

# Möglichkeiten der Verbesserung des EM-Algorithmus für normalverteilte Daten

Dipl. Ing. Luis Huergo

Universität Tübingen  
Lehrstuhl für Statistik, Ökonometrie und Unternehmensforschung  
Prof. Dr. Dr. h.c. mult. Eberhard Schaich  
([luis.huergo@uni-tuebingen.de](mailto:luis.huergo@uni-tuebingen.de))

08. Juli 2008

EBERHARD KARLS  
UNIVERSITÄT  
TÜBINGEN



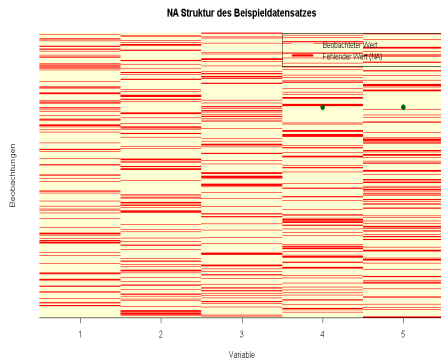
# Motivation

- Neuartige Indikatoren für die wissensbasierte Wirtschaft (knowledge-based economy) sollten zu einem einzigen zusammengesetzten Indikator aggregiert und seine Eigenschaften untersucht werden.
- Es wurden 116 Indikatoren ausgewählt.
- Kein einziges der 25 teilnehmenden europäischen Länder war in der Lage, alle 116 Indikatoren zu liefern.
- Der fertiggestellte Datensatz hatte 42% fehlende Daten.
- Diese fehlenden Daten mussten **imputiert** werden.

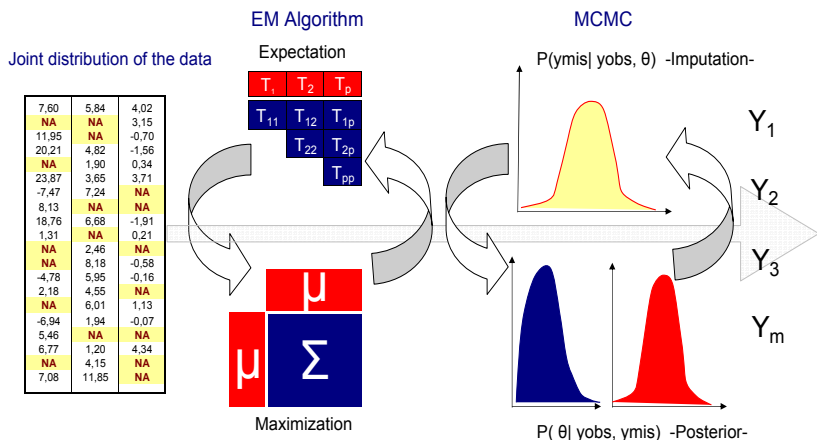


# Imputation

Die (moderne) Imputation eines mit fehlenden Werten behafteten Datensatzes kann mit der Rekonstruktion eines beschädigten Bildes verglichen werden.

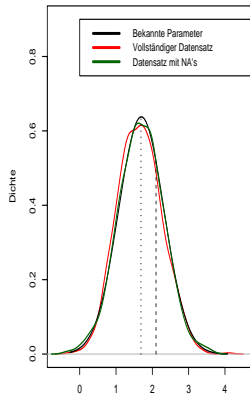


# Moderne Imputationsverfahren: EM-Algorithmus und MCMC-Methoden

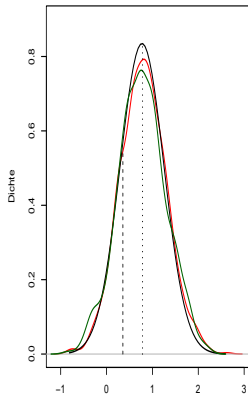


# Ergebnisse der Simulation

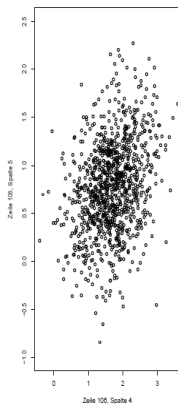
Zeile 106, Spalte 4



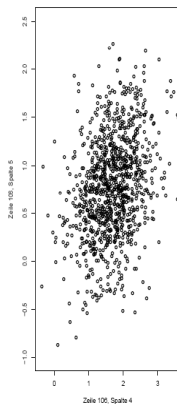
Zeile 106, Spalte 5



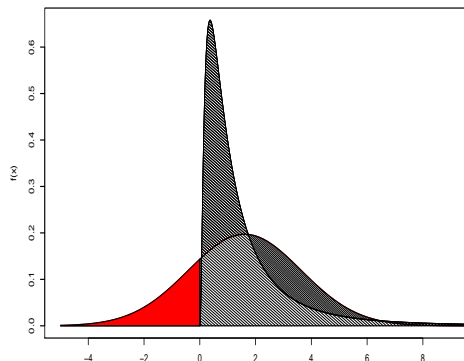
Bekannte Parameter



Datensatz mit NA's



# Probleme bei nicht normalverteilten Daten



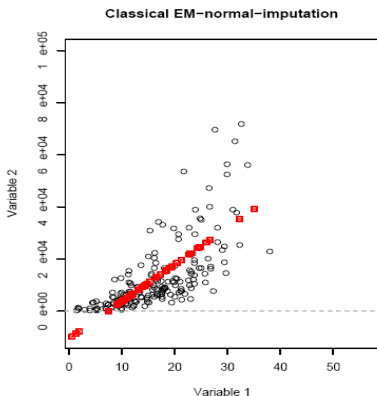
- Beide Verteilungen haben den gleichen Erwartungswert und die gleiche Varianz
- Die Werte der rot markierten Fläche sind unter der rechtsschiefen Verteilung unzulässig.

Der Algorithmus kann jedoch, anhand der ihm verpassten Parameter, ausschließlich diese Normalverteilung „sehen“.



# Stark vereinfachtes Beispiel

4,02	11135,84
3,15	9011,95
0,70	20,21
1,56	554,82
0,34	1,90
3,71	8113,65
NA	0,24
NA	4,55
1,91	6,68
0,21	6,01
NA	2,46
0,58	8,18
0,16	5,95
NA	954,55
1,13	6,01
0,07	1,94
NA	0,84
4,34	20201,20
NA	4,15
NA	11,85



- Aufgrund der Schiefe der Verteilungen weist die Punktwolke einen nichtlinearen Verlauf auf.
- Die roten Punkte unterhalb der gestrichelten grauen Linie sind unter der ersten Verteilung unzulässig.



# Momentenbedingungen im Falle einer multivariaten Normalverteilung

Gegeben sei  $Y := \mathcal{Y}^\theta$  normalverteilt.

$$\left. \begin{array}{l} E[(Y - \mu)^{3+2k}] = 0 \quad \text{für } k = 0, 1, \dots \\ E[(Y - \mu)^4] - 3\sigma^4 = 0 \\ E[(\varepsilon)^{3+2k}] = 0 \quad \text{für } k = 0, 1, \dots \\ E[(\varepsilon)^4] - 3\sigma_\varepsilon^4 = 0 \\ E[X\beta(\varepsilon^2 - \sigma_\varepsilon^2)] = 0 \end{array} \right\} \begin{array}{l} \text{Randverteilungen} \\ \text{Residuen} \end{array} \right\} =: g$$

Diese Bedingungen sind im Falle einer multivariaten Normalverteilung stets erfüllt.

Die Schätzstrategie besteht nun darin, dass eine geeignet gewichtete quadratische Form der empirischen Momentenbedingungen bezüglich des Potenzparameters minimiert wird:

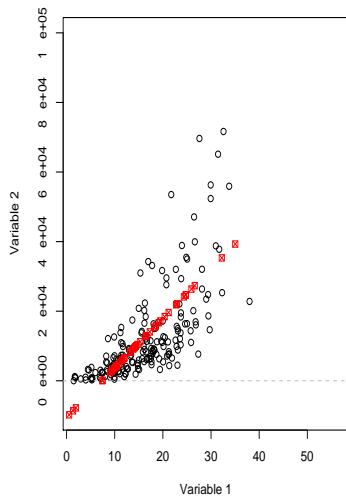
$$\min_{\theta} \hat{g}(\theta)' W \hat{g}(\theta)$$



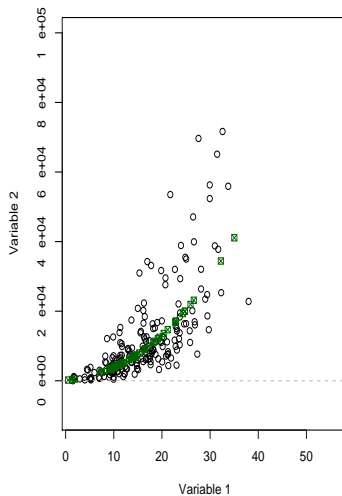


# Vergleich beider Methoden

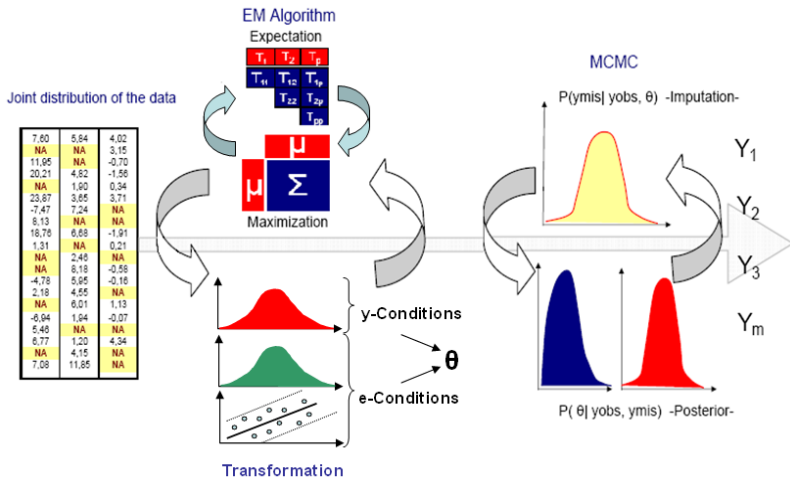
Classical EM-normal-imputation



Proposed EM-general-imputation



# EM-Algorithmus/Potenztransformation + MCMC



Vielen Dank für Ihre Aufmerksamkeit.

