehat{\Sigma}^{(\ell)})^{-1} \right) (x_i - \bar{x}^{(\ell)})$ .
3. Compute  $(e_i^{S(\ell)})^2 = (e_i^{(\ell)})^T \left( \widehat{\Psi}^{(\ell)} (\widehat{\Sigma}^{(\ell)})^{-1} \widehat{\Psi}^{(\ell)} \right)^{-1} e_i^{(\ell)}$ .
4. Take  $\ell + 1$  observations with smallest summarized observational residuals.



# Robust confirmatory factor analysis [4]

## 2. Add observations to the basic set.

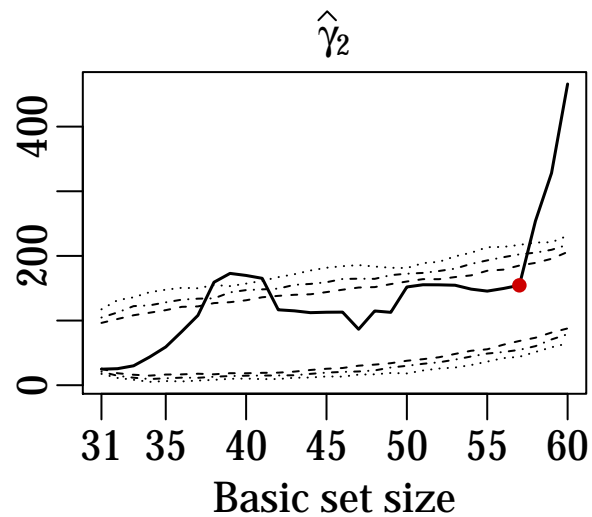
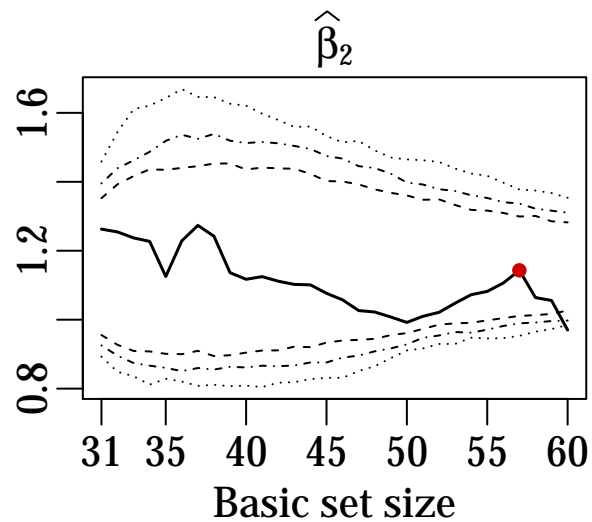
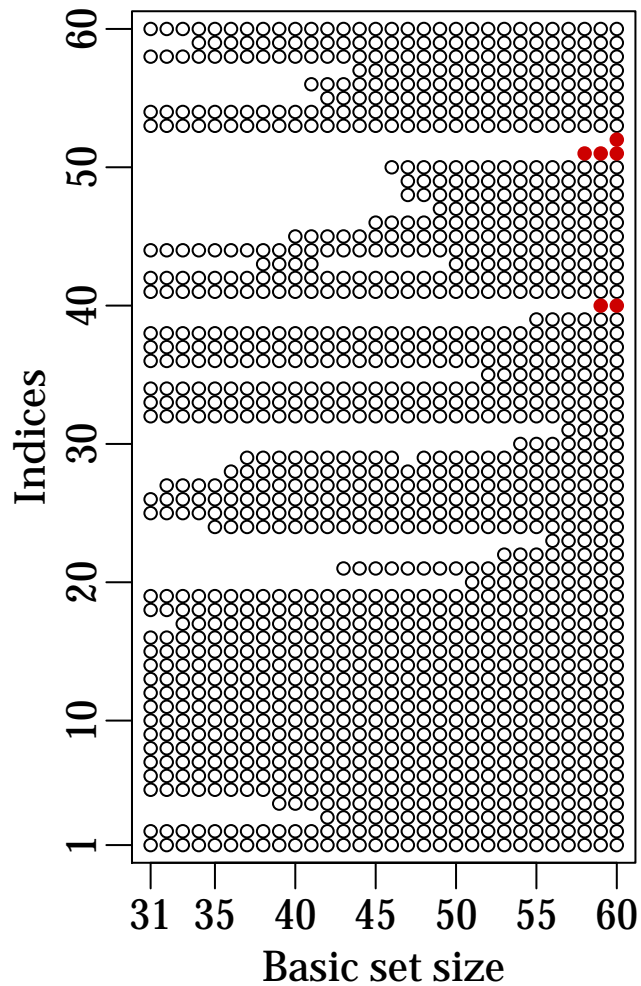
$X^{(\ell)}$  the basic set with  $\ell$  observations: we wish to include one more.

1. Estimate the confirmatory factor model on the basic set.
2. Compute  $e_i^{(\ell)} = \left( I - \widehat{\Lambda}^{(\ell)} \widehat{\Phi}^{(\ell)} (\widehat{\Lambda}^{(\ell)})^T (\widehat{\Sigma}^{(\ell)})^{-1} \right) (x_i - \bar{x}^{(\ell)})$ .
3. Compute  $(e_i^{S(\ell)})^2 = (e_i^{(\ell)})^T \left( \widehat{\Psi}^{(\ell)} (\widehat{\Sigma}^{(\ell)})^{-1} \widehat{\Psi}^{(\ell)} \right)^{-1} e_i^{(\ell)}$ .
4. Take  $\ell + 1$  observations with smallest summarized observational residuals.

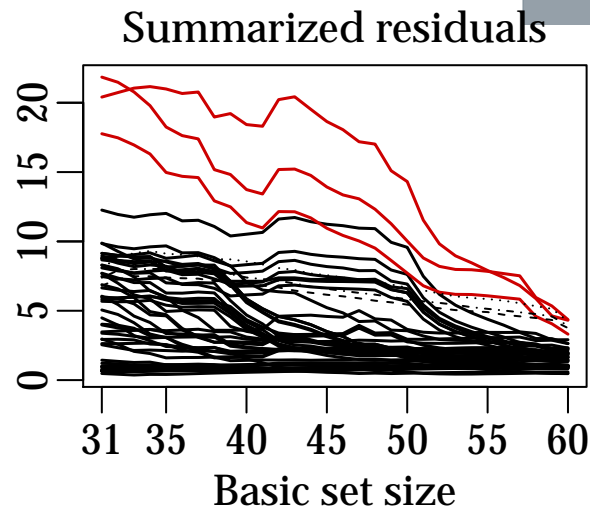
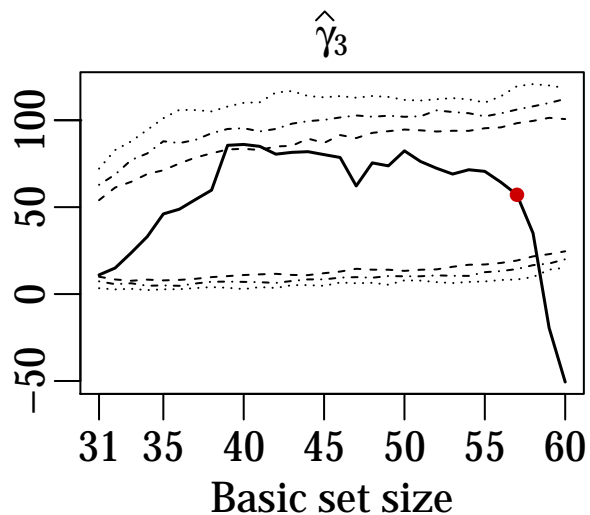
## 3. Use forward plots to show the dynamics of estimates.

- ▶ Ordering of the data.
- ▶ Parameter estimates, fit indices, summarized obs. residuals.

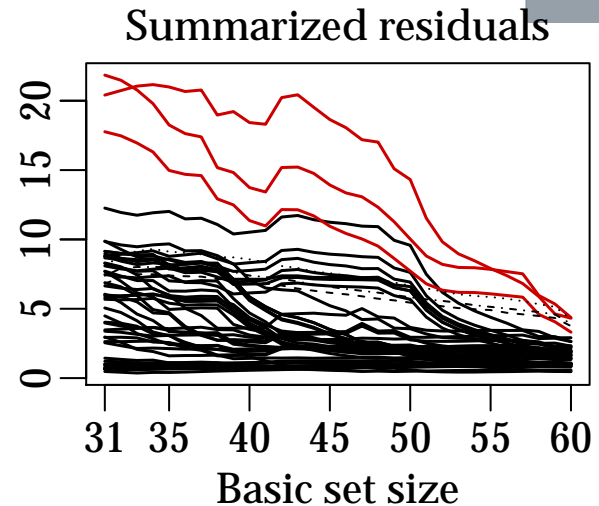
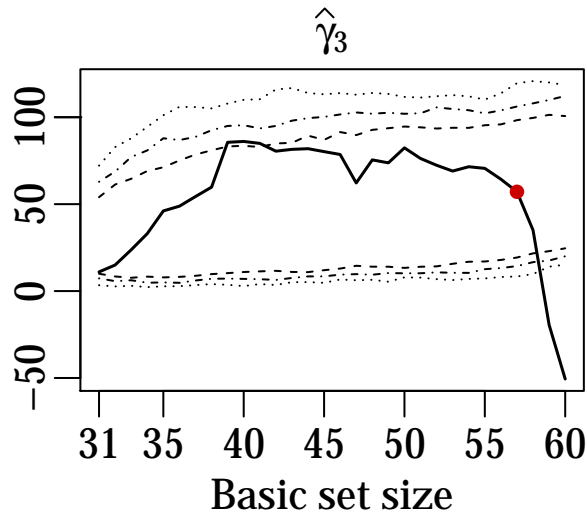
# Ordering of the data



# Robust confirmatory factor analysis [4]



# Robust confirmatory factor analysis [4]

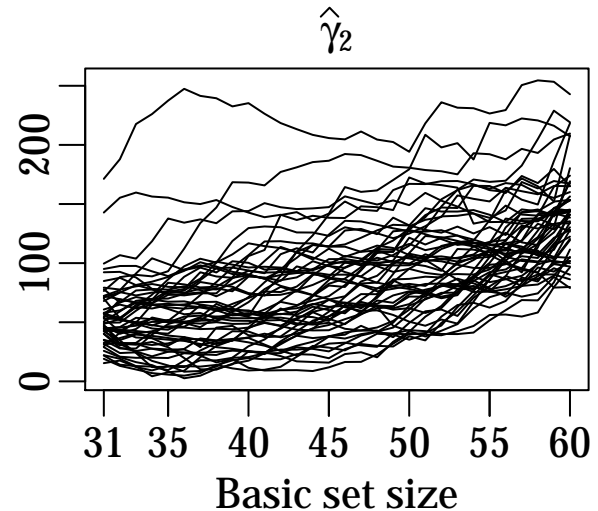
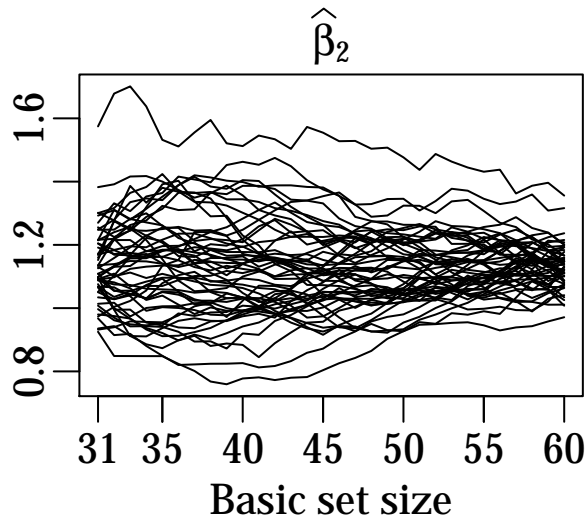


## Conclusions

- ▶ Important to consider the model in robust inference.
- ▶ The algorithm enables us to explore the data.
- ▷ Rarely used in applied research.

# Simulation envelopes

- ▶ Identify the step  $\ell^*$  of the first major changes.
- ▶ Take  $\hat{\Lambda}^{(\ell^*)}$ ,  $\hat{\Phi}^{(\ell^*)}$ ,  $\hat{\Psi}^{(\ell^*)}$  and compute  $\hat{\Sigma}^{(\ell^*)}$ .
- ▶ Simulate 1000 samples from  $N_p(0, \hat{\Sigma}^{(\ell^*)})$ .
- ▶ Repeat forward search analysis on each of the samples.



- ▶ Find pointwise confidence interval.

# References

- [1] **A. C. Atkinson, M. Riani, A. Cerioli:** *Exploring multivariate data with the forward search*, Springer-Verlag, New York, 2004.
- [2] **T. A. Brown:** *Confirmatory factor analysis for applied research*, The Guilford Press, New York, 2006.
- [3] **W.-Y. Poon, Y.-K. Wong:** *A forward search procedure for identifying influential observations in the estimation of a covariance matrix*, Structural Equation Modeling **11** (2004) 357–374.
- [4] **A. Toman:** *Robust confirmatory factor analysis based on the forward search algorithm*, Statistical Papers, **55** (2014) 233–252.

**Thank you for your attention!**

**[ales.toman@ef.uni-lj.si](mailto:ales.toman@ef.uni-lj.si)**

23<sup>rd</sup> International Workshop on Matrices and Statistics  
Ljubljana, June 9, 2014