

LANDWIRTSCHAFTLICHE FAKULTÄT  
DER  
RHEINISCHEN FRIEDRICH-WILHELMS-UNIVERSITÄT ZU BONN



Institut für Photogrammetrie

---

# Der EM-Algorithmus

## bei Schätz- und Klassifikationsproblemen

**Diplomarbeit**

zur  
Erlangung des Grades  
Diplom-Ingenieur  
(Dipl.-Ing.)  
der  
Fachrichtung  
Vermessungswesen

Vorgelegt am 13. November 2000 von  
**Marc Luxen**  
aus Aremberg



LANDWIRTSCHAFTLICHE FAKULTÄT  
DER  
RHEINISCHEN FRIEDRICH-WILHELMS-UNIVERSITÄT ZU BONN



Institut für Photogrammetrie

---

# Der EM-Algorithmus

## bei Schätz- und Klassifikationsproblemen

**Diplomarbeit**

zur  
Erlangung des Grades  
Diplom-Ingenieur  
(Dipl.-Ing.)  
der  
Fachrichtung  
Vermessungswesen

Vorgelegt am 13. November 2000 von  
**Marc Luxen**  
aus Aremberg

Betreuer: Prof. Dr.-Ing. Wolfgang Förstner  
Dipl.-Ing. Ansgar Brunn

Ausgegeben am: 20. Oktober 2000  
Vorgelegt am: 13. November 2000



## *Diplomaufgabe*

für Herrn cand. geod. Marc Luxen

### **Der EM-Algorithmus bei Schätz- und Klassifikationsproblemen**

Eine zentrale Aufgabe der Photogrammetrie besteht heute in der Entwicklung von Verfahren zur automatischen Wiedererkennung von Objekten in digitalen Bildern. Im Zusammenhang hiermit treten häufig Schätz- und Klassifikationsprobleme auf, die mit traditionellen geodätischen Parameterschätzverfahren nicht oder nur mit verhältnismäßig großem Aufwand gelöst werden können. Beziehen sich beispielsweise die Beobachtungen, die zur Bestimmung von objektrelevanten Parametern aus einem Bild abgeleitet werden, auf unterschiedliche Objekte und geht aus dem Beobachtungsmaterial nicht hervor, welche Beobachtung zu welchem Objekt gehört, so muß zur Schätzung der Objektparameter eine Zuordnung der Beobachtungen zu den einzelnen Objekten erfolgen. Mit klassischen Methoden der Parameterschätzung ist es nicht ohne weiteres möglich, dieses Zuordnungsproblem *gemeinsam* mit der Parameterschätzung zu lösen.

Im Jahre 1977 wurde in einem von der ROYAL STATISTICAL SOCIETY veröffentlichten Aufsatz der amerikanischen Professoren A.P. DEMPSTER, N.M. LAIRD und D.B. RUBIN<sup>1</sup> der sogenannte *Expectation Maximization (EM) Algorithmus* vorgestellt. Hierbei handelt es sich um ein Verfahren zur Schätzung von Parametern aus unvollständigen Beobachtungsdaten, das u.a. die Lösung des Zuordnungsproblems gemeinsam mit der Parameterschätzung erlaubt und daher für die automatische Bildinterpretation von großem Interesse ist.

Aufgabe des Diplomanden ist es, die theoretischen Grundlagen des EM-Algorithmus aufzuarbeiten und anhand eines frei wählbaren Beispiels aufzuzeigen, wie der EM-Algorithmus eingesetzt werden kann, wenn im Zusammenhang mit einer Parameterschätzung gleichzeitig Zuordnungsprobleme zu lösen sind.

(Prof. Dr.-Ing. W. Förstner)

Betreuer: W. Förstner, A. Brunn

Ausgegeben am: 20. Oktober 2000

Abgabetermin: 20. Januar 2001

Abgegeben am: 13. November 2000

---

<sup>1</sup>A.P. Dempster, N.M. Laird, D.B. Rubin (1977): *Maximum Likelihood from Incomplete Data via the EM Algorithm*, Journal of the Royal Statistical Society, Series B (Methodological), Bd. 39, Nr.1, 1977, S. 1-38

**E r k l ä r u n g :**

Hiermit versichere ich, die vorliegende Arbeit selbständig und nur unter Nutzung der im Literaturverzeichnis aufgeführten Quellen angefertigt zu haben.

Aremberg, 13. November 2000

# Inhaltsverzeichnis

<b>1</b>	<b>Einleitung</b>	<b>1</b>
11	Statistische Methoden in der Photogrammetrie . . . . .	1
12	Klassische Schätzverfahren und deren Grenzen . . . . .	3
121	Prinzip der klassischen geodätischen Ausgleichung . . . . .	3
122	Grenzen der klassischen geodätischen Ausgleichungsrechnung . . . . .	4
13	Parameterschätzung aus unvollständigen Daten . . . . .	5
14	Zusammenfassung / Überblick über die Arbeit . . . . .	6
<b>2</b>	<b>Der EM-Algorithmus</b>	<b>7</b>
21	Begriffe und Notation . . . . .	7
22	Maximum-Likelihood-Methode . . . . .	9
221	Definition . . . . .	9
222	Schätzung aus unvollständigen Beobachtungsdaten mittels Maximum-Likelihood-Methode . . . . .	11
23	Definition des (G)EM-Algorithmus . . . . .	12
231	Schlüsselgleichung des (G)EM-Algorithmus . . . . .	12
232	Definition des EM-Algorithmus . . . . .	16
233	Zusammenfassung . . . . .	17
<b>3</b>	<b>Eigenschaften des EM-Algorithmus</b>	<b>19</b>
31	Vorbemerkungen . . . . .	19
32	Eigenschaften des (G)EM-Algorithmus . . . . .	20
321	Zum Verhalten der Likelihoodfunktion während der (G)EM-Iterationen	21
322	Betrachtungen zur Konvergenz der logarithmierten Likelihoodfunktion . . . . .	23
323	Betrachtungen zur Konvergenz der Folge geschätzter Parameter . . . . .	29
324	Betrachtungen zur Konvergenzgeschwindigkeit . . . . .	34
33	Zusammenfassung . . . . .	41

<b>4</b>	<b>Beispiel: Parameterschätzung und Klassifikation</b>	<b>43</b>
41	Aufgabe: Bestimmung der Parameter zweier ausgleichender Kurven in einem Bild . . . . .	43
42	Lösung mittels des EM-Algorithmus: . . . . .	45
421	Formulierung der Aufgabe als Schätzproblem aus unvollständigen Beobachtungsdaten . . . . .	45
422	Ableitung der EM-Iterationsschritte . . . . .	47
423	Ergebnisse . . . . .	59
<b>5</b>	<b>Zusammenfassung</b>	<b>65</b>
<b>A</b>	<b>Anhang</b>	<b>69</b>
A1	C-Quellcode des Programms EM.c (Auszug) . . . . .	69
A2	Datensatz <code>Linien.dat</code> . . . . .	79



# Kapitel 1

## Einleitung

Die vorliegende Diplomarbeit beschäftigt sich mit dem sogenannten *Expectation Maximization (EM)-Algorithmus*, einem Verfahren zur Schätzung unbekannter Parameter aus unvollständigen Beobachtungsdaten. Die Arbeit stellt eine Zusammenstellung der Ergebnisse verschiedener Veröffentlichungen zu diesem Thema dar und zeigt darüber hinaus anhand eines einfachen Beispiels auf, wie der EM-Algorithmus zur Lösung von Schätz- und Klassifikationsproblemen eingesetzt werden kann.

Angesichts der Tatsache, daß Schätzverfahren primär im Rahmen des Fachbereiches Ausgleichsrechnung und Statistik behandelt werden, stellt sich die Frage, warum die Auseinandersetzung mit einem Verfahren zur Parameterschätzung aus unvollständigen Beobachtungsdaten in einer Diplomarbeit aus dem Bereich der *Photogrammetrie* erfolgt. In dieser Einleitung wird daher die Motivation zur vorliegenden Arbeit aus der Sicht der Photogrammetrie erläutert. Hierzu wird zunächst in Abschnitt 11 in allgemeiner Weise die enge Verknüpfung der Photogrammetrie mit der Ausgleichsrechnung und Statistik beschrieben. Anschließend werden kurz einige für diese Arbeit wesentliche Aspekte der klassischen Methoden der geodätischen Ausgleichsrechnung erläutert und es werden anhand eines Beispiels Grenzen dieser klassischen Methoden aufgezeigt. Abschließend erfolgt eine Abgrenzung der Parameterschätzung aus unvollständigen Daten von den klassischen Verfahren und es wird anhand einer konkreten Aufgabe skizziert, inwiefern Verfahren zur Parameterschätzung aus unvollständigen Beobachtungsdaten im Rahmen der Photogrammetrie eingesetzt werden können.

## 11 Statistische Methoden in der Photogrammetrie

Im Zusammenhang mit photogrammetrischen Auswertungen spielen statistische Methoden und Parameterschätzverfahren eine wesentliche Rolle. Dies betrifft sowohl die klassische Photogrammetrie als auch moderne Verfahren der digitalen Bildverarbeitung. Im Hinblick auf das Ziel, bei photogrammetrischen Auswertungen Resultate zu erzielen, die von zufälligen Meßunsicherheiten und groben Fehlern möglichst wenig beeinträchtigt werden, werden die Zielgrößen einer Auswertung in der Regel mittels mehr oder weniger robusten<sup>1</sup> Schätzverfahren aus z.T. hochredundanten Beobachtungen bestimmt.

Im Bereich der *klassischen Photogrammetrie* sind in diesem Zusammenhang beispielhaft

---

<sup>1</sup>Unter der Robustheit eines Schätzverfahrens wird hier die Unanfälligkeit der Schätzung gegenüber Ausreißern in den Beobachtungen verstanden

die bekannten Orientierungs- und Kalibrierungsverfahren<sup>2</sup> zu nennen, also

- der Räumlicher Rückwärtsschnitt (RRS),
- der Räumlicher Vorwärtsschnitt (RVS),
- die Relative Orientierung,
- die Absolute Orientierung,
- die Bündelblockausgleichung im Rahmen einer Aerotriangulation und
- die Testfeldkalibrierung.

Bei jedem dieser Verfahren werden Parameterschätzungen durchgeführt, um jeweils eine Menge  $\mathcal{P}$  unbekannter Parameter<sup>3</sup> mittels einer Menge  $\mathcal{V}$  gemessener Größen unter Berücksichtigung der Meßunsicherheiten zu bestimmen. Zwar existieren zu einigen der Verfahren auch direkte Lösungen oder zumindest direkte Näherungslösungen, hierbei handelt es sich jedoch meist um Spezialfälle hinsichtlich der Anzahl der verarbeiteten Bilder oder Beobachtungen. Daher werden die direkten Lösungen in der überwiegenden Zahl der Anwendungen allenfalls zur Bestimmung von Näherungswerten für die unbekannt Parameter verwendet, um die endgültigen Schätzwerte in einer Ausgleichung zu ermitteln.

Im Rahmen der *digitalen Bildverarbeitung* auf Grauwertbildern finden statistische Methoden ebenfalls umfangreiche Anwendung. Hervorzuheben sind in diesem Zusammenhang insbesondere Verfahren zur *automatischen Bildinterpretation*, also Verfahren zur automatischen (Wieder-) Erkennung von Objekten in Grauwertbildern. Bei diesen Verfahren wird vor allem auf statistische Klassifikationsverfahren zurückgegriffen, die i.d.R. entsprechend dem sogenannten *dreistufigen Konzept der Bildinterpretation* auf unterschiedlichen Abstraktionsebenen zur Anwendung kommen (vgl.[FUCHS 1998]). Nach einer Vorverarbeitung des Bildes auf der *unteren Ebene* der Bildinterpretation<sup>4</sup> wird bei der sogenannten *Merkmalsextraktion* auf der *mittleren Ebene der Bildinterpretation* jedes Bildpixel mittels statistischer Klassifikationsverfahren nach seiner Zugehörigkeit zu einer Punktstruktur, einer linienhaften Struktur oder einer homogenen Fläche im Bild klassifiziert. Im Hinblick auf diese Klassifikation werden die Eigenschaften der einzelnen Strukturelemente statistisch modelliert, um für jedes Pixel mittels Hypothesentests eine Entscheidung bezüglich der Zugehörigkeit zur Merkmalsklasse „Punkt“, „Linie“ oder „Fläche“ treffen zu können (vgl. [FUCHS 1998]). Im Anschluß daran wird auf der *oberen Ebene der Bildinterpretation* den einzelnen Strukturelementen, bzw. Gruppen von Strukturelementen jeweils die Bezeichnung einer Objektklasse zugeordnet, wobei wieder statistische Klassifikationsverfahren zum Einsatz kommen. Aus einer Menge von im Objektraum möglicherweise vorkommender Objekte wird also das Objekt bestimmt, das im Bild mit größter Wahrscheinlichkeit vorkommt. Im Ergebnis werden geschätzte Bezeichnungen der im Bild dargestellten Objekte erhalten(vgl.[HORNEGGER 1996]).

Zusammenfassend kann man sagen, daß es aufgrund der hier angedeuteten Vielzahl von Anwendungen statistischer Methoden in der Photogrammetrie durchaus gerechtfertigt ist, sich im Rahmen der Photogrammetrie mit den theoretischen Grundlagen der Ausgleichungsrechnung und Statistik zu beschäftigen, zumal die klassischen Methoden der geodätischen Ausgleichungsrechnung bei bestimmten Aufgaben aus der Photogrammetrie

<sup>2</sup>Zur detaillierten Erläuterung der Verfahren wird auf die Standardliteratur der Photogrammetrie (z.B. [KRAUS 1994]) verwiesen

<sup>3</sup>Die Menge  $\mathcal{P}$  der unbekannt Parameter umfasst je nach Aufgabe Parameter der Inneren und Äußeren Orientierung der verwendeten Kameras sowie die Koordinaten von Neupunkten.

<sup>4</sup>Die Vorverarbeitung kann u.a. aus einer geometrischen Transformation (etwa einer Entzerrung), einer radiometrischen Korrektur (etwa einer Kontrastanpassung) und-/oder einer Glättung bestehen.

nicht direkt angewendet werden können, so daß sie an die konkrete Aufgabe angepasst, ergänzt oder erweitert werden müssen.

Im folgenden Abschnitt werden kurz die herkömmlichen Methoden der klassischen geodätischen Ausgleichsrechnung beschrieben und es wird anhand eines Beispiels gezeigt, daß die klassische Ausgleichsrechnung im Zusammenhang mit bestimmten Schätz- und Klassifikationsproblemen an ihre Grenzen stößt. Da insbesondere die Beschreibung der Schätzverfahren sehr allgemein gehalten ist, sei an dieser Stelle im Hinblick auf eine detaillierte Beschreibung auf [KOCH 1998] verwiesen.

## 12 Klassische Schätzverfahren und deren Grenzen

### 121 Prinzip der klassischen geodätischen Ausgleichung

Den klassischen geodätischen Parameterschätzverfahren ist gemeinsam, daß eine Menge  $\mathcal{P}$  unbekannter Parameter aufgrund einer Menge  $\mathcal{V}$  gemessener Größen (Beobachtungen) bestimmt wird, die mit den unbekanntem Parametern in einem bestimmten Zusammenhang stehen. Der Zusammenhang zwischen den Beobachtungen und den unbekanntem Parametern geht dabei aus einem mathematischen Modell hervor, von dem vorausgesetzt wird, das es aufgrund gewisser Gesetzmäßigkeiten, denen die Beobachtungen in Abhängigkeit von den unbekanntem Parametern genügen, a priori bekannt ist. Es wird also für jede Beobachtung vorausgesetzt, daß von vorneherein bekannt, in welchem mathematischen Zusammenhang sie mit den unbekanntem Parametern steht. Die Gleichungen, die den Zusammenhang zwischen den Beobachtungen und den unbekanntem Parametern beschreiben, werden *Beobachtungsgleichungen* genannt. Die Beobachtungen werden aufgrund der Meßunsicherheiten als stochastische Variable aufgefasst und im Rahmen der Ausgleichung gemäß ihrer Präzision gewichtet.

Im Hinblick auf die Abgrenzung der üblichen Schätzverfahren von den Schätzverfahren aus unvollständigen Daten kommt folgender Aussage eine wesentliche Bedeutung zu:

*Bei den Standardschätzverfahren aus vollständigen Beobachtungsdaten handelt es sich bei allen im Modell eingeführten Observablen um tatsächlich zu messende bzw. tatsächlich gemessene Größen.*

Üblicherweise wird ein linearer Zusammenhang zwischen den Beobachtungen und den Parametern hergestellt, was auf Parameterschätzungen in linearen Modellen führt<sup>5</sup>. Das bekannte *Gauß-Markoff-Modell* stellt ein solches lineares Modell dar, in dem die unbekanntem Parameter mittels unterschiedlicher Schätzmethoden ermittelt werden können. Meist wird hier die beste lineare erwartungstreue Schätzung der unbekanntem Parameter bestimmt. Im Gauß-Markoff-Modell stimmt diese Schätzung mit der Schätzung nach der Methode der kleinsten Quadrate und der Maximum-Likelihood-Methode im Falle normalverteilter Beobachtungen überein (vgl. [KOCH 1998]).

---

<sup>5</sup>ggf. folgt der lineare Zusammenhang aus einem nichtlinearen Zusammenhang durch Taylor-Entwicklung mit Näherungswerten für die unbekanntem Parameter und Abbruch nach dem linearen Term

## 122 Grenzen der klassischen geodätischen Ausgleichsrechnung

Die klassische geodätische Ausgleichsrechnung stößt u.a dann an ihre Grenzen, wenn der Zusammenhang zwischen den Beobachtungen und den unbekanntem Parametern nicht eindeutig angegeben werden kann, d.h. wenn das der Ausgleichung zugrundeliegende Beobachtungsmodell a priori nicht eindeutig festliegt.

Um dies zu verdeutlichen, sei das *Beispiel* betrachtet, bei dem die Abszissen  $u_i$  und Ordinaten  $y_i$  von Punkten zweier in einem ebenen  $uy$ -Koordinatensystem dargestellten Geraden  $\mathcal{K}_1$  und  $\mathcal{K}_2$  gemessen worden seien, um hieraus die Parameter *Steigung* und *Nullablage* der beiden Geraden zu bestimmen (vgl. Abbildung 122). Dabei sei es versäumt worden, festzuhalten, welche Beobachtungen zu welcher der beiden Geraden gehören, so daß diese Informationen für die Ausgleichung fehlen. Offensichtlich ist es nicht möglich, aus den verbleibenden Informationen die

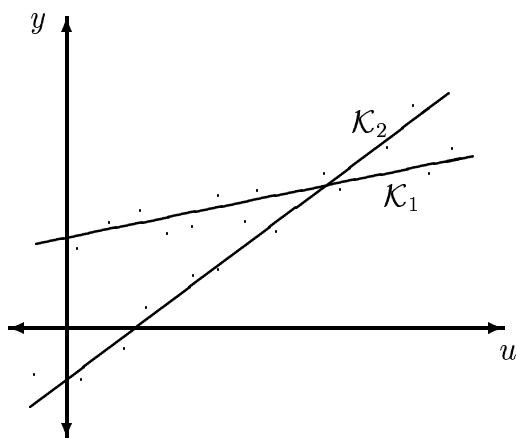


Abbildung 1:

Parameter beider Geraden mittels herkömmlicher geodätischer Schätzmethoden zu bestimmen, da (zumindest ohne Visualisierung der Daten) keine der Beobachtungen eindeutig einer der beiden Geraden zugeordnet werden kann. Es ist zwar bekannt, daß jede Beobachtung entweder zur Gerade  $\mathcal{K}_1$  oder zur Gerade  $\mathcal{K}_2$  gehört, aber nicht zu welcher von beiden. Somit ist der Zusammenhang zwischen den unbekanntem Parametern der Ausgleichung und den Beobachtungen a priori nicht eindeutig angebbar und damit eine klassische Ausgleichsrechnung nicht durchführbar.

Zur Lösung des Problems bedarf es der Anwendung eines Verfahrens zur sogenannten *Parameterschätzung aus unvollständigen Beobachtungsdaten*, wie sie im nächsten Abschnitt beschrieben ist. Der *Expectation-Maximization Algorithmus* stellt ein solches Verfahren dar; mit ihm ist es möglich, zum einen die einzelnen Beobachtungen hinsichtlich ihrer Zugehörigkeit zu einer der beiden Geraden zu klassifizieren und dabei gleichzeitig Schätzwerte für die Geradenparameter zu bestimmen (vgl. hierzu Kapitel 4)

Dem gegebenen Beispiel scheint zunächst der Bezug zur Photogrammetrie zu fehlen; es lässt sich allerdings auf Aufgaben der digitalen Bildverarbeitung verallgemeinern, bei denen bestimmte Kenngrößen (d.h. unbekannte Parameter) mehrerer im Bild erscheinender Objekte vollautomatisch bestimmt werden sollen, wobei die Zusammenhänge zwischen den hierzu herangezogenen Beobachtungen (z.B. den Ergebnissen einer Merkmalsextraktion) und den gesuchten Parametern bis auf die Zuordnung der Beobachtungen zu den in der Bildszene dargestellten Objekten bekannt ist. Da sich -wie beschrieben- solche Aufgaben mit den herkömmlichen Methoden der geodätischen Ausgleichsrechnung nicht

ohne weiteres lösen lassen, braucht man Verfahren zur Parameterschätzung aus unvollständigen Beobachtungsdaten. Das Prinzip der Parameterschätzung aus unvollständigen Daten soll daher im nächsten Abschnitt näher beschrieben werden.

## 13 Parameterschätzung aus unvollständigen Daten

Im Folgenden wird der Unterschied zwischen Parameterschätzungen aus unvollständigen Beobachtungsdaten und den herkömmlichen Schätzverfahren erläutert.

### Abgrenzung der Schätzung aus unvollständigen Beobachtungsdaten

Wie in Abschnitt 12 erläutert, wird bei den klassischen geodätischen Parameterschätzverfahren vorausgesetzt, daß als Beobachtungen nur *tatsächlich gemessene* Werte verwendet werden. Es kann u.U. allerdings auch von Vorteil sein, in die Ausgleichung sogenannte „*fehlende Beobachtungen*“ einzuführen. Fehlende Beobachtungen sind Zufallsparameter<sup>6</sup>, die wie *Beobachtungen* in eine Ausgleichung eingeführt werden, die aber tatsächlich *nicht* beobachtet worden sind. Hierbei kann es sich um Größen handeln, die entweder prinzipiell durch Messungen nicht zugänglich sind oder um meßbare Größen, die im konkreten Fall aus bestimmten Gründen nicht gemessen worden sind. Bei fehlenden Beobachtung handelt es sich also um Variablen einer Parameterschätzung, die aufgrund einer *Überparametrisierung des Modells* entstehen, also durch Einführung zusätzlicher Parameter in eine Ausgleichung, die durch bestehende Gesetzmäßigkeiten zwischen den tatsächlichen Beobachtungen und den eigentlichen Modellparametern nicht erklärbar sind.

Eine solche Überparametrisierung erscheint zunächst nicht plausibel, da es als wenig sinnvoll erscheint, nicht beobachtete Größen trotzdem formal als Beobachtungen in die Ausgleichung einzubringen. Bei näherer Betrachtung sieht man jedoch ein, daß die Einführung fehlender Beobachtungen insbesondere dann Sinn macht, wenn bestimmte Größen des mathematischen Modells einer Ausgleichung zwar nicht direkt meßbar sind, aber bekannt ist, daß zwischen ihnen und den durch Messungen zugänglichen Beobachtungen Korrelationen bestehen und wie groß diese sind. Würden diese unzugänglichen Beobachtungen nicht mit in die Ausgleichung einbezogen, so bliebe auch das Wissen um diese Wechselbeziehungen ungenutzt. Werden im Gegensatz dazu durch eine Überparametrisierung des Schätzproblems Zufallsparameter für die nicht zugänglichen Beobachtungen eingeführt, so kann die Kenntnis um die Korrelationen zwischen den versteckten und tatsächlichen Beobachtungen in der Ausgleichung verwertet werden.

Ein anschauliches, aus dem Bereich der Photogrammetrie stammendes Beispiel für Parameterschätzungen mit fehlenden Beobachtungen stellt die Klassifikation dreidimensionaler Objekte aus (nur) zweidimensionalen Bildern im Rahmen einer Einzelbildauswertung dar (vgl. [HORNEGGER 1996]): Ein dreidimensionales Objekt des Objektraumes wird im Bild nur zweidimensional dargestellt, die Tiefeninformation geht also bei der Abbildung verloren. Trotz dieses Informationsverlustes bei der Abbildung soll aber das räumliche Objekt aus dem ebenen Bild erkannt werden. Es liegt also nahe, im Rahmen einer Überparametrisierung die fehlenden Tiefeninformationen als fehlende Beobachtungen mit in die Parameterschätzung für die Bestimmung der Objektbezeichnung einzuführen.

---

<sup>6</sup>Unter Zufallsparametern werden hier unbekannte Parameter einer Ausgleichung verstanden, die nicht als feste Größen, sondern als Zufallsvariablen aufgefasst werden.

Es stellt sich also die Frage nach Parameterschätzverfahren, die fehlende Beobachtungen zulassen, also auf unvollständigen Beobachtungsdaten arbeiten. Der in dieser Arbeit untersuchte EM-Algorithmus stellt ein solches Schätzverfahren dar. Neben diesem Algorithmus existieren weitere Verfahren, wie z.B. den *DA (Data Augmentation Algorithms)* oder den *PMDA (Poor Man's Data Augmentation Algorithm)*. In Bezug auf diese Verfahren sei hier auf die Literatur [TANNER 1993] verwiesen.

## 14 Zusammenfassung / Überblick über die Arbeit

In dieser Einleitung wurde anhand von einiger Beispiele aufgezeigt, daß in der Photogrammetrie Methoden der Ausgleichsrechnung und Statistik eine sehr starke Rolle spielen, so daß es nahe liegt, sich auch im Rahmen der Photogrammetrie mit diesen Methoden auseinanderzusetzen. Außerdem wurden die herkömmlichen, auf vollständigen Beobachtungsdaten arbeitenden Verfahren der klassischen geodätischen Ausgleichsrechnung kurz skizziert und anhand eines Beispiels deren Grenzen aufgezeigt. Es wurde angedeutet, das in bestimmten Fällen, in denen eine klassische Ausgleichsrechnung nicht möglich ist, eine Parameterschätzung aus unvollständigen Beobachtungsdaten zum Ziel führen kann. Die in den Abschnitten 12) und 13 gegebenen Beispiele zeigen, daß die Parameterschätzung aus unvollständigen Daten für die Photogrammetrie durchaus relevant ist.

Im Folgenden wird der EM-Algorithmus als ein iteratives Verfahren zur Schätzung von Parametern aus unvollständigen Daten vorgestellt. Nach einer allgemeinen Definition dieses Algorithmus erfolgt eine ausführliche Auseinandersetzung mit seinen Eigenschaften. Hier werden insbesondere die Bedingungen betrachtet, unter denen die EM-Iterationen konvergieren, und die Geschwindigkeit des Algorithmus. Anschließend wird der EM-Algorithmus zur Lösung einer konkreten Aufgabe in einem Programm realisiert. Den Abschluss der Arbeit bildet eine Zusammenfassung.

# Kapitel 2

## Der EM-Algorithmus

Der in dieser Arbeit diskutierte Ansatz zur Parameterschätzung aus unvollständigen Daten ist der EM-Algorithmus. Er soll in diesem Abschnitt in seiner allgemeinen Form definiert werden. Im Hinblick auf die Notation sind hierzu einige Vorbemerkungen notwendig.

### 21 Begriffe und Notation

In der Einleitung wurde im Zusammenhang mit Parameterschätzverfahren abstrakt von einer Menge  $\mathcal{P}$  unbekannter Parameter gesprochen, die mit Hilfe einer Menge  $\mathcal{V}$  von Beobachtungen bestimmt werden. Im Folgenden werden nun die  $u$  unbekannt Parameter  $\beta_1, \beta_2, \dots, \beta_u$  einer Parameterschätzung in dem  $u \times 1$  Vektor

$$\boldsymbol{\beta} = [\beta_1, \beta_2, \dots, \beta_u]^T \quad (210.1)$$

zusammengefasst. Die unbekannt Parameter werden zunächst als fest vorausgesetzt, sie stellen also keine Zufallsvariablen dar. Der *Parameterraum*  $\mathcal{B}$  ist eine Untermenge des  $u$ -dimensionalen euklidischen Vektorraumes  $\mathbb{R}^u$  und umfasst alle im Rahmen des Schätzproblems theoretisch möglichen numerischen Realisierungen des Vektors  $\boldsymbol{\beta}$ . In dem  $n_y \times 1$  Zufallsvektor

$$\mathbf{y} = [y_1, y_2, \dots, y_{n_y}]^T \quad (210.2)$$

sind die  $n_y$  *tatsächlichen* Beobachtungen, also die durch Messungen zugänglichen Beobachtungen  $y_1, y_2, \dots, y_{n_y}$  enthalten. Die  $n_x$  durch Messungen nicht zugänglichen Beobachtungen  $x_1, x_2, \dots, x_{n_x}$  werden in dem  $n_x \times 1$  Zufallsvektor

$$\mathbf{x} = [x_1, x_2, \dots, x_{n_x}]^T \quad (210.3)$$

zusammengefasst. Die Mengen aller numerischen Realisierungen von  $\mathbf{x}$  bzw.  $\mathbf{y}$  bilden die *Stichprobenräume*  $\mathcal{X}$  bzw.  $\mathcal{Y}$ .

Der  $n_z \times 1$  Zufallsvektor

$$\mathbf{z} = [z_1, z_2, \dots, z_{n_z}]^T \quad (210.4)$$

ist der sogenannte *vollständige* Beobachtungsvektor. Hierbei handelt es sich *nicht* zwangsläufig um die einfache Zusammenfassung der Vektoren  $\mathbf{x}$  und  $\mathbf{y}$  in einem gemeinsamen

Vektor. Der Vektor  $\mathbf{z}$  umfaßt zwar stets alle fehlenden Beobachtungen aus  $\mathbf{x}$ , muß aber im allgemeinen nicht unbedingt auch tatsächliche Beobachtungen aus  $\mathbf{y}$  enthalten. Anhand des folgenden einfachen Beispiels soll die Struktur des Vektors  $\mathbf{z}$  verdeutlicht werden:

Betrachtet werde eine Parameterschätzung mit den (gemessenen) Beobachtungen  $y_i, i \in \{1 \dots n_y\}$ .

Fall 1: Es soll angenommen werden, daß die Beobachtung  $y_1$  der Summe zweier Beobachtungen  $x_1$  und  $x_2$  entspricht, die aber direkt nicht meßbar sind, also  $y_1 = x_1 + x_2$ . Der Vektor  $\mathbf{y}$  der tatsächlichen Beobachtungen, der Vektor  $\mathbf{x}$  der fehlenden Beobachtungen und der Vektor  $\mathbf{z}$  der vollständigen Beobachtungen haben dann die Form

$$\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \dots \\ y_{n_y} \end{bmatrix}, \quad \mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \quad \text{und} \quad \mathbf{z} = \begin{bmatrix} x_1 \\ x_2 \\ y_1 \\ y_2 \\ \dots \\ y_{n_y} \end{bmatrix}$$

Der vollständige Beobachtungsvektor enthält also neben den fehlenden Beobachtungen alle tatsächlichen Beobachtungen.

Fall 2: Werden sämtliche tatsächlichen Beobachtungen als Summe zweier versteckter Beobachtungen betrachtet, also  $y_i = x_{2i} + x_{2i-1}$ , so ergibt sich

$$\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \dots \\ y_{n_y} \end{bmatrix}, \quad \mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \dots \\ x_{2n_y} \end{bmatrix} \quad \text{und} \quad \mathbf{z} = \begin{bmatrix} x_1 \\ x_2 \\ \dots \\ x_{2n_y} \end{bmatrix} = \mathbf{x},$$

der vollständige Beobachtungsvektor enthält also außer den versteckten Beobachtungen keine weiteren Elemente.

Die Gesamtheit aller numerischer Realisierungen des Vektors  $\mathbf{z}$  bildet den Stichprobenraum  $\mathcal{Z}$  der vollständigen Beobachtungen.

Allgemein kann davon ausgegangen werden, daß der tatsächliche Beobachtungsvektor  $\mathbf{y}$  und der vollständige Beobachtungsvektor  $\mathbf{z}$  in einem mathematischen Zusammenhang stehen, da beide funktional von den selben unbekanntem Parametern abhängen. Es wird daher vorausgesetzt, daß der Stichprobenraum  $\mathcal{Y}$  durch eine Abbildung

$$B : \mathcal{Z} \rightarrow \mathcal{Y}$$

aus dem Stichprobenraum  $\mathcal{Z}$  hervorgeht, so daß sich  $\mathbf{y}$  durch Abbildung

$$B : \mathbf{z} \mapsto \mathbf{y} = B(\mathbf{z}) \tag{210.5}$$

aus  $\mathbf{z}$  ergibt<sup>1</sup>.

Im Fall 1 des oben gegebenen Beispiels ist die Abbildung  $B$  beispielsweise mit

$$B : \mathbf{z} \mapsto \mathbf{y} = \begin{bmatrix} \mathbf{0} \\ \mathbf{I} \end{bmatrix} \cdot \mathbf{z}$$

---

<sup>1</sup>Man kann sich vorstellen, daß sich die unbekanntem Parameter mit Hilfe der Elemente des Vektors  $\mathbf{z}$  bestimmen ließen, wenn  $\mathbf{z}$  vollständig beobachtet werden könnte. Die so ermittelten Parameter bestimmen unmittelbar die Elemente des tatsächlichen Beobachtungsvektors  $\mathbf{y}$ . Somit kann  $\mathbf{y}$  als Abbild von  $\mathbf{z}$  verstanden werden.



gegeben, worin  $\mathbf{0}$  eine Nullmatrix der Dimension  $2 \times n_y$  und  $\mathbf{I}$  eine  $n_y \times n_y$  Einheitsmatrix bedeuten. Im Fall 2 dieses Beispiels gilt

$$B : \mathbf{z} \mapsto \mathbf{y} = \begin{bmatrix} 1 & 1 & 0 & 0 & \dots & 0 & 0 \\ 0 & 0 & 1 & 1 & \dots & 0 & 0 \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ 0 & 0 & 0 & 0 & \dots & 1 & 1 \end{bmatrix} \cdot \mathbf{z}.$$

Die Abbildung  $B$  ist i.a. nicht eindeutig, d.h. die Gleichung  $\mathbf{y} = B(\mathbf{z})$  ist bei gegebenem  $\mathbf{z}$  in der Regel für mehrere Vektoren  $\mathbf{z}$  erfüllt. Daher wird die Menge

$$\mathcal{Z}(\mathbf{y}) = \{\mathbf{z} \mid \mathbf{y} = B(\mathbf{z})\} \quad (210.6)$$

eingeführt. Sie besteht aus allen Vektoren  $\mathbf{z}$ , die für ein gegebenes  $\mathbf{y}$  die Gleichung  $\mathbf{y} = B(\mathbf{z})$  erfüllen.

In Fall 1 des Beispiels besteht  $\mathcal{Z}(\mathbf{y})$  beispielsweise aus allen Vektoren  $\mathbf{z}$ , die die Bedingung  $z_2 = y_1 - z_1$  erfüllen.

Nachdem nun die in der Diplomarbeit verwendeten Bezeichnungen und die Notation erläutert wurden, wird kurz die Maximum-Likelihood-Methode beschrieben. Einerseits kann die Maximum-Likelihood-Methode für sich allein bereits zur Schätzung von Parametern aus unvollständigen Beobachtungen verwendet werden, andererseits stellt sie den Ausgangspunkt für den EM-Algorithmus dar.

## 22 Maximum-Likelihood-Methode

Im Zusammenhang mit dem Expectation-Maximization-Algorithmus spielt der Maximum-Likelihood-Schätzer eine wesentliche Rolle. So wird z.B. innerhalb der einzelnen Iterationsschritte der Erwartungswert der logarithmierten Likelihoodfunktion berechnet und maximiert (vgl. Abschnitt 232). Dies führt im Optimalfall auf eine Maximum-Likelihood-Schätzung der unbekannt Parameter. Wie in Abschnitt 222 gezeigt wird, eignet sich außerdem die Maximum-Likelihood-Methode als solche bereits zur Parameterschätzung aus unvollständigen Daten, wenn als Dichtefunktion eine von den tatsächlichen Beobachtungen  $\mathbf{y}$  abhängige Randdichte verwendet wird. Daher wird die Maximum-Likelihood-Schätzung in diesem Abschnitt kurz erläutert. Im nächsten Kapitel wird dann der EM-Algorithmus definiert.

### 221 Definition

Die Maximum-Likelihood-Methode zählt zu den klassischen Methoden der Parameterschätzung. Im Unterschied zu anderen herkömmlichen Verfahren, wie z.B. der besten linearen erwartungstreuen Schätzung oder der Methode der kleinsten Quadrate kann eine Parameterschätzung nach der Maximum-Likelihood-Methode nur durchgeführt werden, wenn die parametrische Dichtefunktion der Beobachtungen bekannt ist. Liegen normalverteilte<sup>2</sup> Beobachtungen vor, so führt die Maximum-Likelihood-Schätzung im Gauß-

<sup>2</sup>Aufgrund der Gültigkeit des Zentralen Grenzwertsatzes der mathematischen Statistik kann in einer Vielzahl von Anwendungen von normalverteilten Beobachtungen ausgegangen werden.

Markoff-Modell auf die gleichen Schätzwerte wie die Methode der kleinsten Quadrate und die beste lineare erwartungstreue Schätzung.(vgl. [KOCH 1998])

Die Definition der Maximum-Likelihood-Schätzung vorbereitend wird zunächst die sogenannte *Likelihoodfunktion* definiert<sup>3</sup>.

**Definition:** [LIKELIHOODFUNKTION]

Der Zufallsvektor  $\mathbf{z}$  der Beobachtungen besitze die von den unbekanntem, festen Parametern abhängige Dichte  $f(\boldsymbol{\beta})$ , dann ist die Likelihoodfunktion  $l(\mathbf{z} | \boldsymbol{\beta})$  definiert durch

$$l(\mathbf{z} | \boldsymbol{\beta}) = f(\boldsymbol{\beta}) \quad (221.1)$$

Die Bestimmung der Schätzwerte für die unbekanntem Parameter erfolgt durch Maximierung der Likelihoodfunktion. D.h. es wird der Parametervektor  $\boldsymbol{\beta}$  bestimmt, für den die Dichtefunktion  $l(\mathbf{z} | \boldsymbol{\beta})$  der Beobachtungen  $\mathbf{z}$  maximal wird. Dies folgt aus der Definition

**Definition:** [MAXIMUM-LIKELIHOOD-METHODE]

Der Zufallsvektor  $\mathbf{z}$  der Beobachtungen besitze die Likelihoodfunktion  $l(\mathbf{z} | \boldsymbol{\beta})$ , dann bezeichnet man als Maximum-Likelihood-Methode die Schätzung  $\bar{\boldsymbol{\beta}}$  der festen Parameter  $\boldsymbol{\beta}$ , die maximale Werte für  $l(\mathbf{z} | \boldsymbol{\beta})$  liefert, also

$$\bar{\boldsymbol{\beta}} = \arg \sup_{\boldsymbol{\beta}} l(\mathbf{z} | \boldsymbol{\beta}) \quad (221.2)$$

Durch Variation von  $\boldsymbol{\beta}$  ergibt sich also als Schätzung  $\bar{\boldsymbol{\beta}}$  der Vektor, der zusammen mit dem Beobachtungsvektor  $\mathbf{z}$  die maximale Dichte erzielt. Umgangssprachlich kann man sagen, es werden als Schätzwerte die Parameter gewählt, die gemäß der Dichtefunktion am besten zu den Beobachtungen passen.

Anstatt die Likelihoodfunktion  $l(\mathbf{z} | \boldsymbol{\beta})$  selbst zu maximieren wird meist die logarithmierte Likelihoodfunktion

$$L(\boldsymbol{\beta}) := L(\mathbf{z} | \boldsymbol{\beta}) := \log l(\mathbf{z} | \boldsymbol{\beta}) \quad (221.3)$$

maximiert. Dies vereinfacht im allgemeinen die Berechnungen und führt im Zusammenhang mit der EDV auch auf numerisch sicherere Ergebnisse. Aufgrund der strengen Monotonie und damit der Umkehrbarkeit der Logarithmusfunktion ist die Maximierung der logarithmierten Likelihoodfunktion gleichbedeutend mit der Maximierung der Likelihoodfunktion selbst. In diesem Zusammenhang ist auch der aus der Informationstheorie stammende Begriff der *Information* von Interesse. Unter der Information  $I(\mathbf{z})$ , die mit der Realisierung  $\mathbf{z}$  eines Zufallsvektors verbunden ist, wird der „Grad der Überraschung“ verstanden, die das Auftreten des Vektors  $\mathbf{z}$  auslöst (vgl. [FÖRSTNER 1991]). Mathematisch ist die Information  $I(\mathbf{z})$  eines Zufallsvektors  $\mathbf{z}$  definiert als negativer Logarithmus seiner Dichte (vgl. [FÖRSTNER 1991]):

$$I(\mathbf{z}) = \log \frac{1}{f(\mathbf{z})} = -\log f(\mathbf{z}). \quad (221.4)$$

Entsprechend gilt

$$I(\mathbf{z} | \boldsymbol{\beta}) = -\log l(\mathbf{z} | \boldsymbol{\beta}) = -L(\mathbf{z} | \boldsymbol{\beta}) \quad (221.5)$$

Somit entspricht die Maximierung der logarithmierten Likelihoodfunktion der Minimierung der zum Vektor  $\mathbf{z}$  gehörenden Information. Für die geschätzten Parameter wird also der Grad der Überraschung im Hinblick auf die Beobachtungen  $\mathbf{z}$  minimal.

<sup>3</sup>Die Definitionen dieses Abschnittes entsprechen denen in in [KOCH 1998]

## 222 Schätzung aus unvollständigen Beobachtungsdaten mittels Maximum-Likelihood-Methode

Wie zu Beginn des Abschnitts 22 angedeutet, kann die Maximum-Likelihood-Methode für sich bereits zur Schätzung unbekannter Parameter aus unvollständigen Beobachtungsdaten verwendet werden. Die Vorgehensweise hierzu soll im folgenden erläutert werden:

Gemäß der Definition (221.2) ist zur Bestimmung der unbekannt Parameter  $\beta$  mittels der Maximum-Likelihood-Methode die Likelihoodfunktion  $L(\mathbf{z} \mid \beta)$  zu maximieren. Im Falle eines Schätzproblems aus unvollständigen Daten steht man nun vor dem Problem, daß  $\mathbf{z}$  fehlende, also tatsächlich nicht gemessene Beobachtungen enthält. Somit kann, trotzdem die Dichtefunktion  $f(\mathbf{z} \mid \beta)$  der vollständigen Beobachtungen<sup>4</sup> in Abhängigkeit von den unbekannt Parameter als bekannt gilt, die Likelihoodfunktion  $L(\mathbf{z} \mid \beta)$  nicht unmittelbar ausgewertet und maximiert werden.

Der Weg zur Berechnung der Schätzwerte führt über die tatsächlichen Beobachtungen  $\mathbf{y}$ , die aus den vollständigen Beobachtungen durch die Abbildung  $\mathbf{y} = B(\mathbf{z})$  hervorgehen: Die von den unbekannt Parameter abhängige Dichtefunktion  $g(\mathbf{y} \mid \beta)$  der tatsächlichen Beobachtungen  $\mathbf{y}$  ergibt sich aus der (bekannt) Dichtefunktion  $f(\mathbf{z} \mid \beta)$  der vollständigen Beobachtungen als verallgemeinerte<sup>5</sup> Randdichte (vgl. [DEMPSTER ET AL. 1968]) zu

$$g(\mathbf{y} \mid \beta) = \int \cdots \int_{\mathcal{Z}(\mathbf{y})} f(\mathbf{z} \mid \beta) d\mathbf{z}. \quad (222.1)$$

Der Integrationsbereich ist hierbei die Menge  $\mathcal{Z}(\mathbf{y})$ , es wird also über sämtliche Vektoren  $\mathbf{z}$  integriert, die durch die Abbildung (210.5) auf  $\mathbf{y}$  abgebildet werden.

Die Randdichte  $g(\mathbf{y} \mid \beta)$  aus (222.1) kann nun als Likelihoodfunktion verwendet werden, so daß die unbekannt Parameter aufgrund der tatsächlichen Beobachtungen  $\mathbf{y}$  nach der Maximum-Likelihood-Methode bestimmt werden können.

Auf diese Weise kann prinzipiell mittels Maximum-Likelihood-Schätzung ein Schätzproblem mit unvollständigen Beobachtungsdaten gelöst werden. Allerdings werden in der Regel andere Verfahren zur Parameterschätzung aus unvollständigen Beobachtungsdaten vorgezogen, was im allgemeinen 2 Gründe hat: Zum einen ist die Bestimmung der Randdichte  $g(\mathbf{y} \mid \beta)$  mittels numerischer Integration sehr *rechenintensiv*. Zum anderen bleibt bei der beschriebenen Methode das *Wissen um die Spezialität der Dichtefunktion*  $f(\mathbf{z} \mid \beta)$  *ungenutzt*. Im allgemeinen wird eine ganze Menge  $\mathcal{F}$  von Dichtefunktionen  $f^g(\mathbf{z} \mid \beta)$  existieren, die durch Integration über  $\mathcal{Z}(\mathbf{y})$  zu der gleichen Randdichte führen wie  $f(\mathbf{z} \mid \beta)$ . Das vorliegende Wissen darüber, um welche dieser Dichtefunktionen es sich im konkreten Fall handelt, bleibt also bei der vorgestellten Methode unberücksichtigt. Es liegt hingegen die Vermutung nahe, daß sich dieses Wissen in der Parameterschätzung vorteilhaft einbringen lässt. Dies leistet der EM-Algorithmus. Er wird im folgenden nach einer Ableitung seiner Schlüsselgleichung allgemein definiert.

<sup>4</sup>Unter den vollständigen Beobachtungen werden im Rahmen dieser Arbeit die Elemente des vollständigen Beobachtungsvektors verstanden.

<sup>5</sup>Die in [KOCH 1998] auf Seite 98 angegebene Randverteilung stellt einen Spezialfall der Gleichung (222.1) dar. Sie ergibt sich aus (222.1), falls  $\mathbf{z} = (x_1, x_2, \dots, x_n)$  und  $\mathbf{y} = (x_{i+1}, \dots, x_n)$  gilt. Daher wird (222.1) hier als verallgemeinerte Randdichte bezeichnet.

## 23 Definition des (G)EM-Algorithmus

In diesem Abschnitt wird der *Expectation-Maximization-Algorithmus* (EM-Algorithmus) bzw. der allgemeinere *Generalized Expectation-Maximization-Algorithmus* (GEM-Algorithmus) allgemein definiert. Zunächst wird jedoch eine Gleichung abgeleitet, die als Schlüsselgleichung des (G)EM-Algorithmus<sup>6</sup> angesehen werden kann, da sie maßgeblich zu dessen Verständnis beiträgt.

### 231 Schlüsselgleichung des (G)EM-Algorithmus

Zur Ableitung der Schlüsselgleichung wird zunächst die bedingte Dichte  $f(\mathbf{z} \mid \mathbf{y}, \boldsymbol{\beta})$  des Vektors  $\mathbf{z}$  der vollständigen Beobachtungen eingeführt für den Fall, daß  $\mathbf{y}$  und  $\boldsymbol{\beta}$  gegeben sind. Nach [DEMPSTER ET AL. 1968], Gl. (2.5) gilt

$$f(\mathbf{z} \mid \mathbf{y}, \boldsymbol{\beta}) = \frac{f(\mathbf{z} \mid \boldsymbol{\beta})}{g(\mathbf{y} \mid \boldsymbol{\beta})}. \quad (231.1)$$

Somit folgt für die in Abschnitt 222 benutzte Likelihoodfunktion

$$g(\mathbf{y} \mid \boldsymbol{\beta}) = \frac{f(\mathbf{z} \mid \boldsymbol{\beta})}{f(\mathbf{z} \mid \mathbf{y}, \boldsymbol{\beta})}$$

und für die hierzu gehörende logarithmierte Likelihoodfunktion

$$L(\boldsymbol{\beta}) := L(\mathbf{y} \mid \boldsymbol{\beta}) = \log g(\mathbf{y} \mid \boldsymbol{\beta}) = \log f(\mathbf{z} \mid \boldsymbol{\beta}) - \log f(\mathbf{z} \mid \mathbf{y}, \boldsymbol{\beta}). \quad (231.2)$$

Ziel ist es nun, die logarithmierte Likelihoodfunktion (231.2) durch Variation von  $\boldsymbol{\beta}$  zu maximieren, um eine Maximum-Likelihood-Schätzung der unbekannt Parameter  $\boldsymbol{\beta}$  zu erhalten. Die Dichtefunktion  $f(\mathbf{z} \mid \boldsymbol{\beta})$  des vollständigen Beobachtungsvektors  $\mathbf{z}$  wird dabei bis auf die unbekannt Parameter als bekannt vorausgesetzt. Da die logarithmierte Likelihoodfunktion unter Ausnutzung dieser Kenntnis erfolgen soll, wird auf die Berechnung der Likelihoodfunktion nach (222.1) verzichtet und stattdessen die rechte Seite von (231.2) maximiert. Eine Hürde stellt in diesem Zusammenhang die Tatsache dar, daß der vollständige Beobachtungsvektor  $\mathbf{z}$  fehlende, also numerisch unbekannt Beobachtungen enthält, so daß weder  $f(\mathbf{z} \mid \boldsymbol{\beta})$  noch  $f(\mathbf{z} \mid \mathbf{y}, \boldsymbol{\beta})$  in Abhängigkeit von  $\boldsymbol{\beta}$  direkt ausgewertet werden können. Die Grundidee zur Überwindung dieser Hürde besteht darin, die fehlenden Beobachtungen aus  $\mathbf{z}$  durch deren Erwartungswerte<sup>7</sup> zu ersetzen, um dann die rechte Seite von (231.2) durch Variation von  $\boldsymbol{\beta}$  zu maximieren. Dies ist das Grundprinzip des EM-Algorithmus, der als iteratives Verfahren definiert ist (vgl. Abschnitt 232).

Zunächst wird vorausgesetzt, daß mit  $\hat{\boldsymbol{\beta}}^{(i)}$  vorläufige Schätzwerte für die unbekannt Parameter vorliegen. Mit Hilfe der Gleichung

$$L(\mathbf{y} \mid \hat{\boldsymbol{\beta}}^{(i+1)}) = \log g(\mathbf{y} \mid \hat{\boldsymbol{\beta}}^{(i+1)}) = \log f(\mathbf{z} \mid \hat{\boldsymbol{\beta}}^{(i+1)}) - \log f(\mathbf{z} \mid \mathbf{y}, \hat{\boldsymbol{\beta}}^{(i+1)}) \quad (231.3)$$

sollen nun neue Schätzwerte  $\hat{\boldsymbol{\beta}}^{(i+1)}$  für  $\boldsymbol{\beta}$  berechnet werden, wozu die logarithmierte Likelihoodfunktion  $L(\mathbf{y} \mid \hat{\boldsymbol{\beta}}^{(i+1)})$  als Zufallsvariable aufgefasst wird. Es kann der Erwartungswert dieser Zufallsvariable über  $\mathcal{Z}(\mathbf{y})$  berechnet werden, wobei zur Berechnung der

<sup>6</sup>Soweit in dieser Arbeit der Ausdruck *(G)EM-Algorithmus* verwendet wird, beziehen sich die entsprechenden Aussagen sowohl auf den EM-Algorithmus als auch auf den GEM-Algorithmus

<sup>7</sup>Die Berechnung dieser Erwartungswerte erfolgt iterativ und stützt sich jeweils auf vorher ermittelte Schätzwerte für die unbekannt Parameter. Das Verfahren wird noch ausführlich erläutert.

benötigten Dichtefunktionen  $f(\mathbf{z} \mid \mathbf{y}, \boldsymbol{\beta})$  die vorläufigen Schätzwerte  $\hat{\boldsymbol{\beta}}^{(i)}$  herangezogen werden. Aufgrund der Linearität des Erwartungswertoperators erhält man zunächst

$$E[L(\mathbf{y} \mid \hat{\boldsymbol{\beta}}^{(i+1)}) \mid \mathbf{y}, \hat{\boldsymbol{\beta}}^{(i)}] = E[\log f(\mathbf{z} \mid \hat{\boldsymbol{\beta}}^{(i+1)}) \mid \mathbf{y}, \hat{\boldsymbol{\beta}}^{(i)}] \\ - E[\log f(\mathbf{z} \mid \mathbf{y}, \hat{\boldsymbol{\beta}}^{(i+1)}) \mid \mathbf{y}, \hat{\boldsymbol{\beta}}^{(i)}]. \quad (231.4)$$

Hierin gilt für die linke Seite

$$E[L(\mathbf{y} \mid \hat{\boldsymbol{\beta}}^{(i+1)}) \mid \mathbf{y}, \hat{\boldsymbol{\beta}}^{(i)}] = E[\log g(\mathbf{y} \mid \hat{\boldsymbol{\beta}}^{(i+1)}) \mid \mathbf{y}, \hat{\boldsymbol{\beta}}^{(i)}] \\ = \int \cdots \int_{\mathbf{z}(\mathbf{y})} \log g(\mathbf{y} \mid \hat{\boldsymbol{\beta}}^{(i+1)}) f(\mathbf{z} \mid \mathbf{y}, \hat{\boldsymbol{\beta}}^{(i)}) d\mathbf{z} \\ = \log g(\mathbf{y} \mid \hat{\boldsymbol{\beta}}^{(i+1)}) \int \cdots \int_{\mathbf{z}(\mathbf{y})} f(\mathbf{z} \mid \mathbf{y}, \hat{\boldsymbol{\beta}}^{(i)}) d\mathbf{z} \quad (231.5) \\ = \log g(\mathbf{y} \mid \hat{\boldsymbol{\beta}}^{(i+1)}) \\ = L(\mathbf{y} \mid \hat{\boldsymbol{\beta}}^{(i+1)}).$$

Durch Einsetzen von (231.5) in (231.4) erhält man die

**Schlüsselgleichung des EM-Algorithmus:**

$$L(\mathbf{y} \mid \hat{\boldsymbol{\beta}}^{(i+1)}) = E[\log f(\mathbf{z} \mid \hat{\boldsymbol{\beta}}^{(i+1)}) \mid \mathbf{y}, \hat{\boldsymbol{\beta}}^{(i)}] \\ - E[\log f(\mathbf{z} \mid \mathbf{y}, \hat{\boldsymbol{\beta}}^{(i+1)}) \mid \mathbf{y}, \hat{\boldsymbol{\beta}}^{(i)}] \quad (231.6)$$

bzw.

$$L(\mathbf{y} \mid \hat{\boldsymbol{\beta}}^{(i+1)}) = Q(\hat{\boldsymbol{\beta}}^{(i+1)} \mid \hat{\boldsymbol{\beta}}^{(i)}) - H(\hat{\boldsymbol{\beta}}^{(i+1)} \mid \hat{\boldsymbol{\beta}}^{(i)}). \quad (231.7)$$

Hierin bedeuten  $Q(\hat{\boldsymbol{\beta}}^{(i+1)} \mid \hat{\boldsymbol{\beta}}^{(i)})$  und  $H(\hat{\boldsymbol{\beta}}^{(i+1)} \mid \hat{\boldsymbol{\beta}}^{(i)})$  Abkürzungen für die beiden Erwartungswerte in (231.6). Die Größe  $Q(\hat{\boldsymbol{\beta}}^{(i+1)} \mid \hat{\boldsymbol{\beta}}^{(i)})$  wird in diesem Zusammenhang als *Kullback-Leibler-Statistik* und  $-H(\hat{\boldsymbol{\beta}}^{(i+1)} \mid \hat{\boldsymbol{\beta}}^{(i)})$  als *Entropie*<sup>8</sup> bezeichnet. (Vgl. hierzu u.a. [HORNEGGER 1996],[FÖRSTNER 1991])

Explizit gilt für die **Kullback-Leibler-Statistik**

$$Q(\hat{\boldsymbol{\beta}}^{(i+1)} \mid \hat{\boldsymbol{\beta}}^{(i)}) := E[\log f(\mathbf{z} \mid \hat{\boldsymbol{\beta}}^{(i+1)}) \mid \mathbf{y}, \hat{\boldsymbol{\beta}}^{(i)}] \\ = \int \cdots \int_{\mathbf{z}(\mathbf{y})} f(\mathbf{z} \mid \hat{\boldsymbol{\beta}}^{(i+1)}) f(\mathbf{z} \mid \mathbf{y}, \hat{\boldsymbol{\beta}}^{(i)}) d\mathbf{z} \quad (231.8)$$

und für die **Entropie**

$$H(\hat{\boldsymbol{\beta}}^{(i+1)} \mid \hat{\boldsymbol{\beta}}^{(i)}) := E[\log f(\mathbf{z} \mid \mathbf{y}, \hat{\boldsymbol{\beta}}^{(i+1)}) \mid \mathbf{y}, \hat{\boldsymbol{\beta}}^{(i)}] \\ = \int \cdots \int_{\mathbf{z}(\mathbf{y})} \log f(\mathbf{z} \mid \mathbf{y}, \hat{\boldsymbol{\beta}}^{(i+1)}) f(\mathbf{z} \mid \mathbf{y}, \hat{\boldsymbol{\beta}}^{(i)}) d\mathbf{z}, \quad (231.9)$$

<sup>8</sup>Der Begriff der Entropie stammt aus der Informationstheorie. Unter der Entropie versteht man den Mittleren Wert der zu übermittelnden Information. (vgl. [FÖRSTNER 1991])

wobei  $f(\mathbf{z} | \mathbf{y}, \hat{\boldsymbol{\beta}}^{(i+1)})$  mit (231.1) berechnet werden kann<sup>9</sup>.

Bei gegebenem Vektor  $\hat{\boldsymbol{\beta}}^{(i+1)}$  lassen sich die beiden Erwartungswerte der rechten Seite von (231.6) bzw. (231.7) mittels der Integrale (231.8) und (231.9) und damit die logarithmierte Likelihoodfunktion  $L(\mathbf{y} | \hat{\boldsymbol{\beta}}^{(i+1)})$  (ggf. durch numerische Integration) berechnen.

Gleichung (231.6) bzw. (231.7) kann als *Schlüsselgleichung* des EM - Algorithmus angesehen werden, da sie die Grundlage für das nachfolgend vorgestellte iterative Verfahren zur Bestimmung der unbekannt Parameter bereitstellt. Aus dem beschriebenen Verfahren folgt unmittelbar der EM-Algorithmus.

### Iteratives Verfahren:

- 1.) Es wird ein Startvektor  $\hat{\boldsymbol{\beta}}^{(0)}$  als vorläufige Schätzung für die unbekannt Parameter gewählt. Die Wahl erfolgt mehr oder weniger willkürlich.
- 2.) Ausgehend von der vorläufigen Schätzung  $\hat{\boldsymbol{\beta}}^{(i)}$  der unbekannt Parameter wird (231.6) in Verbindung mit (231.8) und (231.9) dazu verwendet, die logarithmierte Likelihoodfunktion  $L(\mathbf{y} | \boldsymbol{\beta}^{(i+1)})$  für alternative Parameter  $\boldsymbol{\beta}^{(i+1)}$  zu berechnen.<sup>10</sup> Es werden diejenigen Parameter  $\hat{\boldsymbol{\beta}}^{(i+1)}$  bestimmt, für die  $L(\mathbf{y} | \boldsymbol{\beta}^{(i+1)})$  maximal wird. Diese Parameter stellen neue (verbesserte) Schätzwerte für die unbekannt Parameter dar.
- 3.) Falls die Abweichung  $\|\hat{\boldsymbol{\beta}}^{(i+1)} - \hat{\boldsymbol{\beta}}^{(i)}\|$  der neuen Schätzung  $\hat{\boldsymbol{\beta}}^{(i+1)}$  von der vorherigen Schätzung  $\hat{\boldsymbol{\beta}}^{(i)}$  einen Schwellwert  $\epsilon$  überschreitet, wird das Verfahren mit den verbesserten Schätzwerten bei 1.) fortgesetzt.
- 4.) Die endgültigen Schätzwerte  $\hat{\boldsymbol{\beta}}^*$  werden ausgegeben.

Das iterative Vorgehen ist notwendig, da in jeder Iteration die erhaltenen vorläufigen Schätzwerte von den Schätzwerten des vorherigen Schrittes abhängen.

### Übergang auf den EM-Algorithmus:

Beim oben beschriebenen iterativen Verfahren wird in jeder Iteration das Maximum der logarithmierten Likelihoodfunktion entweder analytisch oder numerisch bestimmt. Wegen

$$L(\mathbf{y} | \hat{\boldsymbol{\beta}}^{(i+1)}) = Q(\hat{\boldsymbol{\beta}}^{(i+1)} | \hat{\boldsymbol{\beta}}^{(i)}) - H(\hat{\boldsymbol{\beta}}^{(i+1)} | \hat{\boldsymbol{\beta}}^{(i)}).$$

müssen dazu die Kullback-Leibler-Statistik  $Q(\hat{\boldsymbol{\beta}}^{(i+1)} | \hat{\boldsymbol{\beta}}^{(i)})$  und die negative Entropie  $H(\hat{\boldsymbol{\beta}}^{(i+1)} | \hat{\boldsymbol{\beta}}^{(i)})$  in Abhängigkeit von  $\hat{\boldsymbol{\beta}}^{(i+1)}$  analytisch oder numerisch berechnet werden. Wie aus (231.1) und den Erläuterungen zu (222.1) und (231.9) hervorgeht, führt dabei insbesondere die numerische Berechnung der Entropie  $-H$  zu erheblichem Rechenaufwand.

Zur Verringerung des Rechenaufwandes macht man sich im Hinblick auf den EM-Algorithmus folgende Eigenschaft der Entropie zunutze: Es kann gezeigt werden, daß im  $(i + 1)$ -ten Iterationsschritt die negative Entropie  $H(\hat{\boldsymbol{\beta}}^{(i+1)} | \hat{\boldsymbol{\beta}}^{(i)})$  für jeden beliebigen Vektor  $\boldsymbol{\beta}'$

<sup>9</sup>Wie weiter unten erläutert, wird im Rahmen des EM-Algorithmus bei der Auswertung von (231.6) bzw. (231.7) auf die Berechnung der Entropie verzichtet. Somit braucht die aufwendige Berechnung von  $f(\mathbf{z} | \mathbf{y}, \hat{\boldsymbol{\beta}}^{(i+1)})$  und die zugehörige Integration tatsächlich nicht zu erfolgen

<sup>10</sup>Die Berechnung der logarithmierten Likelihoodfunktion und die anschließende Maximierung kann in vielen Fällen analytisch erfolgen. Ist dies nicht möglich, so muß mit numerischen Verfahren gearbeitet werden, bei denen die logarithmierte Likelihoodfunktion (231.6) numerisch berechnet.

mindestens so groß wird wie die sich mit den „alten“ Parametern  $\hat{\beta}^{(i)}$  ergebende Größe  $H(\hat{\beta}^{(i)} | \hat{\beta}^{(i)})$ . Die entsprechende Aussage liefert der

**Satz:** *Es seien  $\hat{\beta}^{(i)}$  und  $\beta'$  zwei Vektoren des Parameterraumes  $\mathcal{B}$ , dann gilt*

$$H(\beta' | \hat{\beta}^{(i)}) \leq H(\hat{\beta}^{(i)} | \hat{\beta}^{(i)}). \quad (231.10)$$

*Identität herrscht genau dann, wenn überall in  $\mathbf{z}$  gilt*

$$f(\mathbf{z} | \mathbf{y}, \hat{\beta}^{(i)}) = f(\mathbf{z} | \mathbf{y}, \beta') \quad (231.11)$$

BEWEIS: Betrachte die Differenz

$$\begin{aligned} H(\beta' | \hat{\beta}^{(i)}) - H(\hat{\beta}^{(i)} | \hat{\beta}^{(i)}) &= E[\log f(\mathbf{z} | \mathbf{y}, \beta') | \mathbf{y}, \hat{\beta}^{(i)}] - E[\log f(\mathbf{z} | \mathbf{y}, \hat{\beta}^{(i)}) | \mathbf{y}, \hat{\beta}^{(i)}] \\ &= E[\log f(\mathbf{z} | \mathbf{y}, \beta') - \log f(\mathbf{z} | \mathbf{y}, \hat{\beta}^{(i)}) | \mathbf{y}, \hat{\beta}^{(i)}] \\ &= E\left[\log \left( \frac{f(\mathbf{z} | \mathbf{y}, \beta')}{f(\mathbf{z} | \mathbf{y}, \hat{\beta}^{(i)})} \right) | \mathbf{y}, \hat{\beta}^{(i)}\right] \\ &= \int \cdots \int_{\mathbf{z}(\mathbf{y})} \log \left\{ \frac{f(\mathbf{z} | \mathbf{y}, \beta')}{f(\mathbf{z} | \mathbf{y}, \hat{\beta}^{(i)})} \right\} f(\mathbf{z} | \mathbf{y}, \hat{\beta}^{(i)}) d\mathbf{z} \end{aligned} \quad (231.12)$$

Durch Taylor-Entwicklung der Logarithmusfunktion  $\log x$  an der Stelle  $x_0 = 1$  ergibt sich bei Verwendung des Lagrange-Restgliedes vom Grad 2 (vgl. [HELFRICH II 1996])

$$\log(x) = (x - 1) - \underbrace{\frac{1}{2\xi^2} \cdot (x - 1)^2}_{\geq 0},$$

mit  $\xi$  zwischen 1 und  $x$ , also  $\xi \in [x; 1]$  falls  $x < 1$  bzw.  $\xi \in [1; x]$  falls  $x \geq 1$ . Daraus folgt  $\log(x) \leq x - 1$ , wobei Gleichheit genau dann gilt, wenn  $x = 1$ . Hiermit ergibt sich

$$\begin{aligned} H(\beta' | \hat{\beta}^{(i)}) - H(\hat{\beta}^{(i)} | \hat{\beta}^{(i)}) &\leq \int \cdots \int_{\mathbf{z}(\mathbf{y})} \left\{ \frac{f(\mathbf{z} | \mathbf{y}, \beta')}{f(\mathbf{z} | \mathbf{y}, \hat{\beta}^{(i)})} - 1 \right\} f(\mathbf{z} | \mathbf{y}, \hat{\beta}^{(i)}) d\mathbf{z} \\ &= \int \cdots \int_{\mathbf{z}(\mathbf{y})} \{f(\mathbf{z} | \mathbf{y}, \beta') - f(\mathbf{z} | \mathbf{y}, \hat{\beta}^{(i)})\} d\mathbf{z} \\ &= \int \cdots \int_{\mathbf{z}(\mathbf{y})} f(\mathbf{z} | \mathbf{y}, \beta') d\mathbf{z} - \int \cdots \int_{\mathbf{z}(\mathbf{y})} f(\mathbf{z} | \mathbf{y}, \hat{\beta}^{(i)}) d\mathbf{z} = 1 - 1 = 0, \end{aligned}$$

da nach (222.1) und (231.1) gilt

$$\int \cdots \int_{\mathbf{z}(\mathbf{y})} f(\mathbf{z} | \mathbf{y}, \hat{\beta}^{(\phi)}) d\mathbf{z} = \frac{1}{g(\mathbf{y} | \hat{\beta}^{(\phi)})} \int \cdots \int_{\mathbf{z}(\mathbf{y})} f(\mathbf{z} | \mathbf{y}, \hat{\beta}^{(\phi)}) d\mathbf{z} = 1 \quad \text{für } \phi = (i, i + 1).$$

Hieraus folgt schließlich die erste Aussage  $H(\beta' | \hat{\beta}^{(i)}) - H(\hat{\beta}^{(i)} | \hat{\beta}^{(i)}) \leq 0$ .

Entsprechend der im Zusammenhang mit der Taylorentwicklung der Logarithmusfunktion getroffenen Aussagen gilt das Gleichheitszeichen genau dann, wenn in (231.12) und damit überall in  $\mathbf{z}$  gilt  $f(\mathbf{z} | \mathbf{y}, \beta') = f(\mathbf{z} | \mathbf{y}, \hat{\beta}^{(i)})$ . Somit sind beide Aussagen des Satzes bewiesen.  $\square$

Aus der Eigenschaft  $H(\boldsymbol{\beta}' | \hat{\boldsymbol{\beta}}^{(i)}) \leq H(\hat{\boldsymbol{\beta}}^{(i)} | \hat{\boldsymbol{\beta}}^{(i)})$  resultiert im Hinblick auf den EM-Algorithmus nun folgende Überlegung: Wird bei dem oben beschriebenen iterativen Verfahren nur die Kullback-Leibler-Statistik  $Q$  betrachtet und werden dabei im  $(i + 1)$ -ten Iterationsschritt Schätzwerte  $\hat{\boldsymbol{\beta}}^{(i+1)}$  gefunden, für die  $Q(\hat{\boldsymbol{\beta}}^{(i+1)} | \hat{\boldsymbol{\beta}}^{(i)}) \geq Q(\hat{\boldsymbol{\beta}}^{(i)} | \hat{\boldsymbol{\beta}}^{(i)})$  gilt, so ist damit wegen  $H \geq 0$  nach (231.9) und  $H(\hat{\boldsymbol{\beta}}^{(i+1)} | \hat{\boldsymbol{\beta}}^{(i)}) \leq H(\hat{\boldsymbol{\beta}}^{(i)} | \hat{\boldsymbol{\beta}}^{(i)})$  gleichzeitig eine Vergrößerung der logarithmierten Likelihoodfunktion

$$L(\mathbf{y} | \hat{\boldsymbol{\beta}}^{(i+1)}) = Q(\hat{\boldsymbol{\beta}}^{(i+1)} | \hat{\boldsymbol{\beta}}^{(i)}) - H(\hat{\boldsymbol{\beta}}^{(i+1)} | \hat{\boldsymbol{\beta}}^{(i)}) \geq L(\mathbf{y} | \hat{\boldsymbol{\beta}}^{(i)})$$

im Vergleich zu  $L(\mathbf{y} | \hat{\boldsymbol{\beta}}^{(i)})$  verbunden. Aus diesem Grund verzichtet man auf die (rechenintensive) Auswertung der negativen Entropie  $H$  und beschränkt sich bei den Iterationen des EM-Algorithmus auf die Maximierung der Kullback-Leibler-Statistik  $Q$ <sup>11</sup>.

Im folgenden wird der EM-Algorithmus formal definiert.

## 232 Definition des EM-Algorithmus

Wie im letzten Abschnitt bereits erläutert, besteht der EM-Algorithmus in der iterativen Berechnung und Maximierung der Kullback-Leibler-Statistik  $Q(\hat{\boldsymbol{\beta}}^{(i+1)} | \hat{\boldsymbol{\beta}}^{(i)})$ . Die Definition des EM-Algorithmus in seiner allgemeinen Form lautet (vgl. [DEMPSTER ET AL. 1968]):

**Definition:** [EM-ALGORITHMUS]

Für den Zufallsvektor  $\mathbf{z}$  der vollständigen Beobachtungen sei die von den unbekanntesten, festen Parametern  $\boldsymbol{\beta}$  abhängige Wahrscheinlichkeitsdichte  $f(\mathbf{z} | \boldsymbol{\beta})$  bekannt. In  $\mathbf{z}$  seien fehlende, also durch Messungen nicht direkt zugängliche Beobachtungen wie ggf. auch tatsächliche Beobachtungen enthalten. Die tatsächlichen Beobachtungen seien im Vektor  $\mathbf{y}$  und die unzugänglichen Beobachtungen im Vektor  $\mathbf{x}$  zusammengefaßt. Dann besteht der *Expectation Maximization - Algorithmus (EM-Algorithmus)* aus der iterativen Anwendung des *Expectation-Schrittes (E-Schritt)* und des *Maximization-Schrittes (M-Schritt)*:

*E-Schritt:* Berechne die Kullback-Leibler-Statistik

$$Q(\hat{\boldsymbol{\beta}}^{(i+1)} | \hat{\boldsymbol{\beta}}^{(i)}) = E[\log f(\mathbf{z} | \hat{\boldsymbol{\beta}}^{(i+1)}) | \mathbf{y}, \hat{\boldsymbol{\beta}}^{(i)}] \quad \text{nach (231.8)}$$

(232.1)

*M-Schritt:* Maximiere die Kullback-Leibler-Statistik, um die neue Schätzung  $\hat{\boldsymbol{\beta}}^{(i+1)}$  zu erhalten.

$$\hat{\boldsymbol{\beta}}^{(i+1)} := \arg \max_{\hat{\boldsymbol{\beta}}^{(i+1)} \in \mathcal{B}} Q(\hat{\boldsymbol{\beta}}^{(i+1)} | \hat{\boldsymbol{\beta}}^{(i)}) \quad (232.2)$$

Hierin bezeichnet  $\mathcal{B}$  den Parameterraum der unbekanntesten Parameter.

Beim EM-Algorithmus wird im M-Schritt die Maximierung der Kullback-Leibler-Statistik gefordert. Es muß also gelten

$$Q(\hat{\boldsymbol{\beta}}^{(i+1)} | \hat{\boldsymbol{\beta}}^{(i)}) \geq Q(\boldsymbol{\beta}' | \hat{\boldsymbol{\beta}}^{(i)}) \quad \text{für alle Paare } (\boldsymbol{\beta}', \hat{\boldsymbol{\beta}}^{(i)}) \in \mathcal{B} \times \mathcal{B}. \quad (232.3)$$

<sup>11</sup>Die Beschränkung auf die Maximierung der Kullback-Leibler-Statistik führt dazu, daß die Werte  $\hat{\boldsymbol{\beta}}^{(i+1)}$ , für die  $Q$  maximal wird, nicht unbedingt auch die logarithmierte Likelihoodfunktion maximieren. Daher sind für die Kullback-Leibler-Statistik und die logarithmierte Likelihoodfunktion differenzierte Konvergenzbetrachtungen erforderlich (vgl. Kap 3).



Fordert man hingegen nur, daß in der Iteration der neue Schätzwert  $\hat{\beta}^{(i+1)}$  eine *mindestens so große* Kullback-Leibler-Statistik liefert wie der vorherige Schätzwert  $\hat{\beta}^{(i)}$ , so erhält man den sogenannten *Generalized EM -Algorithmus (GEM)*.

**Definition:**[GEM-ALGORITHMUS]

Wird im M-Schritt lediglich gefordert, daß die im  $(i + 1)$ -ten Iterationsdurchgang mit der neuen Schätzung  $\hat{\beta}^{(i+1)}$  berechnete Kullback-Leibler-Statistik mindestens so groß ist wie die mit der alten Schätzung  $\hat{\beta}^{(i)}$  berechnete Kullback-Leibler-Statistik, also

$$Q(\hat{\beta}^{(i+1)} | \hat{\beta}^{(i)}) \geq Q(\hat{\beta}^{(i)} | \hat{\beta}^{(i)}) \quad (232.4)$$

so bezeichnet man den zugehörigen Algorithmus als *Generalized Expectation-Maximization-Algorithmus (GEM-Algorithmus)*.

Diese Definitionen stellen die Grundlage für die nachfolgenden Untersuchungen der Eigenschaften des EM-Algorithmus bzw. des GEM-Algorithmus dar. Der Ablauf der Parameterschätzung aus unvollständigen Beobachtungsdaten mittels des EM-Algorithmus ist in Abbildung 1 schematisch dargestellt.

### 233 Zusammenfassung

In diesem Kapitel wurde die Maximum-Likelihood-Methode kurz vorgestellt und erläutert, wie sich damit Parameterschätzungen aus unvollständigen Beobachtungsdaten durchführen lassen. Im Anschluß daran wurde die Schlüsselgleichung des EM-Algorithmus abgeleitet und das sich daraus ergebende iterative Verfahren zur Parameterschätzung aus unvollständigen Daten beschrieben. Es wurde aufgezeigt, wie durch aus diesem iterativen Verfahren der EM-Algorithmus entsteht. Abschließend wurden der EM-Algorithmus und der GEM-Algorithmus formal definiert und deren Unterschied wurde erläutert.

Im nächsten Kapitel werden die Eigenschaften des (G)EM-Algorithmus beschrieben. Für die konkrete Anwendung dieses Algorithmus ist dabei insbesondere die Frage zu stellen, unter welchen Bedingungen der Algorithmus wie schnell zu welchem Ergebnis führt.

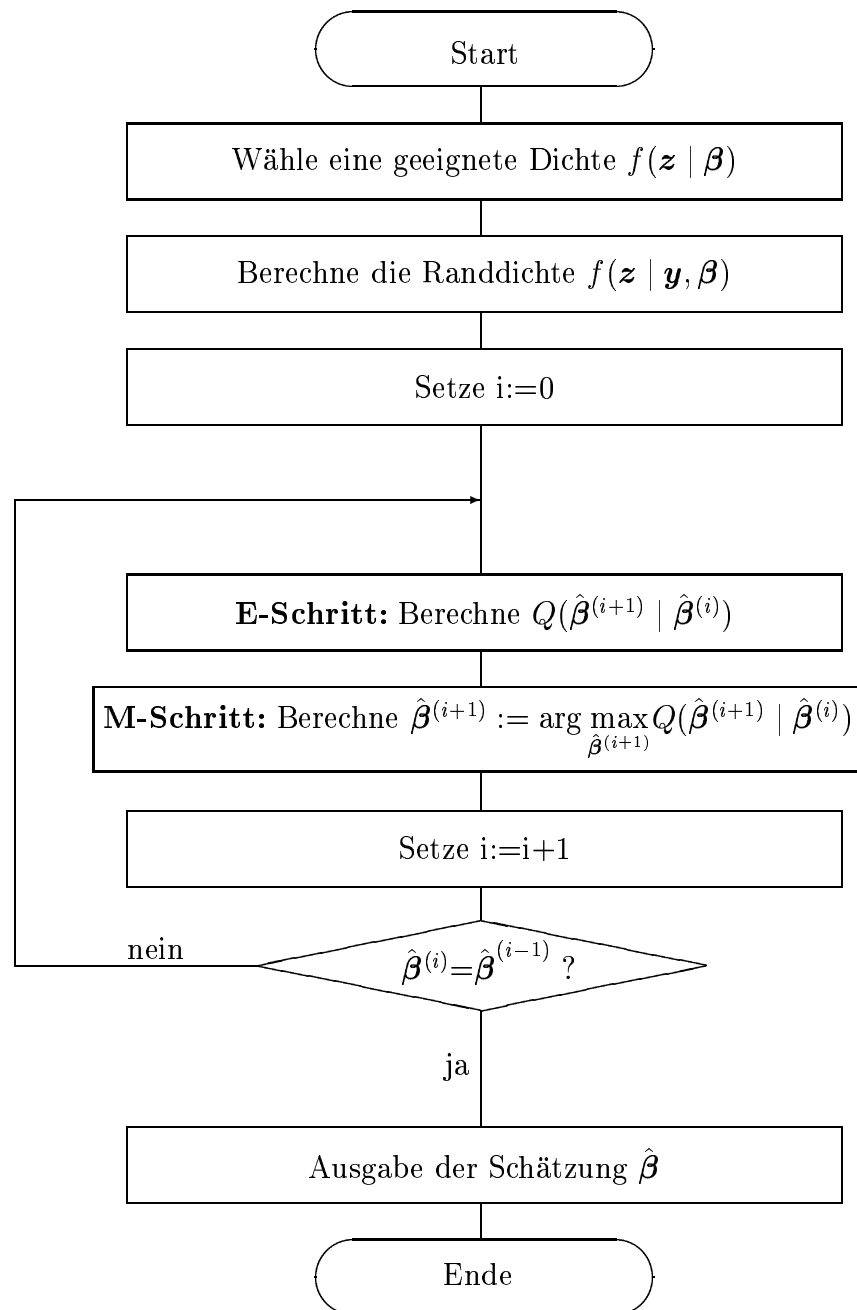


Abbildung 1: Ablaufschema des EM-Algorithmus

## Kapitel 3

# Eigenschaften des EM-Algorithmus

In diesem Kapitel werden die Eigenschaften des EM-Algorithmus beschrieben. Insbesondere ist hierbei von Interesse, unter welchen Bedingungen der Algorithmus konvergiert und ob er tatsächlich wie gewünscht die logarithmierte Likelihoodfunktion maximiert.

Wie sich aus den Definitionen (232.2) bzw. (232.4) ergibt, stellt der EM-Algorithmus einen Spezialfall des allgemeineren GEM-Algorithmus dar. Daher werden im Weiteren vor allem die Eigenschaften des GEM-Algorithmus wiedergegeben. Für den EM-Algorithmus gelten diese Eigenschaften dann entsprechend.

### 31 Vorbemerkungen

Im Rahmen des EM-Algorithmus werden in jeder Iteration neue (verbesserte) Schätzwerte  $\hat{\beta}^{(i+1)}$  für die unbekannt Parameter unter Zuhilfenahme der in der vorherigen Iteration erhaltenen Schätzwerte  $\hat{\beta}^{(i)}$  bestimmt. Es entsteht eine Folge von Parametervektoren

$$\{\hat{\beta}^{(i)}\}_{i \geq 0} = \hat{\beta}^{(0)}, \hat{\beta}^{(1)}, \hat{\beta}^{(2)}, \dots, \quad (310.1)$$

deren einzelne Folgenglieder jeweils vom vorherigen Folgenglied abhängen. Aufgrund dieser Abhängigkeit wird davon ausgegangen, daß der Parametervektor  $\hat{\beta}^{(i+1)}$  der  $(i+1)$ -ten Iteration formal aus dem Parametervektor  $\hat{\beta}^{(i)}$  der  $i$ -ten Iteration durch eine Abbildung  $A: \mathcal{B} \rightarrow \mathcal{B}$  hervorgeht; der Maximization-Schritt wird also beschrieben durch

$$A: \hat{\beta}^{(i)} \mapsto \hat{\beta}^{(i+1)} = A(\hat{\beta}^{(i)}). \quad (310.2)$$

Im Allgemeinen ergeben sich für verschiedene Startwerte  $\hat{\beta}^{(0)}$  in der  $i$ -ten Iteration unterschiedliche Schätzwerte  $\hat{\beta}^{(i)}$ . Die zu den Schätzungen  $\hat{\beta}^{(i)}$  gehörenden Werte  $L(\hat{\beta}^{(i)}) = L(\mathbf{y} | \hat{\beta}^{(i)})$  der logarithmierten Likelihoodfunktion bilden die Folge

$$\{L(\hat{\beta}^{(i)})\}_{i \geq 0} = L(\hat{\beta}^{(0)}), L(\hat{\beta}^{(1)}), L(\hat{\beta}^{(2)}), \dots. \quad (310.3)$$

Entsprechend der Definition (232.2) kommen beim EM-Algorithmus in jeder Iteration  $(i+1)$  als neue Schätzung  $\hat{\beta}^{(i+1)}$  alle Vektoren  $\beta$  in Frage, die die Kullback-Leibler-Statistik  $Q(\beta | \hat{\beta}^{(i)})$  in dieser Iteration maximieren; sie bilden die Lösungsmege  $N^{(i+1)}$  des M-Schrittes der  $(i+1)$ -ten Iteration. Wegen der Abhängigkeit der Kullback-Leibler-Statistik von der Schätzung  $\hat{\beta}^{(i)}$  hängt die Lösungsmenge der  $(i+1)$ -ten Iteration von der

Schätzung  $\hat{\beta}^{(i)}$  aus der  $i$ -ten Iteration ab. Diesem Umstand Rechnung tragend, wird die Abbildung

$$N : \mathcal{B} \rightarrow \mathcal{U} \subset \mathcal{B} \quad \text{bzw.} \quad (310.4)$$

$$N : \hat{\beta}^{(i)} \mapsto N(\hat{\beta}^{(i)}) = \{\beta : Q(\beta \mid \hat{\beta}^{(i)}) \geq Q(\bar{\beta} \mid \hat{\beta}^{(i)}) \forall \bar{\beta} \in \mathcal{B}\} \quad (310.5)$$

des Parameterraumes  $\mathcal{B}$  auf die Menge  $\mathcal{U}$  aller Teilmengen des Parameterraums eingeführt. Durch sie wird jeder Schätzung  $\hat{\beta}^{(i)}$  der  $i$ -ten Iteration die sich hieraus ergebende Lösungsmenge  $N^{(i+1)} = N(\hat{\beta}^{(i)})$  des M-Schrittes der  $(i+1)$ -ten Iteration zugeordnet<sup>1</sup>. Analog wird entsprechend der Definition (232.4) für den GEM-Algorithmus die Abbildung

$$M : \mathcal{B} \rightarrow \mathcal{U} \subset \mathcal{B} \quad \text{bzw.} \quad (310.6)$$

$$M : \hat{\beta}^{(i)} \mapsto M(\hat{\beta}^{(i)}) = \{\beta : Q(\beta \mid \hat{\beta}^{(i)}) \geq Q(\hat{\beta}^{(i)} \mid \hat{\beta}^{(i)})\} \quad (310.7)$$

des Schätzwertes  $\hat{\beta}^{(i)}$  auf die Teilmenge  $M(\hat{\beta}^{(i)})$  des Parameterraums eingeführt. Mit ihr ergibt sich die Lösungsmenge  $M^{(i+1)}$  des M-Schrittes in der  $(i+1)$ -ten EM-Iteration zu  $M^{(i+1)} = M(\hat{\beta}^{(i)})$ <sup>1</sup>.

Bei den Abbildungen  $N$  und  $M$  wird jeweils ein Punkt  $\hat{\beta}^{(i)}$  auf eine Menge  $N(\hat{\beta}^{(i)})$  bzw.  $M(\hat{\beta}^{(i)})$  abgebildet. Eine solche Abbildung bezeichnet man als eine *Punkt-Menge-Abbildung (PMA)*. Im Zusammenhang mit Punkt-Menge-Abbildungen wird häufig der Begriff der *Abgeschlossenheit* verwendet.

**Definition:** [ABGESCHLOSSENHEIT EINER PUNKT-MENGE-ABBILDUNG]

Man bezeichnet eine auf der Menge  $\mathcal{B}$  erklärte Punkt-Menge-Abbildung  $M : \mathcal{B} \rightarrow \mathcal{U}$  von Punkten  $\beta \in \mathcal{B}$  auf Mengen  $M(\beta) \in \mathcal{U}$  als abgeschlossen an der Stelle  $\beta^*$ , falls sie folgende Eigenschaft besitzt:

Ist  $\{\beta^{(k)}\}_{k \geq 0}$  mit  $\beta^{(k)} \in \mathcal{B}$  eine im Urbild  $\mathcal{B}$  definierte Folge mit dem Grenzwert  $\beta^*$ , d.h.  $\beta^{(k)} \rightarrow \beta^*$  und  $\{L^{(k)}\}_{k \geq 0}$  eine auf der Menge  $\mathcal{U}$  definierte Folge mit dem Grenzwert  $L^*$ , d.h.  $L^{(k)} \rightarrow L^*$ , deren Folgenglieder  $L^{(k)}$  jeweils in der Menge  $M(\beta^{(k)})$  liegen, also  $L^{(k)} \in M(\beta^{(k)})$ , dann folgt  $L^* \in M(\beta^*)$ . (310.8)

Die Abgeschlossenheitsbedingung für Punkt-Menge-Abbildungen ist vergleichbar mit der Stetigkeitsbedingung für Abbildungen, bei denen Punkte auf Punkte abgebildet werden (Vgl. [KERNER ET AL. 1995], S. 277).

In den folgenden Abschnitten wird bei der Beschreibung der Eigenschaften des (G)EM-Algorithmus auf die Vorbemerkungen dieses Abschnittes zurückgegriffen.

## 32 Eigenschaften des (G)EM-Algorithmus

Bereits aus den Erläuterungen zur Schlüsselgleichung des EM-Algorithmus in Abschnitt 231 und der Definition (232.4) des GEM-Algorithmus gehen zwei Eigenschaften des (G)EM-Algorithmus hervor:

- 1.) In der  $(i+1)$ -ten Iteration liefert die neue Schätzung  $\hat{\beta}^{(i+1)}$  eine mindestens gleich große Kullback-Leibler-Statistik wie die Schätzung  $\hat{\beta}^{(i)}$  der vorherigen Iteration  $i$ ,  $Q(\hat{\beta}^{(i+1)} \mid \hat{\beta}^{(i)}) \geq Q(\hat{\beta}^{(i)} \mid \hat{\beta}^{(i)})$ . Dies ergibt sich unmittelbar aus der Definition (232.4) des GEM-Algorithmus.

<sup>1</sup>Im folgenden ist also zu berücksichtigen, daß das Symbol  $N$  entweder im Zusammenhang mit der Lösungsmenge  $N^{(i+1)}$  eines M-Schrittes oder mit der Abbildung  $N(\hat{\beta}^{(i)})$  auftritt. Entsprechendes gilt für  $M$ .

- 2.) Die negative Entropie  $H$  wird in der  $(i+1)$ -ten Iteration mit der neuen Schätzung  $\hat{\beta}^{(i+1)}$  höchstens so groß wie mit der Schätzung  $\hat{\beta}^{(i)}$  der vorherigen Iteration  $i$ ,  
 $H(\hat{\beta}^{(i+1)} | \hat{\beta}^{(i)}) \leq H(\hat{\beta}^{(i)} | \hat{\beta}^{(i)})$ . (vgl. (231.11))

Über das aus diesen Eigenschaften resultierende Verhalten der logarithmierten Likelihoodfunktion  $L(\beta)$  während der EM-Iterationen wurde bisher noch keine explizite Aussage getroffen. Das Verhalten von  $L(\beta)$  ist aber besonders von Interesse, da das Ziel des EM-Algorithmus ja darin besteht, zu einer Maximum-Likelihood-Schätzung zu gelangen. Im nächsten Abschnitt wird daher das Verhalten der logarithmierten Likelihoodfunktion im Rahmen des GEM-Algorithmus beschrieben und es werden Bedingungen genannt, unter denen die logarithmierte Likelihoodfunktion zu einem lokalen Maximum bzw. zu einem stationären Punkt konvergiert. Im darauf folgenden Abschnitt erfolgen dann Betrachtungen zur Konvergenz der Parameterfolge  $\{\hat{\beta}^{(i)}\}_{i \geq 0}$ .

### 321 Zum Verhalten der Likelihoodfunktion während der (G)EM-Iterationen

Eine erste Aussage bezüglich des Verhaltens der logarithmierten Likelihoodfunktion in den (G)EM-Iterationen liefert der

**Satz:** *Im Rahmen des GEM-Algorithmus wächst die Folge  $L(\hat{\beta}^{(0)}), L(\hat{\beta}^{(1)}), \dots$  der in den Iterationen erhaltenen logarithmierten Likelihoodfunktionen  $L(\beta) = \log g(\mathbf{y} | \beta)$  monoton, d.h. es gilt für jede Iteration  $(i+1)$*

$$L(\hat{\beta}^{(i+1)}) \geq L(\hat{\beta}^{(i)}), \quad (321.1)$$

wobei das Gleichheitszeichen genau dann gültig ist, wenn sowohl

$$Q(\hat{\beta}^{(i+1)} | \hat{\beta}^{(i)}) = Q(\hat{\beta}^{(i)} | \hat{\beta}^{(i)}) \quad (321.2)$$

als auch

$$f(\mathbf{z} | \mathbf{y}, \hat{\beta}^{(i+1)}) = f(\mathbf{z} | \mathbf{y}, \hat{\beta}^{(i)}) \quad \text{für alle } \mathbf{z} \in \mathcal{Z} \quad (321.3)$$

gilt.

BEWEIS: Nach (231.7) gilt für die Differenz der logarithmierten Likelihoodfunktion zweier aufeinanderfolgender Iterationen  $i$  und  $(i+1)$

$$L(\hat{\beta}^{(i+1)}) - L(\hat{\beta}^{(i)}) = Q(\hat{\beta}^{(i+1)} | \hat{\beta}^{(i)}) - H(\hat{\beta}^{(i+1)} | \hat{\beta}^{(i)}) \\ - [Q(\hat{\beta}^{(i)} | \hat{\beta}^{(i-1)}) - H(\hat{\beta}^{(i)} | \hat{\beta}^{(i-1)})] \quad (321.4)$$

Zum Beweis des Satzes muß gezeigt werden, daß dieser Ausdruck nicht negativ ist und daß Gleichheit genau unter den genannten Bedingungen gilt. Hierzu wird zunächst die rechte Seite von (321.4) umgeformt. Es gilt

$$Q(\hat{\beta}^{(i)} | \hat{\beta}^{(i-1)}) - H(\hat{\beta}^{(i)} | \hat{\beta}^{(i-1)}) \\ = E[\log f(\mathbf{z} | \hat{\beta}^{(i)}) | \mathbf{y}, \hat{\beta}^{(i-1)}] - E[\log f(\mathbf{z} | \mathbf{y}, \hat{\beta}^{(i)}) | \mathbf{y}, \hat{\beta}^{(i-1)}] \\ = E[\log f(\mathbf{z} | \hat{\beta}^{(i)}) - \log f(\mathbf{z} | \mathbf{y}, \hat{\beta}^{(i)}) | \mathbf{y}, \hat{\beta}^{(i-1)}]$$

$$\begin{aligned}
&= E\left[\log \frac{f(\mathbf{z} \mid \hat{\boldsymbol{\beta}}^{(i)})}{f(\mathbf{z} \mid \mathbf{y}, \hat{\boldsymbol{\beta}}^{(i)})} \mid \mathbf{y}, \hat{\boldsymbol{\beta}}^{(i-1)}\right], \quad \text{mit } \frac{f(\mathbf{z} \mid \hat{\boldsymbol{\beta}}^{(i)})}{f(\mathbf{z} \mid \mathbf{y}, \hat{\boldsymbol{\beta}}^{(i)})} \stackrel{(231.1)}{=} g(\mathbf{y} \mid \hat{\boldsymbol{\beta}}^{(i)}) \\
&= E[\log g(\mathbf{y} \mid \hat{\boldsymbol{\beta}}^{(i)}) \mid \mathbf{y}, \hat{\boldsymbol{\beta}}^{(i-1)}] = \int \cdots \int_{\mathcal{Z}(\mathbf{y})} \log g(\mathbf{y} \mid \hat{\boldsymbol{\beta}}^{(i)}) f(\mathbf{z} \mid \mathbf{y}, \hat{\boldsymbol{\beta}}^{(i-1)}) d\mathbf{z} \\
&= \log g(\mathbf{y} \mid \hat{\boldsymbol{\beta}}^{(i)}) \int \cdots \int_{\mathcal{Z}(\mathbf{y})} f(\mathbf{z} \mid \mathbf{y}, \hat{\boldsymbol{\beta}}^{(i-1)}) d\mathbf{z} \\
&= \log g(\mathbf{y} \mid \hat{\boldsymbol{\beta}}^{(i)}).
\end{aligned}$$

Nach (231.3) und (231.7) gilt aber auch  $\log g(\mathbf{y} \mid \hat{\boldsymbol{\beta}}^{(i)}) = Q(\hat{\boldsymbol{\beta}}^{(i)} \mid \hat{\boldsymbol{\beta}}^{(i)}) - H(\hat{\boldsymbol{\beta}}^{(i)} \mid \hat{\boldsymbol{\beta}}^{(i)})$ , so daß für die rechte Seite von (321.4) folgt

$$L(\hat{\boldsymbol{\beta}}^{(i)}) = Q(\hat{\boldsymbol{\beta}}^{(i)} \mid \hat{\boldsymbol{\beta}}^{(i-1)}) - H(\hat{\boldsymbol{\beta}}^{(i)} \mid \hat{\boldsymbol{\beta}}^{(i-1)}) = Q(\hat{\boldsymbol{\beta}}^{(i)} \mid \hat{\boldsymbol{\beta}}^{(i)}) - H(\hat{\boldsymbol{\beta}}^{(i)} \mid \hat{\boldsymbol{\beta}}^{(i)}). \quad (321.5)$$

Zum Beweis des Satzes ist somit die Beziehung

$$Q(\hat{\boldsymbol{\beta}}^{(i+1)} \mid \hat{\boldsymbol{\beta}}^{(i)}) - H(\hat{\boldsymbol{\beta}}^{(i+1)} \mid \hat{\boldsymbol{\beta}}^{(i)}) - [Q(\hat{\boldsymbol{\beta}}^{(i)} \mid \hat{\boldsymbol{\beta}}^{(i)}) - H(\hat{\boldsymbol{\beta}}^{(i)} \mid \hat{\boldsymbol{\beta}}^{(i)})] \geq 0$$

zu verifizieren. Aus der Definition (232.4) des GEM-Algorithmus und der Beziehung  $H(\hat{\boldsymbol{\beta}}^{(i)} \mid \hat{\boldsymbol{\beta}}^{(i)}) \geq H(\hat{\boldsymbol{\beta}}^{(i+1)} \mid \hat{\boldsymbol{\beta}}^{(i)})$  aus (231.11) folgt

$$\begin{aligned}
&Q(\hat{\boldsymbol{\beta}}^{(i+1)} \mid \hat{\boldsymbol{\beta}}^{(i)}) - H(\hat{\boldsymbol{\beta}}^{(i+1)} \mid \hat{\boldsymbol{\beta}}^{(i)}) - Q(\hat{\boldsymbol{\beta}}^{(i)} \mid \hat{\boldsymbol{\beta}}^{(i)}) + H(\hat{\boldsymbol{\beta}}^{(i)} \mid \hat{\boldsymbol{\beta}}^{(i)}) \\
&= \underbrace{[Q(\hat{\boldsymbol{\beta}}^{(i+1)} \mid \hat{\boldsymbol{\beta}}^{(i)}) - Q(\hat{\boldsymbol{\beta}}^{(i)} \mid \hat{\boldsymbol{\beta}}^{(i)})]}_{=: R \geq 0} + \underbrace{[H(\hat{\boldsymbol{\beta}}^{(i)} \mid \hat{\boldsymbol{\beta}}^{(i)}) - H(\hat{\boldsymbol{\beta}}^{(i+1)} \mid \hat{\boldsymbol{\beta}}^{(i)})]}_{=: S \geq 0} \geq 0.
\end{aligned}$$

Hieraus folgt die erste Aussage.

Gleichheit besteht genau dann, wenn  $R=S=0$  gilt. Nach (231.11) ist dies genau dann der Fall, wenn (321.2) und (321.3) erfüllt sind. Damit ist der Satz vollständig bewiesen.  $\square$

Das monotone Wachstum der logarithmierten Likelihoodfunktion während der (G)EM-Iterationen stellt eine wichtige Eigenschaft des (G)EM-Algorithmus dar: Zusammen mit der unten geforderten Kompaktheit des Parameterraumes  $\mathcal{B}$  (vgl. (322.2)) gewährleistet sie die Konvergenz der logarithmierten Likelihoodfunktion im Rahmen des (G)EM-Algorithmus (vgl. die Betrachtungen zur Konvergenz der logarithmierten Likelihoodfunktion in Abschnitt 322).

Als erste, direkte Konsequenz der Monotonie der Folge  $\left\{L(\hat{\boldsymbol{\beta}}^{(i+1)})\right\}_{i \geq 0}$  ergibt sich der

**Satz:** Das Maximum  $L^{(max)}$  der logarithmierten Likelihoodfunktion stellt einen Fixpunkt (321.6) des GEM-Algorithmus dar.

Denn falls in der  $i$ -ten Iteration des (G)EM-Algorithmus mit  $\hat{\boldsymbol{\beta}}^{(i)}$  vorläufige Schätzwerte für die unbekannt Parameter gefunden werden, für die die logarithmierte Likelihoodfunktion  $L(\boldsymbol{\beta})$  maximal wird, also  $L(\hat{\boldsymbol{\beta}}^{(i)}) = L^{(max)} \geq L(\boldsymbol{\beta})$  für alle  $\boldsymbol{\beta} \in \mathcal{B}$ , so folgt aus (321.1) für die darauf folgende Schätzung  $\hat{\boldsymbol{\beta}}^{(i+1)} \in \mathcal{B}$  in der  $(i+1)$ -ten Iteration

- 1.)  $L(\hat{\boldsymbol{\beta}}^{(i+1)}) = L(\hat{\boldsymbol{\beta}}^{(i)}) = L^{(max)}$ , die logarithmierte Likelihoodfunktion nimmt also für die neuen Schätzwerte  $\hat{\boldsymbol{\beta}}^{(i+1)}$  wieder den gleichen Wert an wie im  $i$ -ten Schritt, und somit aus (321.2) und (321.3)
- 2.)  $Q(\hat{\boldsymbol{\beta}}^{(i+1)} \mid \hat{\boldsymbol{\beta}}^{(i)}) = Q(\hat{\boldsymbol{\beta}}^{(i)} \mid \hat{\boldsymbol{\beta}}^{(i)})$ , die Kullback-Leibler-Statistik nimmt in der  $(i+1)$ -ten Iteration den gleichen Wert an wie in der  $i$ -ten Iteration,

- 3.)  $f(\mathbf{z} \mid \mathbf{y}, \hat{\boldsymbol{\beta}}^{(i+1)}) = f(\mathbf{z} \mid \mathbf{y}, \hat{\boldsymbol{\beta}}^{(i)})$  für alle  $\mathbf{z}$ , die Dichtefunktion  $f(\mathbf{z} \mid \mathbf{y}, \hat{\boldsymbol{\beta}}^{(i+1)})$  ist überall in  $\mathbf{z}$  identisch mit der Dichtefunktion  $f(\mathbf{z} \mid \mathbf{y}, \hat{\boldsymbol{\beta}}^{(i)})$ , die mit Hilfe der Schätzwerte der  $i$ -ten Iteration berechnet wird.

Obwohl sich der Wert der logarithmierten Likelihoodfunktion in nachfolgenden Iterationen nicht mehr ändert, wenn sie in einer Iteration  $i$  einmal ihr Maximum erreicht hat, bedeutet dies jedoch *nicht*, daß damit auch ihr Argument, also die Schätzung  $\hat{\boldsymbol{\beta}}$  für die unbekannt Parameter festliegt. Denn nimmt die logarithmierte Likelihoodfunktion ihren Maximalwert für mehrere, voneinander verschiedene  $\boldsymbol{\beta} \in \mathcal{B}$  an, so besteht die Möglichkeit, daß die Schätzungen der nachfolgenden Iterationen zwischen diesen Vektoren alterniert, ohne daß die logarithmierte Likelihoodfunktion ihren Wert ändert. Setzt man jedoch voraus, daß die logarithmierte Likelihoodfunktion ihr Maximum nur in einem Punkt  $\boldsymbol{\beta}^*$  annimmt und daß dieses Maximum im  $i$ -ten EM-Iterationsschritt erreicht wird, also  $L(\hat{\boldsymbol{\beta}}^{(i)}) = L(\boldsymbol{\beta}^*) > L(\boldsymbol{\beta})$  für alle  $\boldsymbol{\beta} \in \mathcal{B} \setminus \{\hat{\boldsymbol{\beta}}^{(i+1)}\}$ , so folgt für die nächste Iteration des GEM-Algorithmus

$$\hat{\boldsymbol{\beta}}^{(i+1)} = \hat{\boldsymbol{\beta}}^{(i)} = \boldsymbol{\beta}^*.$$

(Dies läßt sich durch Widerspruch beweisen. Nimmt man an, daß  $\hat{\boldsymbol{\beta}}^{(i+1)} \neq \hat{\boldsymbol{\beta}}^{(i)}$ , so gilt  $L(\hat{\boldsymbol{\beta}}^{(i+1)}) > L^{(max)} > L(\hat{\boldsymbol{\beta}}^{(i+1)}) \in \mathcal{B}$ , was einen Widerspruch zur Annahme darstellt.)

Zusammenfassend kann also der folgende Satz festgehalten werden:

**Satz:** *Hat die logarithmierte Likelihoodfunktion  $L(\boldsymbol{\beta})$  innerhalb der (G)EM-Iterationen einmal ihren Maximalwert  $L^{(max)}$  erreicht, so ändert sich ihr Wert in den nachfolgenden Iterationen nicht mehr. Besitzt die logarithmierte Likelihoodfunktion nur an einer Stelle  $\boldsymbol{\beta}^* \in \mathcal{B}$  dieses Maximum, so liegt dann auch der endgültige Schätzwert  $\hat{\boldsymbol{\beta}} = \boldsymbol{\beta}^*$  fest, sonst liegt eine mehrdeutige Lösung des Optimierungsproblems vor.*

## 322 Betrachtungen zur Konvergenz der logarithmierten Likelihoodfunktion

Die Tatsache, daß die Folge  $\{L(\hat{\boldsymbol{\beta}}^{(i)})\}_{i \geq 0}$  der logarithmierten Likelihoodfunktion während der (G)EM-Iterationen monoton wächst und daß das globale Maximum  $L^*$  einen Fixpunkt des Algorithmus darstellt, läßt noch keine Aussage darüber zu, ob und ggf. unter welchen Bedingungen  $\{L(\hat{\boldsymbol{\beta}}^{(i)})\}_{i \geq 0}$  überhaupt konvergiert. Wird in einer Iteration das globale Maximum der logarithmierten Likelihoodfunktion erhalten, so folgt zwar aus den bisherigen Überlegungen, daß  $\{L(\hat{\boldsymbol{\beta}}^{(i)})\}_{i \geq 0}$  konvergiert. Es ist bisher jedoch nicht geklärt, ob das Maximum  $L^*$  in den (G)EM-Iterationen überhaupt erreicht wird.

Für die Betrachtungen zur Konvergenz der logarithmierten Likelihoodfunktion wird von folgenden Annahmen ausgegangen:

**Annahme 1:** Die Menge  $\mathcal{B}_{\beta_0} = \{\boldsymbol{\beta} \in \mathcal{B} : L(\boldsymbol{\beta}) \geq L(\beta_0)\}$  ist für jedes  $L(\beta_0) \in \mathbb{R}$  kompakt.

**Annahme 2:** Die Funktion  $L(\boldsymbol{\beta})$  ist stetig in  $\mathcal{B}$  und differenzierbar im Inneren von  $\mathcal{B}$ .

**Annahme 3:** Die logarithmierte Likelihoodfunktion nimmt ihren Maximalwert  $L^*$  im Innern des Parameterraums  $\mathcal{B}$  an, d.h. das globale Maximum ist gleichzeitig auch ein lokales Maximum der logarithmierten Likelihoodfunktion  $L(\boldsymbol{\beta})$ .

*Bemerkungen:*

Entsprechend der Ausführungen in Abschnitt 21 handelt es sich bei dem Parameterraum  $\mathcal{B}$  um eine Teilmenge des  $u$ -dimensionalen euklidischen Vektorraumes,  $\mathcal{B} \subset \mathbb{R}^u$ . Entsprechend den Ausführungen in [KERNER ET AL. 1995] ist eine solche Teilmenge des  $\mathbb{R}^u$  genau dann kompakt, wenn sie abgeschlossen<sup>3</sup> und beschränkt<sup>4</sup> ist. Mit Annahme 1 wird also davon ausgegangen, daß es sich bei der Lösungsmenge des Maximization-Schrittes in jeder Iteration des GEM-Algorithmus jeweils um eine abgeschlossene und beschränkte Menge handelt. Dies ist eine sehr schwache Forderung, so daß Annahme 1 nahezu immer richtig ist<sup>5</sup>. Auch die mit den Annahmen 2 und 3 verbundenen Forderungen nach Stetigkeit und Differenzierbarkeit der logarithmierten Likelihoodfunktion bzw. nach dem Auftreten des globalen Maximums der logarithmierten Likelihoodfunktion im Inneren des Parameterraumes, sind für viele Dichtefunktionen kontinuierlicher Zufallsvariable - insbesondere für die Normalverteilung und die Exponentialverteilung (vgl. [DEMPSTER ET AL. 1968]) - erfüllt. Somit sind die nachfolgend unter den Annahmen 1, 2 und 3 getroffenen Aussagen praktisch von sehr allgemeiner Gültigkeit.

Die beiden Annahmen (322.1) und (322.2) lassen nachstehende *Folgerung* zu:

Nach (322.1) entstammt jedes Glied der Folge  $\{\hat{\beta}^{(i)}\}_{i \geq 0}$  der beschränkten Menge  $\mathcal{B}_{\hat{\beta}^{(0)}}$ .

Somit ist  $\{\hat{\beta}^{(i)}\}_{i \geq 0}$  eine beschränkte Folge. Wegen der angenommenen Stetigkeit der logarithmierten Likelihoodfunktion  $L(\beta)$  nach (322.2) folgt hieraus, daß auch die Folge  $\{L(\hat{\beta}^{(i)})\}_{i \geq 0}$  beschränkt ist. Insgesamt handelt es sich also bei  $\{L(\hat{\beta}^{(i)})\}_{i \geq 0}$  um eine monoton wachsende und unter den (praktisch sehr häufig zutreffenden) Annahmen (322.1) und (322.2) um eine nach oben beschränkte Folge. Nach dem Satz von BOLZANO-WEIYERSTRASS (vgl. [KERNER ET AL. 1995], S. 35) ist jede solche monotone und beschränkte Folge konvergent, womit die Konvergenz von  $\{L(\hat{\beta}^{(i)})\}_{i \geq 0}$  folgt. Dies führt auf den

**Satz:** *Mit Ausnahme der Fälle, in denen die Annahmen (322.1) und (322.2) nicht zutreffen, ist beim (G)EM-Algorithmus die Folge  $\{L(\hat{\beta}^{(i)})\}_{i \geq 0}$  der Werte der logarithmierten Likelihoodfunktion, die sich mit den in den (G)EM-Iterationen erhaltenen Schätzungen  $\hat{\beta}^{(i)}$  für die unbekannt Parameter ergeben, konvergent. Der Grenzwert wird mit  $L^*$  bezeichnet.* (322.4)

Nachdem die Konvergenz der logarithmierten Likelihoodfunktion im Rahmen des (G)EM-Algorithmus nachgewiesen ist, stellt sich die Frage, ob es sich bei dem Grenzwert  $L^*$  um das gesuchte globale Maximum der logarithmierten Likelihoodfunktion handelt; die bisherigen Überlegungen lassen diesbezüglich noch keine Aussage zu.

Treffen die Annahmen 1, 2 und 3 zu, so handelt es sich bei dem Grenzwert  $L^*$  entweder

- um einen Wert, der zu einem lokalen Maximum der logarithmierten Likelihoodfunktion gehört, das mit ihrem globalen Maximum übereinstimmt, *oder*
- um einen Wert, der zu einem sonstigen lokalen Maximum der logarithmierten Likelihoodfunktion gehört *oder*
- um einen Wert, der zu einem stationären Punkt<sup>6</sup> gehört.

<sup>3</sup>Eine Teilmenge  $\mathcal{E}$  des  $\mathbb{R}^u$  heißt abgeschlossen, wenn ihr Komplement  $\bar{\mathcal{E}} = \mathbb{R}^u \setminus \mathcal{E}$  offen ist, d.h. wenn zu jedem Punkt  $\bar{e} \in \bar{\mathcal{E}}$  eine Kugel mit Mittelpunkt in  $\bar{e}$  existiert, die vollständig in  $\bar{\mathcal{E}}$  liegt.

<sup>4</sup>Eine Teilmenge des  $\mathbb{R}^u$  heißt beschränkt, wenn es eine Kugel mit endlichem Radius gibt, die die Menge vollständig einschließt.

<sup>5</sup>Zumindest dem Verfasser ist kein Fall bekannt, in dem Annahme 1 nicht erfüllt wäre.



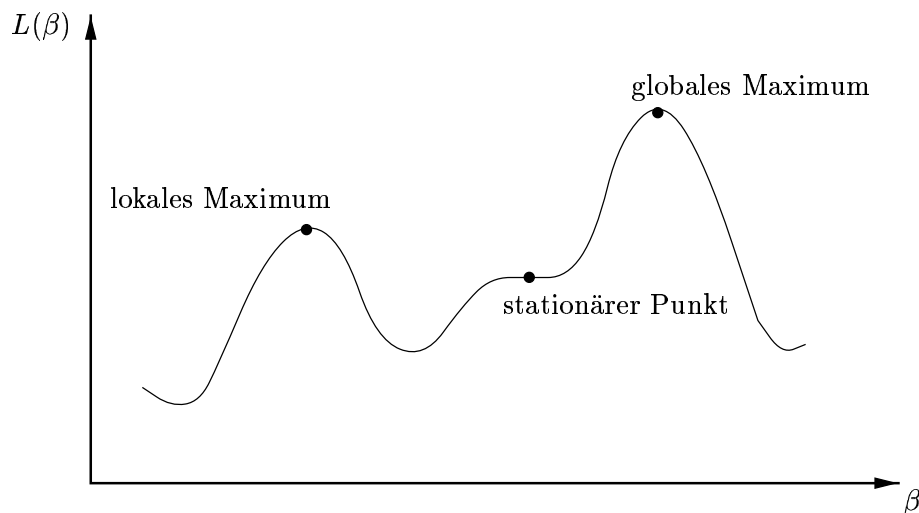


Abbildung 1: Theoretisch mögliche Konvergenzpunkte des EM-Algorithmus

Für den Fall, daß bei einer Parameterschätzung nur ein unbekannter Parameter  $\beta$  zu bestimmen ist, sind die möglichen Fälle in Abbildung 1 dargestellt.

Selbst wenn im M-Schritt des (G)EM-Algorithmus globale Optimierungsverfahren zur Maximierung der Kullback-Leibler-Statistik eingesetzt werden, ist es i.a. nicht sicher, daß der (G)EM-Algorithmus zum gesuchten globalen Maximum der logarithmierten Likelihoodfunktion führt. Das liegt daran, daß beim (G)EM-Algorithmus die negative Entropie  $H$  unberücksichtigt bleibt: Nach der Schlüsselgleichung des (G)EM-Algorithmus 231.7 gilt mit der Schätzung  $\hat{\beta}^{(i)}$  der  $i$ -ten (G)EM-Iteration die Beziehung  $L(\beta) = Q(\beta | \hat{\beta}^{(i)}) - H(\beta | \hat{\beta}^{(i)})$ . Im Rahmen des (G)EM-Algorithmus wird die hierin enthaltene negative Entropie  $H$  nicht weiter berücksichtigt, sondern im M-Schritt wird lediglich die Kullback-Leibler-Statistik maximiert. Diese Vernachlässigung der Größe  $H$  führt dazu, daß der (G)EM-Algorithmus auf eine Schätzung  $\hat{\beta}$  führt, die i.a. *nicht* zugleich die Größe  $L = Q + H$  global maximiert. So kommt es zur Konvergenz der logarithmierten Likelihoodfunktion gegen einen Wert, der zu einem (vom globalen Maximum verschiedenen) lokalen Maximum oder zu einem stationären Punkt gehört.

Im folgenden soll nun erklärt werden, unter welchen Bedingungen im Rahmen des (G)EM-Algorithmus mit der Konvergenz der logarithmierten Likelihoodwerte zu einem globalen Maximum, einem lokalen Maximum oder zu einem stationären Punkt (z.B. Sattelpunkt) zu rechnen ist.

Grundlegend für die hier angestellten Konvergenzbetrachtungen ist das folgende

### Theorem über globale Konvergenz

**Satz:** *Es sei  $\{\mathbf{a}_i\}_{i=0}^{\infty}$ ,  $\mathbf{a}_i \in \mathcal{A}$  eine Folge von Punkten mit  $\mathbf{a}_{i+1} \in F(\mathbf{a}_i)$ , wobei  $F$  eine auf der Menge  $\mathcal{A}$  definierte Punkt-Menge-Abbildung bezeichnet. Weiter sei mit  $\Gamma \subset \mathcal{A}$*

<sup>6</sup>Ein stationärer Punkt der logarithmierten Likelihoodfunktion ist ein Punkt, in dem der Gradient der logarithmierten Likelihoodfunktion zwar gleich dem Nullvektor ist (d.h. die Notwendige Bedingung für ein lokales Extremum ist in diesem Punkt erfüllt), bei dem es sich jedoch um kein lokales Extremum der logarithmierten Likelihoodfunktion handelt.

eine Lösungsmenge<sup>7</sup> gegeben und es gelte:

- 1.) Alle Punkte  $\mathbf{a}_i$  sind Elemente einer kompakten Menge  $\mathcal{S} \subset \mathcal{A}$
- 2.) Die Punkt-Menge-Abbildung  $F$  ist abgeschlossen in  $\mathcal{A} \setminus \Gamma$
- 3.) Es ist eine stetige Funktion  $\gamma(\mathbf{a})$  auf  $\mathcal{A}$  definiert, für die gilt
  - a) gehört  $\mathbf{a}$  nicht zur Lösungsmenge, also  $\mathbf{a} \notin \Gamma$ , dann gilt für alle Punkte  $\mathbf{b}$  des Bildes  $F(\mathbf{a})$  die Beziehung  $\gamma(\mathbf{b}) > \gamma(\mathbf{a})$
  - b) gehört  $\mathbf{a}$  zur Lösungsmenge, also  $\mathbf{a} \in \Gamma$ , dann gilt für alle Punkte  $\mathbf{b}$  des Bildes  $F(\mathbf{a})$  die Beziehung  $\gamma(\mathbf{b}) \geq \gamma(\mathbf{a})$

Dann liegen alle Konvergenzpunkte der Folge  $\{\mathbf{a}_i\}_0^\infty$  in der Lösungsmenge  $\Gamma$  und  $\{\gamma(\mathbf{a}_i)\}_{i \geq 0}$  konvergiert monoton mit  $\lim_{i \rightarrow \infty} \gamma(\mathbf{a}_i) = \gamma(\mathbf{a})$  für ein  $\mathbf{a} \in \Gamma$

Auf den Beweis dieses Satzes wird hier verzichtet, da er den Rahmen dieser Arbeit sprengen würde. Er kann in ([ZANGWILL 1969]) gefunden werden.

Nach dem Theorem über globale Konvergenz konvergiert also eine Folge  $\{\gamma(\mathbf{a}_i)\}_{i \geq 0}$ , deren Folgenglieder  $\gamma(\mathbf{a}_i)$  sich als Funktion der Glieder  $\mathbf{a}_i$  einer Folge  $\{\mathbf{a}_i\}_{i \geq 0}$  ergeben, wenn für  $\{\mathbf{a}_i\}$  eine Lösungsmenge  $\Gamma$  vorgegeben ist und die Bedingungen 1.), 2.), 3.a) und 3.b) erfüllt sind. Dies gilt für allgemeine Folgen  $\{\mathbf{a}_i\}$ , Punkt-Menge-Abbildungen  $F$ , Lösungsmengen  $\Gamma$  und Funktionen  $\gamma$ . Somit gilt dieses Theorem auch in dem Spezialfall des (G)EM-Algorithmus, in dem

- es sich bei  $\{\mathbf{a}_i\}_{i \geq 0}$  um die Folge der beim (G)EM-Algorithmus jeweils im M-Schritt erhaltenen Schätzungen  $\hat{\beta}^{(i)}$  für die Unbekannten Parameter handelt, also  $\mathbf{a}_i = \hat{\beta}^{(i)}$
- $\{\gamma(\mathbf{a}_i)\}_{i \geq 0}$  die entsprechende Folge der logarithmierten Likelihoodwerte darstellt, also  $\gamma(\mathbf{a}_i) = L(\hat{\beta}^{(i)})$
- $F$  die Punkt-Menge-Abbildung ist, die die Schätzung  $\hat{\beta}^{(i)}$  der  $i$ -ten (G)EM-Iteration auf die Lösungsmenge  $M^{(i+1)}$  bzw.  $N^{(i+1)}$  der  $(i+1)$ -ten Iteration abbildet, also  $F(\mathbf{a}_i) = M(\hat{\beta}^{(i)})$  beim GEM-Algorithmus und  $F(\mathbf{a}_i) = M(\hat{\beta}^{(i)})$  beim EM-Algorithmus.
- die Lösungsmenge entweder die Menge  $\Gamma^L$  Lokaler Maxima im Inneren von  $\mathcal{B}$ , die Menge  $\Gamma^S$  stationärer Punkte oder die Menge  $\Gamma^G$  globaler Maxima der logarithmierten Likelihoodfunktion ist.

Die Menge  $\mathcal{A}$  entspricht dann dem Parameterraum  $\mathcal{B}$ . Es können mit dem Theorem über globale Konvergenz also die Bedingungen formuliert werden, unter denen im Rahmen des (G)EM-Algorithmus die Folge  $\{L(\hat{\beta}^{(i)})\}_{i \geq 1}$  gegen ein globales oder lokales Maximum oder gegen einen stationären Punkt konvergiert.

Zunächst soll angegeben werden, unter welchen Bedingungen die aus den (G)EM-Iterationen resultierende Folge  $\{L(\hat{\beta}^{(i)})\}_{i \geq 0}$  gegen einen stationären Punkt der logarithmierten Likelihoodfunktion konvergiert. Als Lösungsmenge  $\Gamma$  wird hierzu die Menge  $\Gamma^S$  der stationären Punkte von  $L(\beta)$  im Inneren von  $\mathcal{B}$  gewählt. Wegen Annahme (322.1) ist die Bedingung 1 des Theorems über globale Konvergenz für den (G)EM-Algorithmus erfüllt. Genau so ist wegen des monotonen Wachstums der logarithmierten Likelihoodfunktion nach (321.1) für den (G)EM-Algorithmus die Bedingung (3b) erfüllt. Als Konsequenz aus dem Theorem über globale Konvergenz ergibt sich also unmittelbar der

**Satz:** Sei  $\{\hat{\beta}^{(i)}\}_{i \geq 0}$  eine Folge von Schätzungen der unbekannt Parameter aus den Iterationen des GEM-Algorithmus, so daß gilt  $\hat{\beta}^{(i+1)} \in M(\hat{\beta}^{(i)})$ . Weiter sei

<sup>7</sup>Die Lösungsmenge enthält diejenigen Punkte, mit denen alle bis auf endlich viele Folgenglieder der Folge  $\{\mathbf{a}_i\}_{i=0}^\infty$  zusammenfallen sollen. Sie wird aufgrund gewisser Vorgaben festgelegt.

1. die Punkt-Menge-Abbildung  $M$ , die die Schätzung  $\hat{\beta}^{(i)}$  der unbekannt Parameter der  $i$ -ten Iteration auf die Lösungsmenge  $M^{(i+1)} = M(\hat{\beta}^{(i)})$  des  $M$ -Schrittes der  $(i+1)$ -ten Iteration abbildet, für  $\hat{\beta}^{(i)} \in \mathcal{B} \setminus \Gamma^S$  abgeschlossen, wobei  $\Gamma^S$  die Menge aller stationären Punkte der logarithmierten Likelihoodfunktion im Innern von  $\mathcal{B}$  bezeichnet.
2. Es gelte  $L(\hat{\beta}^{(i+1)}) > L(\hat{\beta}^{(i)})$  für alle  $\hat{\beta}^{(i)} \notin \Gamma^S$

Dann sind alle Konvergenzpunkte von  $\{\hat{\beta}^{(i)}\}$  stationäre Punkte von  $L$  und  $L(\hat{\beta}^{(i)})$  konvergiert monoton mit dem Grenzwert  $\lim L(\hat{\beta}^{(i)}) = L^* = L(\beta^*)$  für ein  $\beta^* \in \Gamma^S$  (322.6). Eine entsprechende Aussage in Bezug auf die Konvergenz zu lokalen Maxima erhält man, in dem man mit der Lösungsmenge  $\Gamma$  die Menge  $\Gamma^L$  der lokalen Maxima identifiziert. Man erhält den

**Satz:** Sei  $\{\hat{\beta}^{(i)}\}$  eine Folge von Schätzungen der unbekannt Parameter aus den Iterationen des GEM-Algorithmus, mit  $\hat{\beta}^{(i+1)} \in M(\hat{\beta}^{(i)})$ . Weiter sei

1. die Punkt-Menge-Abbildung  $M$ , die die Schätzung  $\hat{\beta}^{(i)}$  der unbekannt Parameter der  $i$ -ten Iteration auf die Lösungsmenge  $M^{(i+1)} = M(\hat{\beta}^{(i)})$  des  $M$ -Schrittes der  $(i+1)$ -ten Iteration abbildet, für  $\hat{\beta}^{(i)} \in \mathcal{B} \setminus \Gamma^L$  abgeschlossen, wobei  $\Gamma^L$  die Menge aller lokalen Maxima der logarithmierten Likelihoodfunktion im Innern von  $\mathcal{B}$  bezeichnet.
2.  $L(\hat{\beta}^{(i+1)}) > L(\hat{\beta}^{(i)})$  für alle  $\hat{\beta}^{(i)} \notin \Gamma^L$ .

Dann sind alle Konvergenzpunkte von  $\{\hat{\beta}^{(i)}\}$  lokale Maxima von  $L$  und  $L(\hat{\beta}^{(i)})$  konvergiert monoton mit dem Grenzwert  $\lim L(\hat{\beta}^{(i)}) = L^* = L(\beta^*)$  für ein  $\beta^* \in \Gamma^L$  (322.7).

Mit den Sätzen (322.6) bzw. (322.7) sind die Voraussetzungen für die Konvergenz der logarithmierten Likelihoodfunktion gegen lokale Maxima bzw. gegen stationäre Punkte in allgemeiner Weise für den GEM-Algorithmus formuliert. Die Voraussetzung 1. der Abgeschlossenheit der Abbildung  $M$  in  $\mathcal{B} \setminus \Gamma$  ist hierbei allerdings wenig anschaulich und kann für den GEM-Algorithmus in der konkreten Anwendung nicht oder nur schwer nachgeprüft werden. Im Gegensatz dazu ergibt sich für den EM-Algorithmus als Spezialfall des GEM-Algorithmus eine anschauliche und verifizierbare Bedingung für die Abgeschlossenheit der Punkt-Menge-Abbildung  $M$ : Nach [Wu 1982], Gl. (10) kann gezeigt werden, daß die Abgeschlossenheit der zum  $M$ -Schritt des EM-Algorithmus gehörenden Punkt-Menge-Abbildung  $M$  bereits folgt, wenn die Kullback-Leibler-Statistik  $Q(\psi | \phi)$  sowohl in  $\psi$  als auch in  $\phi$  stetig ist. Da die meisten in den Anwendungen verwendeten Dichtefunktionen diese Bedingungen erfüllen, ist mit der Forderung nach Abgeschlossenheit von  $M$  eine nur geringe Einschränkung der Allgemeinheit verbunden.

Wie sich aus nachfolgendem Satz ergibt, folgt aus der Stetigkeit von  $Q(\psi | \phi)$  in  $\psi$  und  $\phi$  nicht nur die Abgeschlossenheit der Abbildung  $M$  sondern auch die Konvergenz der logarithmierten Likelihoodfunktion gegen einen stationären Punkt.

**Satz:** Die Kullback-Leibler-Statistik  $Q(\psi | \phi)$  sei stetig sowohl in  $\psi$  als auch in  $\phi$ . Dann sind alle Konvergenzpunkte jeder Realisierung  $\{\hat{\beta}^{(i)}\}_{i \geq 0}$  eines EM-Algorithmus stationäre Punkte der logarithmierten Likelihoodfunktion  $L(\hat{\beta}^{(i)})$  und  $L(\hat{\beta}^{(i)})$  konvergiert monoton gegen  $L^* = L(\beta^*)$  für einen stationären Punkt  $\beta^*$ . (322.8)

**BEWEIS:** Es muß nachgewiesen werden, daß die Voraussetzungen des Satzes (322.6) erfüllt sind. Da die Stetigkeitsbedingung nach Voraussetzung erfüllt ist, ist die Forderung nach Abgeschlossenheit der Punkt-Menge-Abbildung  $M$  erfüllt. Es muß also noch gezeigt werden, daß  $L(\hat{\beta}^{(i+1)}) > L(\hat{\beta}^{(i)})$  für alle  $\hat{\beta}^{(i)} \notin \Gamma^S$  gilt.

Betrachtet werde ein Punkt  $\hat{\beta}^{(i)} \in \mathcal{B}$ , der nicht zur Menge  $\Gamma^S$  der stationären Punkte gehört, also  $\hat{\beta}^{(i)} \in \mathcal{B} \setminus \Gamma^S$ . Nach (231.11) gilt  $H(\hat{\beta}^{(i)} | \hat{\beta}^{(i)}) \geq H(\beta | \hat{\beta}^{(i)})$  für alle  $\beta \in$

$\mathcal{B} \setminus \{\hat{\beta}^{(i)}\}$ , so daß  $\beta = \hat{\beta}^{(i)}$  die Größe  $H(\beta \mid \hat{\beta}^{(i)})$  maximiert. Als notwendige Bedingung hierfür gilt<sup>8</sup>

$$\left. \frac{\partial H(\beta \mid \hat{\beta}^{(i)})}{\partial \beta} \right|_{\hat{\beta}^{(i)}} =: \mathbf{D}^{10} H(\beta \mid \hat{\beta}^{(i)}) \Big|_{\hat{\beta}^{(i)}} = \mathbf{D}^{10} H(\hat{\beta}^{(i)} \mid \hat{\beta}^{(i)}) = 0$$

Dieses Ergebnis setzt man in

$$\begin{aligned} \mathbf{D}L(\hat{\beta}^{(i)}) &:= \left. \frac{\partial L(\beta)}{\partial \beta} \right|_{\hat{\beta}^{(i)}} = \underbrace{\left. \frac{\partial Q(\beta \mid \hat{\beta}^{(i)})}{\partial \beta} \right|_{\hat{\beta}^{(i)}}}_{=: \mathbf{D}^{10} Q(\hat{\beta}^{(i)} \mid \hat{\beta}^{(i)})} - \underbrace{\left. \frac{\partial H(\beta \mid \hat{\beta}^{(i)})}{\partial \beta} \right|_{\hat{\beta}^{(i)}}}_{=: \mathbf{D}^{10} H(\hat{\beta}^{(i)} \mid \hat{\beta}^{(i)})=0} = \mathbf{D}^{10} Q(\hat{\beta}^{(i)} \mid \hat{\beta}^{(i)}) \end{aligned}$$

ein. Da der Punkt  $\hat{\beta}^{(i)} \notin \Gamma^S$  kein stationärer Punkt der logarithmierten Likelihoodfunktion ist, muß der Gradient  $\mathbf{D}L(\hat{\beta}^{(i)})$  und damit  $\mathbf{D}^{10} Q(\hat{\beta}^{(i)} \mid \hat{\beta}^{(i)})$  in diesem Punkt vom Nullvektor verschieden sein:

$$\mathbf{D}L(\hat{\beta}^{(i)}) = \mathbf{D}^{10} Q(\hat{\beta}^{(i)} \mid \hat{\beta}^{(i)}) \neq \mathbf{0} \quad \text{für jedes } \hat{\beta}^{(i)} \notin \Gamma^S.$$

Nach der Definition (232.2) des EM-Algorithmus nimmt in der  $i$ -ten Iteration die Kullback-Leibler-Statistik  $Q(\beta \mid \hat{\beta}^{(i)})$  für die neue Schätzung  $\beta = \hat{\beta}^{(i+1)}$  der unbekannt Parameter ihr globales Maximum an. An der Stelle  $\beta = \hat{\beta}^{(i)}$  kann  $Q$  diesen Maximalwert nicht annehmen, da wegen  $\mathbf{D}^{10} Q(\hat{\beta}^{(i)} \mid \hat{\beta}^{(i)}) \neq \mathbf{0}$  die notwendige Bedingung für ein lokales Maximum nicht erfüllt ist, so daß es mindestens einen Vektor  $\beta$  mit  $Q(\beta \mid \hat{\beta}^{(i)}) > Q(\hat{\beta}^{(i)} \mid \hat{\beta}^{(i)})$  gibt. Es muß daher für die neue Schätzung  $\hat{\beta}^{(i+1)}$  gelten

$$Q(\hat{\beta}^{(i+1)} \mid \hat{\beta}^{(i)}) > Q(\hat{\beta}^{(i)} \mid \hat{\beta}^{(i)}),$$

womit unter Berücksichtigung von  $H(\hat{\beta}^{(i)} \mid \hat{\beta}^{(i)}) \geq H(\hat{\beta}^{(i+1)} \mid \hat{\beta}^{(i)})$  aus (231.11) und mit (321.5) folgt

$$\begin{aligned} L(\hat{\beta}^{(i+1)}) &= Q(\hat{\beta}^{(i+1)} \mid \hat{\beta}^{(i)}) - H(\hat{\beta}^{(i+1)} \mid \hat{\beta}^{(i)}) \\ &> Q(\hat{\beta}^{(i)} \mid \hat{\beta}^{(i)}) - H(\hat{\beta}^{(i)} \mid \hat{\beta}^{(i)}) \\ &\stackrel{(321.5)}{=} Q(\hat{\beta}^{(i)} \mid \hat{\beta}^{(i-1)}) - H(\hat{\beta}^{(i)} \mid \hat{\beta}^{(i-1)}) = L(\hat{\beta}^{(i)}). \end{aligned}$$

Hieraus folgt die Aussage. □

Ein dem Satz (322.8) entsprechender Satz für die konvergenz des EM-Algorithmus gegen ein lokales Maximum gilt *nicht*.

**BEGRÜNDUNG:** Anhand eines Gegenbeispiels kann gezeigt werden, daß aus  $\hat{\beta}^{(i)} \notin \Gamma^L$  nicht zwangsläufig  $L(\hat{\beta}^{(i+1)} \mid \hat{\beta}^{(i)})$  folgt. Betrachtet sei der Fall, in dem es sich bei  $\hat{\beta}^{(i)}$  zwar um kein lokales Maximum, aber um einen stationären Punkt der logarithmierten Likelihoodfunktion handelt, also  $\hat{\beta}^{(i)} \notin \Gamma^L$  aber  $\hat{\beta}^{(i)} \in \Gamma^S$ . Wegen der Stationarität von  $L$  an der Stelle  $\hat{\beta}^{(i)}$  gilt dann

$$\mathbf{D}L(\hat{\beta}^{(i)}) = \mathbf{D}^{10} Q(\hat{\beta}^{(i)} \mid \hat{\beta}^{(i)}) = 0,$$

und somit ist  $L(\hat{\beta}^{(i+1)}) > L(\hat{\beta}^{(i)})$  in stationären Punkten nicht gegeben. Ausgehend von einem stationären Punkt erfolgt also keine weitere Konvergenz gegen ein lokales Maximum. Damit ist Bedingung 2. von Satz (322.7) nicht erfüllt und ein dem Satz (322.8) entsprechender Satz für die Konvergenz der logarithmierten Likelihoodfunktion gegen ein lokales Maximum gilt nicht.

<sup>8</sup>Im Rahmen dieser Arbeit bezeichnet  $\mathbf{D}^{lm} H(\phi \mid \psi)$  bzw.  $\mathbf{D}^{lm} Q(\phi \mid \psi)$  die gemischt-partielle Ableitung

$$\mathbf{D}^{lm} H(\phi \mid \psi) = \frac{\partial^{(l+m)} H(\phi \mid \psi)}{\partial \phi^l \partial \psi^m} \quad \text{bzw.} \quad \mathbf{D}^{lm} Q(\phi \mid \psi) = \frac{\partial^{(l+m)} Q(\phi \mid \psi)}{\partial \phi^l \partial \psi^m}.$$

Um sicherzustellen, daß der EM-Algorithmus gegen ein lokales Maxima konvergiert, muss eine weitere Bedingung eingeführt werden. Falls innerhalb der EM-Iterationen der Fall eintritt, daß die Schätzwerte  $\hat{\beta}^{(i)}$  einen stationären Punkt der Likelihoodfunktion darstellen, so muss im nächsten Iterationsschritt ein  $\hat{\beta}^{(i+1)}$  gefunden werden, der eine größere Kullback-Leibler-Statistik liefert. Eine Bedingung hierfür ist, daß für alle  $\hat{\beta}^{(i)} \in \Gamma^S \setminus \Gamma^L$  gilt  $\sup_{\beta \in \mathcal{B}} Q(\beta \mid \hat{\beta}^{(i)}) > Q(\hat{\beta}^{(i)} \mid \hat{\beta}^{(i)})$ . Mit dieser Bedingung ergibt sich der

**Satz:** Die Kullback-Leibler-Statistik  $Q(\psi \mid \phi)$  sei stetig sowohl in  $\psi$  als auch in  $\phi$  und es gelte

$$\sup_{\beta \in \mathcal{B}} Q(\beta \mid \hat{\beta}^{(i)}) > Q(\hat{\beta}^{(i)} \mid \hat{\beta}^{(i)}) \quad \text{für jedes } \hat{\beta}^{(i)} \in \Gamma^S \setminus \Gamma^L. \quad (322.9)$$

Dann sind alle Konvergenzpunkte jeder Realisierung  $\{\hat{\beta}^{(i)}\}_{i \geq 0}$  eines EM-Algorithmus lokale Maxima der logarithmierten Likelihoodfunktion  $L(\hat{\beta}^{(i)})$  und  $L(\hat{\beta}^{(i)})$  konvergiert gegen  $L^* = L(\beta^*)$  für ein lokales Maximum  $\beta^* \in \Gamma^L$ .

Die Voraussetzung (322.9) lässt sich praktisch schwer nachprüfen und ist daher überwiegend von theoretischem Interesse. Alternative Bedingungen für die Konvergenz der logarithmierten Likelihoodfunktion gegen ein lokales Maximum werden in der Literatur jedoch nicht gegeben.

Zusammenfassend lassen sich folgende Eigenschaften des EM-Algorithmus in Bezug auf die Konvergenz der logarithmierten Likelihoodfunktion festhalten:

**Zur Konvergenz der logarithmierten Likelihoodfunktion im Rahmen der EM-Iterationen:**

1. Ist die Kullback-Leibler-Statistik  $Q(\psi \mid \phi)$  stetig in beiden Argumenten, so konvergiert die Likelihoodfunktion im Rahmen des EM-Algorithmus gegen einen stationären Punkt.
2. Gilt zusätzlich  $\sup_{\beta \in \mathcal{B}} Q(\beta \mid \hat{\beta}^{(i)}) > Q(\hat{\beta}^{(i)} \mid \hat{\beta}^{(i)})$  für jedes  $\hat{\beta}^{(i)} \in \mathcal{L} \setminus \mathcal{T}$ , so konvergiert die Likelihoodfunktion gegen ein lokales Maximum.

### 323 Betrachtungen zur Konvergenz der Folge geschätzter Parameter

Die bisherigen Konvergenzbetrachtungen beziehen sich ausschließlich auf die Konvergenz der logarithmierten Likelihoodfunktion im Rahmen des EM-Algorithmus. Es wurde gezeigt, daß die Folge der logarithmierten Likelihoodwerte stets konvergiert und daß unter bestimmten Bedingungen Konvergenz gegen einen stationären Punkt oder gegen ein lokales Maximum der logarithmierten Likelihoodfunktion eintritt. Die Konvergenz der Folge  $\{L(\hat{\beta}^{(i)})\}_{i \geq 0}$  stellt jedoch keinesfalls sicher, daß auch die Folge  $\{\hat{\beta}^{(i)}\}_{i \geq 0}$  der Argumente der logarithmierten Likelihoodfunktion, also der Schätzungen  $\hat{\beta}^{(i)}$  für die unbekannt Parameter  $\beta$  konvergiert<sup>9</sup>. Da das Ziel des (G) EM-Algorithmus letztlich darin besteht, optimale Schätzwerte zu bestimmen, ist die Konvergenz der Folge  $\{\hat{\beta}^{(i)}\}_{i \geq 0}$  ebenfalls von

<sup>9</sup>Nach [WU 1982] ist die umgekehrte Aussage jedoch richtig, falls  $D^{10}H(\beta \mid \hat{\beta}^{(i)})$  stetig ist, d.h. aus der Konvergenz der Folge  $\{\hat{\beta}^{(i)}\}_{i \geq 0}$  folgt dann die Konvergenz der Folge  $\{L(\hat{\beta}^{(i)})\}_{i \geq 0}$ .

Interesse. Im folgenden wird daher erläutert, unter welchen Bedingungen die Folge der Schätzungen im Rahmen des (G)EM-Algorithmus konvergiert und um welche Art Grenzwert es sich im Falle der Konvergenz handelt. Hierzu werden zunächst alle stationären Punkte bzw. lokalen Maxima der logarithmierten Likelihoodfunktion, zu denen der selbe Likelihoodwert  $a$  gehört, in einer Menge  $\Gamma^S(a)$  bzw.  $\Gamma^L(a)$  zusammengefasst:

$$\Gamma^S(a) := \{\boldsymbol{\beta} \in \Gamma^S : L(\boldsymbol{\beta}) = a\} \quad (323.1)$$

$$\Gamma^L(a) := \{\boldsymbol{\beta} \in \Gamma^L : L(\boldsymbol{\beta}) = a\} \quad (323.2)$$

Nach den Sätzen (322.6) bzw. (322.7) konvergiert  $L(\hat{\boldsymbol{\beta}}^{(i)})$  im Rahmen der GEM-Iterationen unter bestimmten Bedingungen,  $L(\hat{\boldsymbol{\beta}}^{(i)}) \rightarrow L(\boldsymbol{\beta}^*) = L^*$ , wobei jeder Konvergenzpunkt  $\hat{\boldsymbol{\beta}}^*$  in der Menge  $\Gamma^S$  der stationären Punkte bzw. in der Menge  $\Gamma^L$  der lokalen Maxima liegt. Besteht die Menge  $\Gamma^S(L^*)$  bzw. die Menge  $\Gamma^L(L^*)$  nur aus einem einzigen Punkt  $\boldsymbol{\beta}^*$ , so konvergiert die Folge unter diesen Bedingungen offensichtlich gegen diesen Punkt. Es gilt somit für die Konvergenz der Folge  $\{\hat{\boldsymbol{\beta}}^{(i)}\}_{i \geq 0}$  der aus Satz (322.6) resultierende

**Satz:** *Es sei  $\{\hat{\boldsymbol{\beta}}^{(i)}\}_{i \geq 0}$  eine Folge von Schätzungen  $\hat{\boldsymbol{\beta}}^{(i)}$  für die unbekannt Parameter  $\boldsymbol{\beta}$  beim GEM-Algorithmus und es sei*

1. *die Abbildung  $M$ , die die Schätzung  $\hat{\boldsymbol{\beta}}^{(i)}$  der unbekannt Parameter der  $i$ -ten Iteration auf die Lösungsmenge  $M^{(i+1)} = M(\hat{\boldsymbol{\beta}}^{(i)})$  des  $M$ -Schrittes der  $(i+1)$ -ten Iteration abbildet,  $M$  für  $\hat{\boldsymbol{\beta}}^{(i)} \in \mathcal{B} \setminus \Gamma^S$  abgeschlossen, wobei  $\Gamma^S$  die Menge aller stationären Punkte der logarithmierten Likelihoodfunktion im Innern des Parameterraumes  $\mathcal{B}$  bezeichnet;*
2.  *$L(\hat{\boldsymbol{\beta}}^{(i+1)}) > L(\hat{\boldsymbol{\beta}}^{(i)})$  für alle  $\hat{\boldsymbol{\beta}}^{(i)} \notin \Gamma^L$ ,*

*so daß gemäß (322.6) die Folge  $\{L(\hat{\boldsymbol{\beta}}^{(i)})\}_{i \geq 0}$  der logarithmierten Likelihoodwerte  $L(\hat{\boldsymbol{\beta}}^{(i)})$  für  $i \rightarrow \infty$  gegen den Grenzwert  $L^* = L(\boldsymbol{\beta}^*)$  mit  $\boldsymbol{\beta}^* \in \Gamma^S$  konvergiert.*

*Besteht ferner die Menge  $\Gamma^S(L^*) = \{\boldsymbol{\beta} \in \Gamma^S : L(\boldsymbol{\beta}) = L^*\}$  der stationären Punkte von  $L(\boldsymbol{\beta})$ , deren logarithmierter Likelihoodwert gleich dem Grenzwert  $L^*$  ist, nur aus einem einzigen Element, also  $\Gamma^S(L^*) = \{\boldsymbol{\beta}^*\}$ , so konvergiert  $\{\hat{\boldsymbol{\beta}}^{(i)}\}_{i \geq 0}$  für  $i \rightarrow \infty$  gegen diesen stationären Punkt  $\boldsymbol{\beta}^*$ .* (323.5)

Entsprechend gilt, resultierend aus Satz (322.7) für die Konvergenz der Folge  $\{\hat{\boldsymbol{\beta}}^{(i)}\}_{i \geq 0}$ , falls  $\{L(\hat{\boldsymbol{\beta}}^{(i)})\}_{i \geq 0}$  gegen ein lokales Maximum der logarithmierten Likelihoodfunktion konvergiert

**Satz:** *Es sei  $\{\hat{\boldsymbol{\beta}}^{(i)}\}_{i \geq 0}$  eine Folge von Schätzungen  $\hat{\boldsymbol{\beta}}^{(i)}$  für die unbekannt Parameter  $\boldsymbol{\beta}$  beim GEM-Algorithmus und es sei*

1. *die Abbildung  $M$ , die die Schätzung  $\hat{\boldsymbol{\beta}}^{(i)}$  der unbekannt Parameter der  $i$ -ten Iteration auf die Lösungsmenge  $M^{(i+1)} = M(\hat{\boldsymbol{\beta}}^{(i)})$  des  $M$ -Schrittes der  $(i+1)$ -ten Iteration abbildet, für  $\hat{\boldsymbol{\beta}}^{(i)} \in \mathcal{B} \setminus \Gamma^L$  abgeschlossen, wobei  $\Gamma^L$  die Menge aller lokalen Maxima der logarithmierten Likelihoodfunktion im Innern des Parameterraumes  $\mathcal{B}$  bezeichnet;*
2.  *$L(\hat{\boldsymbol{\beta}}^{(i+1)}) > L(\hat{\boldsymbol{\beta}}^{(i)})$  für alle  $\hat{\boldsymbol{\beta}}^{(i)} \notin \Gamma^L$ ,*

*so daß gemäß (322.7) die Folge  $\{L(\hat{\boldsymbol{\beta}}^{(i)})\}_{i \geq 0}$  der logarithmierten Likelihoodwerte  $L(\hat{\boldsymbol{\beta}}^{(i)})$  für  $i \rightarrow \infty$  gegen den Grenzwert  $L^* = L(\boldsymbol{\beta}^*)$  mit  $\boldsymbol{\beta}^* \in \Gamma^L$  konvergiert.*

*Besteht weiter die Menge  $\Gamma^L(L^*) = \{\boldsymbol{\beta} \in \Gamma^L : L(\boldsymbol{\beta}) = L^*\}$  der lokalen Maxima von  $L(\boldsymbol{\beta})$ , deren logarithmierter Likelihoodwert gleich dem Grenzwert  $L^*$  ist, nur aus einem einzigen Element, also  $\Gamma^L(L^*) = \{\boldsymbol{\beta}^*\}$ , so konvergiert  $\{\hat{\boldsymbol{\beta}}^{(i)}\}_{i \geq 0}$  für  $i \rightarrow \infty$  gegen dieses lokale Maximum  $\boldsymbol{\beta}^*$ .*

(323.4)

Die Bedingungen, daß die Menge  $\Gamma^S(L^*)$  der stationären Punkte bzw. die Menge  $\Gamma^L(L^*)$  der lokalen Maxima von  $L(\beta)$  nur aus einem Element bestehen, sind relativ scharf und für viele multimodale<sup>10</sup> Dichtefunktionen nicht erfüllt. Es können jedoch auch weniger scharfe Bedingungen formuliert werden, unter denen die Parameterfolge  $\{\hat{\beta}^{(i)}\}_{i \geq 0}$  des (G)EM-Algorithmus konvergiert. Beispielsweise kann gezeigt werden, daß – falls  $\Gamma$  eine diskrete<sup>11</sup> Menge ist – die Folge  $\{\hat{\beta}^{(i)}\}_{i \geq 0}$  bereits konvergiert, wenn  $\lim_{i \rightarrow \infty} \|\hat{\beta}^{(i+1)} - \hat{\beta}^{(i)}\| = 0$  gilt (vgl.[WU 1982]). Eine entsprechende Aussage liefert der

**Satz:** *Es sei  $\{\hat{\beta}^{(i)}\}_{i \geq 0}$  eine Folge von Schätzungen  $\hat{\beta}^{(i)}$  für die unbekannt Parameter  $\beta$  beim GEM-Algorithmus und es sei*

1. *die Abbildung  $M$ , die die Schätzung  $\hat{\beta}^{(i)}$  der unbekannt Parameter der  $i$ -ten Iteration auf die Lösungsmenge  $M^{(i+1)} = M(\hat{\beta}^{(i)})$  des  $M$ -Schrittes der  $(i+1)$ -ten Iteration abbildet, für  $\hat{\beta}^{(i)} \in \mathcal{B} \setminus \Gamma^S$  abgeschlossen, wobei  $\Gamma^S$  die Menge aller stationären Punkte der logarithmierten Likelihoodfunktion im Innern des Parameterraumes  $\mathcal{B}$  bezeichnet;*
2.  *$L(\hat{\beta}^{(i+1)}) > L(\hat{\beta}^{(i)})$  für alle  $\hat{\beta}^{(i)} \notin \Gamma^S$ ,*

*so daß gemäß (322.6) die Folge  $\{L(\hat{\beta}^{(i)})\}_{i \geq 0}$  der logarithmierten Likelihoodwerte  $L(\hat{\beta}^{(i)})$  für  $i \rightarrow \infty$  gegen den Grenzwert  $L^* = L(\beta^*)$  mit  $\beta^* \in \Gamma^S$  konvergiert.*

*Gilt weiter  $\|\hat{\beta}^{(i+1)} - \hat{\beta}^{(i)}\| \rightarrow \mathbf{0}$  für  $i \rightarrow \infty$ , dann liegen alle Konvergenzpunkte  $\beta^*$  von  $\{\hat{\beta}^{(i)}\}_{i \geq 0}$  in einer zusammenhängenden und kompakten Teilmenge von  $\Gamma^S(L^*)$ . Handelt es sich bei  $\Gamma^S(L^*)$  um eine diskrete Menge, dann konvergiert  $\hat{\beta}^{(i)}$  gegen ein  $\beta^*$  in  $\Gamma^S(L^*)$ .* (323.5)

BEWEIS: Hinsichtlich dieses Beweises sei auf [WU 1982] verwiesen. □

Wie bereits aus den Vorbemerkungen zu Satz (322.8) hervorgeht ist die Abgeschlossenheitsbedingung 1.) speziell für den EM-Algorithmus erfüllt, wenn die Kullback-Leibler-Statistik  $Q(\psi | \phi)$  sowohl in  $\psi$  als auch in  $\phi$  stetig ist. Weiter geht aus dem Beweis zu Satz (322.8) hervor, daß beim EM-Algorithmus auch Bedingung 2.) von (323.5) erfüllt ist. Damit ergibt sich der

**Satz:** *Beim EM-Algorithmus sei die Kullback-Leibler-Statistik  $Q(\psi | \phi)$  stetig sowohl in  $\psi$  als auch in  $\phi$ . Gilt weiter  $\lim_{i \rightarrow \infty} \|\hat{\beta}^{(i+1)} - \hat{\beta}^{(i)}\| = \mathbf{0}$ , dann liegen alle Konvergenzpunkte  $\beta^*$  von  $\{\hat{\beta}^{(i)}\}_{i \geq 0}$  in einer zusammenhängenden und kompakten Teilmenge von  $\Gamma^S(L^*)$ . Handelt es sich bei  $\Gamma^S(L^*)$  um eine diskrete Menge, dann konvergiert  $\hat{\beta}^{(i)}$  gegen ein  $\beta^*$  in  $\Gamma^S(L^*)$ .* (323.6)

Ein dem Satz (323.5) entsprechender Satz ergibt sich wieder für die Konvergenz gegen lokale Maxima der logarithmierten Likelihoodfunktion:

**Satz:** *Es sei  $\{\hat{\beta}^{(i)}\}_{i \geq 0}$  eine Folge von Schätzungen  $\hat{\beta}^{(i)}$  für die unbekannt Parameter  $\beta$  beim GEM-Algorithmus und es sei*

1. *die Abbildung  $M$ , die die Schätzung  $\hat{\beta}^{(i)}$  der unbekannt Parameter der  $i$ -ten Iteration auf die Lösungsmenge  $M^{(i+1)} = M(\hat{\beta}^{(i)})$  des  $M$ -Schrittes der  $(i+1)$ -ten Iteration*

<sup>10</sup>Eine multimodale Dichtefunktion ist eine Dichtefunktion, die durch Zusammensetzung (beispielsweise als Linearkombination) mehrerer Basisdichtefunktionen entsteht.

<sup>11</sup>Eine Menge wird als *diskrete* Menge bezeichnet, falls ihre Zusammenhangskomponenten jeweils nur aus genau einem Element bestehen.

abbildet, für  $\hat{\beta}^{(i)} \in \mathcal{B} \setminus \Gamma^L$  abgeschlossen, wobei  $\Gamma^L$  die Menge aller lokalen Maxima der logarithmierten Likelihoodfunktion im Innern des Parameterraumes  $\mathcal{B}$  bezeichnet;

$$2. \quad L(\hat{\beta}^{(i+1)}) > L(\hat{\beta}^{(i)}) \text{ für alle } \hat{\beta}^{(i)} \notin \Gamma^L,$$

so daß gemäß (322.6) die Folge  $\{L(\hat{\beta}^{(i)})\}_{i \geq 0}$  der logarithmierten Likelihoodwerte  $L(\hat{\beta}^{(i)})$  für  $i \rightarrow \infty$  gegen den Grenzwert  $L^* = L(\beta^*)$  mit  $\beta^* \in \Gamma^L$  konvergiert.

Gilt weiter  $\|\hat{\beta}^{(i+1)} - \hat{\beta}^{(i)}\| \rightarrow \mathbf{0}$  für  $i \rightarrow \infty$ , dann liegen alle Konvergenzpunkte  $\beta^*$  von  $\{\hat{\beta}^{(i)}\}_{i \geq 0}$  in einer zusammenhängenden und kompakten Teilmenge von  $\Gamma^L(L^*)$ . Handelt es sich bei  $\Gamma^L(L^*)$  um eine diskrete Menge, dann konvergiert  $\hat{\beta}^{(i)}$  gegen ein  $\beta^*$  in  $\Gamma^L(L^*)$ . (323.7)

An dieser Stelle soll ein weiterer Satz angeführt werden, der Bedingungen für die Konvergenz der Parameterfolge  $\{\hat{\beta}^{(i)}\}_{i \geq 0}$  im Rahmen eines (G)EM-Algorithmus nennt. Zunächst wird dazu die Abkürzung

$$\mathcal{L}(L) = \{\beta \in \mathcal{B} : L(\beta) = L\} \quad (323.8)$$

eingeführt für die Menge aller Parametervektoren  $\beta$ , für die die logarithmierte Likelihoodfunktion  $L(\beta)$  den Wert  $L$  annimmt. Hiermit gilt dann der

**Satz:** Es sei  $\{\hat{\beta}^{(i)}\}_{i \geq 0}$  eine mittels des GEM-Algorithmus erhaltene Folge von Schätzungen  $\hat{\beta}^{(i)}$  für die unbekannt Parameter  $\beta$  und es gelte  $\mathbf{D}^{10}Q(\hat{\beta}^{(i+1)} | \hat{\beta}^{(i)}) = \mathbf{0}$ . Weiter sei  $\mathbf{D}^{10}Q(\psi | \phi)$  sowohl in  $\psi$  als auch in  $\phi$  stetig.

Dann konvergiert  $\{\hat{\beta}^{(i)}\}_{i \geq 0}$  gegen einen stationären Punkt  $\beta^*$  der logarithmierten Likelihoodfunktion  $L(\beta)$  mit  $L(\beta^*) = \lim_{i \rightarrow \infty} L(\hat{\beta}^{(i)}) = L^*$ , wenn entweder

- (a) die Menge  $\mathcal{L}(L^*)$  aus einem einzigen Element besteht, also  $\mathcal{L}(L^*) = \{\beta^*\}$  oder  
 (b)  $\lim_{i \rightarrow \infty} \|\hat{\beta}^{(i+1)} - \hat{\beta}^{(i)}\| = \mathbf{0}$  gilt und  $\mathcal{L}(L^*)$  diskret ist. (323.9)

**BEWEIS:** Zunächst muß nachgewiesen werden, daß die Folge  $\{\hat{\beta}^{(i)}\}_{i \geq 0}$  für  $i \rightarrow \infty$  konvergiert, falls entweder Bedingung (a) oder Bedingung (b) erfüllt ist.

Wegen der vorausgesetzten Stetigkeit von  $\mathbf{D}^{10}Q(\psi | \phi)$  sowohl in  $\psi$  als auch in  $\phi$  folgt nach [HELFRICH I 1995] zunächst die Differenzierbarkeit von  $Q(\psi | \phi)$  und schließlich die Stetigkeit von  $Q(\psi | \phi)$  in  $\psi$  und  $\phi$ . Wie aus dem Beweis von Satz (322.8) hervorgeht, folgt aus der Stetigkeit von  $Q$  und der Voraussetzung  $\mathbf{D}^{10}Q(\hat{\beta}^{(i+1)} | \hat{\beta}^{(i)}) = 0$  die Konvergenz der Folge  $\{L(\hat{\beta}^{(i)})\}_{i \geq 0}$  gegen einen Grenzwert  $L^*$ , so daß die Schätzungen  $\hat{\beta}^{(i)}$  für  $i \rightarrow \infty$  der Menge  $\mathcal{L}(L^*)$  angehören müssen.<sup>12</sup> Besteht gemäß Bedingung (a) die Menge  $\mathcal{L}(L^*)$  nur aus einem einzigen Punkt, so folgt hieraus offensichtlich die Konvergenz der Parameterfolge,  $\lim_{i \rightarrow \infty} \hat{\beta}^{(i)} = \beta^* \in \mathcal{L}(L^*)$  gegen diesen Punkt  $\beta^*$ . Sofern Bedingung (b) gilt, folgt die Konvergenz der Parameterfolge  $\{\hat{\beta}^{(i)}\}_{i \geq 0}$  gegen (genau) ein  $\beta^* \in \mathcal{L}(L^*)$  wegen Satz (323.5).

Nachdem die Konvergenz der Parameterfolge bewiesen ist, muß nachgewiesen werden, daß es sich bei dem Grenzwert  $\beta^*$  in beiden Fällen (a) und (b) um einen stationären Punkt der logarithmierten Likelihoodfunktion handelt. Hierzu ist zu zeigen, daß  $\mathbf{D}L(\beta^*)$  identisch mit dem Nullvektor ist. Es gilt zunächst

$$\begin{aligned} \mathbf{D}L(\beta^*) &= \lim_{i \rightarrow \infty} \mathbf{D}L(\hat{\beta}^{(i+1)}) \\ &= \lim_{i \rightarrow \infty} \mathbf{D}^{10}Q(\hat{\beta}^{(i+1)} | \hat{\beta}^{(i)}) - \lim_{i \rightarrow \infty} \mathbf{D}^{10}H(\hat{\beta}^{(i+1)} | \hat{\beta}^{(i)}). \end{aligned} \quad (323.10)$$

<sup>12</sup>An dieser Stelle sei angemerkt, daß hierdurch keineswegs die Konvergenz der Folge  $\{\hat{\beta}^{(i)}\}_{i \geq 0}$  nachgewiesen ist, da die Möglichkeit besteht, daß  $\hat{\beta}^{(i)}$  für  $i \rightarrow \infty$  zwischen mehreren Vektoren der Menge  $\mathcal{L}(L^*)$  alterniert.



Nach Voraussetzung gilt  $\mathbf{D}^{10}Q(\hat{\beta}^{(i+1)} | \hat{\beta}^{(i)}) = \mathbf{0}$ , so daß wegen der vorausgesetzten Stetigkeit von  $\mathbf{D}^{10}Q(\psi | \phi)$  folgt

$$\mathbf{D}^{10}Q(\beta^* | \beta^*) = \lim_{i \rightarrow \infty} \mathbf{D}^{10}Q(\hat{\beta}^{(i+1)} | \hat{\beta}^{(i)}) = \mathbf{0}. \quad (323.11)$$

Da nach (231.11) für alle  $\beta \in \mathcal{B}$  gilt  $H(\hat{\beta}^{(i)} | \hat{\beta}^{(i)}) \geq H(\beta | \hat{\beta}^{(i)})$ , nimmt  $H(\psi | \hat{\beta}^{(i)})$  für  $\psi = \hat{\beta}^{(i)}$  ein lokales Maximum an, es gilt also  $\mathbf{D}^{10}H(\hat{\beta}^{(i)} | \hat{\beta}^{(i)}) = \mathbf{0}$ . Damit folgt

$$\lim_{i \rightarrow \infty} \mathbf{D}^{10}H(\hat{\beta}^{(i+1)} | \hat{\beta}^{(i)}) = \mathbf{D}^{10}H(\beta^* | \beta^*) = \mathbf{0}. \quad (323.12)$$

Einsetzen von (323.12) und (323.11) in (323.10) liefert

$$DL(\beta^*) = \mathbf{0},$$

womit folgt, daß  $\beta^*$  ein stationärer Punkt der logarithmierten Likelihoodfunktion  $L(\beta)$  ist. Damit ist der Satz bewiesen.  $\square$

Der Vorteil von Satz (323.9) gegenüber den vergleichbaren Sätzen (323.3) und (323.5) besteht darin, daß er ohne die (in der konkreten Anwendung schwer verifizierbaren) Bedingungen 1. und 2. dieser Sätze auskommt. An die Stelle dieser Bedingungen treten die Forderung nach Stetigkeit von  $\mathbf{D}^{10}Q(\psi | \phi)$  und die Bedingung  $\mathbf{D}^{10}Q(\hat{\beta}^{(i+1)} | \hat{\beta}^{(i)}) = \mathbf{0}$ . Die Stetigkeitsbedingung kann in diesem Zusammenhang für eine große Anzahl von Dichtefunktionen  $f(z | \beta)$  als erfüllt gelten. Insbesondere für den EM-Algorithmus ist auch die Voraussetzung  $\mathbf{D}^{10}Q(\hat{\beta}^{(i+1)} | \hat{\beta}^{(i)}) = \mathbf{0}$  erfüllt, da die Kullback-Leibler-Statistik  $Q(\beta | \hat{\beta}^{(i)})$  gemäß der Definition des M-Schrittes beim EM-Algorithmus für  $\beta = \hat{\beta}^{(i+1)}$  ein lokales Maximum der logarithmierten Likelihoodfunktion darstellt.

Offensichtlich handelt es sich bei der Menge  $\Gamma^S(L^*)$  um eine Teilmenge der Menge  $\mathcal{L}(L^*)$ , es gilt also  $\Gamma^S(L^*) \subset \mathcal{L}(L^*)$ . Daher stellt die Bedingung (a), daß  $\mathcal{L}(L^*)$  genau ein Element  $\beta^*$  enthält, eine schärfere Bedingung dar als die Voraussetzung des Satzes (323.3), daß  $\Gamma^S(L^*)$  aus genau einem Element besteht. Ebenso ist die Forderung (b) nach Diskretheit der Menge  $\mathcal{L}(L^*)$  schärfer als in (323.5) die Forderung nach Diskretheit der Menge  $\Gamma^S(L^*)$ . Im Vergleich zu den Sätzen (323.3) und (323.5) stellt Satz (323.9) also restriktivere Anforderungen an die Funktion  $L(\beta)$ .

Die Bedingung  $\lim_{i \rightarrow \infty} \|\hat{\beta}^{(i+1)} - \hat{\beta}^{(i)}\| = \mathbf{0}$  kann nur in Einzelfällen und dann auch nur mit verhältnismäßig großem Aufwand nachgewiesen werden. Daher kann (b) nur in Ausnahmefällen als Kriterium für die Konvergenz der Parameterfolge  $\{\hat{\beta}^{(i)}\}_{i \geq 0}$  dienen. Es verbleibt in der überwiegenden Zahl der Anwendungen das Kriterium (a).

Ein wichtiger Spezialfall des Satzes (323.9) liegt vor, wenn die logarithmierte Likelihoodfunktion eine unimodale Funktion ist, d.h. wenn  $L(\beta)$  nur einen einzigen stationären Punkt besitzt, der zugleich ein lokales Maximum von  $L$  ist. Dies ist beispielsweise der Fall, wenn es sich bei der Dichtefunktion  $f(z | \beta)$  der vollständigen Beobachtungen um die Normalverteilung handelt. Wird in einem solchen Fall mit dem EM-Algorithmus gearbeitet, so sind die Bedingungen (a) und  $Q(\hat{\beta}^{(i+1)} | \hat{\beta}^{(i)}) = \mathbf{0}$  von Satz (323.9) erfüllt und die Parameterfolge  $\{\hat{\beta}^{(i)}\}_{i \geq 0}$  konvergiert gegen den stationären Punkt, also das lokale Maximum der logarithmierten Likelihoodfunktion. Somit gilt der

**Satz:** Die logarithmierte Likelihoodfunktion  $L(\beta)$  sei im Innern des Parameterraums  $\mathcal{B}$  eine unimodale Funktion und  $\beta^*$  bezeichne den einzigen stationären Punkt. Weiter sei  $\mathbf{D}^{10}Q(\psi | \phi)$  stetig in  $\psi$  und in  $\phi$ . Dann konvergiert beim EM-Algorithmus die Folge

$\{\hat{\beta}^{(i)}\}_{i \geq 0}$  der in den einzelnen Iterationen bestimmten Schätzungen  $\hat{\beta}^{(i)}$  für die unbekannt Parameter  $\beta$  gegen den stationären Punkt  $\hat{\beta}^*$ , dem einzigen lokalen Maximum der logarithmierten Likelihoodfunktion  $L(\beta)$ .

Dieser Satz wird in der Praxis am häufigsten verwandt, um die Konvergenz der Parameterfolge gegen ein lokales Maximum der logarithmierten Likelihoodfunktion nachzuweisen.

Die bis hierhin abgeleiteten Eigenschaften des EM-Algorithmus in Bezug auf die Konvergenz der Parameterfolge  $\{\hat{\beta}^{(i)}\}_{i \geq 0}$  sind in Übersicht 3.1 zusammengefaßt.

Hierin sind vor allem die Bedingungen von Interesse, unter denen die Parameterfolge gegen ein lokales Maximum der logarithmierten Likelihoodfunktion konvergiert, da ja bei der Parameterschätzung das lokale Maximum der logarithmierten Likelihoodfunktion gesucht wird, das diese Funktion auch global maximiert. Im Allgemeinen werden mit dem EM-Algorithmus für unterschiedliche Startvektoren  $\hat{\beta}^{(0)}$  unterschiedliche Schätzungen  $\beta^*$  erhalten. Daher sollte der EM-Algorithmus zur Lösung eines Parameterschätzproblems stets mehrfach unter Verwendung verschiedener Startwerte durchgeführt werden, um die erhaltenen Ergebnisse einer Auswertung durch weitere Auswertungen mit unterschiedlichen Startwerten zu stützen.

## 324 Betrachtungen zur Konvergenzgeschwindigkeit

Wie bereits in Abschnitt 231 erläutert, besteht das Ziel des (G)EM-Algorithmus darin, Schätzwerte für die unbekannt Parameter zu erhalten, die die logarithmierte Likelihoodfunktion der unvollständigen Beobachtungsdaten maximieren. Nach Möglichkeit sollte die Schätzung eindeutig sein. In den vorangegangenen Abschnitten wurde diskutiert, unter welchen Bedingungen der EM-Algorithmus bzw. der GEM-Algorithmus zu diesem angestrebten Ziel führt: Konvergiert die Folge der logarithmierten Likelihoodwerte gegen ein lokales Maximum, so werden mittels des (G)EM-Algorithmus Schätzwerte für die unbekannt Parameter erhalten. Konvergiert darüber hinaus auch die Parameterfolge gegen ein lokales Maximum, so sind diese Schätzwerte auch eindeutig.

Neben der zunächst im Vordergrund stehenden Frage, wann der (G)EM-Algorithmus konvergiert, ist hier wie bei grundsätzlich allen iterativ ablaufenden Algorithmen von besonderem Interesse, wie rasch sich die iterativ berechneten Größen im Verlauf der Iterationen ihrem Grenzwert annähern. In der Regel entscheidet diese sogenannte *Konvergenzgeschwindigkeit* darüber, welcher Algorithmus aus einer Reihe geeigneter Algorithmen für eine bestimmte Aufgabe ausgewählt wird. In diesem Abschnitt wird daher zunächst ein Maß für die Konvergenzgeschwindigkeit des EM-Algorithmus entwickelt und schließlich werden allgemeine Aussagen über die Konvergenzgeschwindigkeit des (G)EM-Algorithmus getroffen. Die Ausführungen entsprechen denen in [DEMPSTER ET AL. 1968].

Die Angabe eines Maßes für die Konvergenzgeschwindigkeit beschränkt sich auf den EM-Algorithmus, was folgenden Grund hat: Beim EM-Algorithmus müssen die neuen Schätzwerte  $\hat{\beta}^{(i+1)}$  jeder Iteration  $i$  die Kullback-Leibler-Statistik  $Q(\hat{\beta}^{(i+1)} | \hat{\beta}^{(i)})$  maximieren. Da dies eine sehr scharfe Anforderung an die Neuschätzung  $\hat{\beta}^{(i+1)}$  darstellt, ist in jeder Iteration die Menge der als neue Schätzung in Frage kommenden Vektoren  $\beta$  relativ klein. Somit ist ausgehend vom Startwert  $\beta^{(0)}$  die Variationsbreite möglicher Realisierungen der Parameterfolge  $\{\hat{\beta}^{(i)}\}_{i \geq 0}$  relativ eng und in erster Näherung wird die Konvergenzgeschwindigkeit aller dieser Realisierungen die gleiche sein. Im Gegensatz dazu wird beim GEM-Algorithmus in der  $i$ -ten Iteration von der neuen Schätzung  $\hat{\beta}^{(i+1)}$  lediglich verlangt, daß sie keine geringere Kullback-Leibler-Statistik liefert als die Schätzung  $\hat{\beta}^{(i)}$



Offensichtlich wird die Geschwindigkeit, mit der sich die Schätzungen  $\hat{\beta}^{(i)}$  dem Grenzwert  $\beta^*$  annähern, für genügend große  $i$  maßgeblich von der Jacobi-Matrix  $\mathbf{DA}(\beta)|_{\beta=\beta^*} = \mathbf{DA}(\beta^*)$  bestimmt. Für die Norm der Differenz  $\hat{\beta}^{(i+1)} - \beta^*$  ergibt sich zunächst

$$\|\hat{\beta}^{(i+1)} - \beta^*\| = \left\| \mathbf{DA}(\beta^*)(\hat{\beta}^{(i)} - \beta^*) \right\|. \quad (324.6)$$

Aus dem Ausdruck der rechten Seite dieser Gleichung soll nun die Norm  $\|\hat{\beta}^{(i)} - \beta^*\|$  separiert werden. Hierzu muß die zur euklidischen Vektornorm zugehörige *Matrixnorm* eingeführt werden. Nach [KERNER ET AL. 1995] ist eine *Matrixnorm* eine auf einer Matrix definierte Norm und jede Vektornorm  $\|\mathbf{b}_{o \times 1}\|$  besitzt eine *zugehörige Matrixnorm*  $\|\mathbf{B}\|_M$  für quadratische Matrizen  $\mathbf{B}$  der Dimension  $o$ ,

$$\|\mathbf{B}\|_M = \sup_{\mathbf{x} \neq \mathbf{0}} \frac{\|\mathbf{B}\mathbf{x}\|}{\|\mathbf{x}\|}. \quad (324.7)$$

Die der euklidischen Vektornorm zugehörige Matrixnorm wird als *Spektralnorm* bezeichnet. Für sie gilt (vgl. [KERNER ET AL. 1995])

$$\|\mathbf{B}\|_{M, \text{euklid.}} = \sqrt{\lambda_{\max}(\mathbf{B}^T \mathbf{B})} =: \sqrt{\lambda_{\max}^{(\mathbf{B}^T \mathbf{B})}}, \quad (324.8)$$

wobei  $\lambda_{\max}^{(\mathbf{B}^T \mathbf{B})}$  den größten Eigenwert der Matrix  $\mathbf{B}^T \mathbf{B}$  bedeutet. Im folgenden wird bei der Matrixnorm auf die Indizes  $M$  und *euklid.* verzichtet, da im konkreten Fall zwischen der Matrix- und der Vektornorm anhand der Darstellung von Matrizen mittels Großbuchstaben und von Vektoren mittels Kleinbuchstaben unterschieden werden kann und weil hier ausschließlich die Spektralnorm verwendet wird. Wegen der Definition des *Supremums* als kleinste obere Schranke gilt

$$\|\mathbf{B}\| = \sup_{\mathbf{x} \neq \mathbf{0}} \frac{\|\mathbf{B}\mathbf{x}\|}{\|\mathbf{x}\|} \geq \frac{\|\mathbf{B}\mathbf{x}\|}{\|\mathbf{x}\|} \quad \text{und damit} \quad \|\mathbf{B}\| \cdot \|\mathbf{x}\| = \sqrt{\lambda_{\max}^{(\mathbf{B}^T \mathbf{B})}} \cdot \|\mathbf{x}\| \geq \|\mathbf{B}\mathbf{x}\|. \quad (324.9)$$

Überträgt man dieses Ergebnis auf Gl. (324.6), so ergibt sich

$$\begin{aligned} \|\hat{\beta}^{(i+1)} - \beta^*\| &= \left\| \mathbf{DA}(\beta^*)(\hat{\beta}^{(i)} - \beta^*) \right\| \\ &\leq \|\mathbf{DA}(\beta^*)\| \cdot \|\hat{\beta}^{(i)} - \beta^*\| = \sqrt{\lambda_{\max}^{([\mathbf{DA}(\beta^*)]^T \mathbf{DA}(\beta^*))}} \cdot \|\hat{\beta}^{(i)} - \beta^*\| \end{aligned} \quad (324.10)$$

und als Maß für die Konvergenzgeschwindigkeit des EM-Algorithmus erhält man die Quadratwurzel des größten Eigenvektors der Matrix  $\mathbf{DA}(\beta^*)$ ,

$$v = \frac{\|\hat{\beta}^{(i+1)} - \beta^*\|}{\|\hat{\beta}^{(i)} - \beta^*\|} = \sqrt{\lambda_{\max}^{([\mathbf{DA}(\beta^*)]^T \mathbf{DA}(\beta^*))}}. \quad (324.11)$$

(Vergleiche hierzu [DEMPSTER ET AL. 1968] und [HORNEGGER 1996].)<sup>13</sup>

Es es gilt also mit  $\sqrt{\lambda_{\max}^{([\mathbf{DA}(\beta^*)]^T \mathbf{DA}(\beta^*))}} = \lambda_{\max}^{\mathbf{DA}(\beta^*)}$  die Beziehung

$$v = \lambda_{\max}^{\mathbf{DA}(\beta^*)}. \quad (324.12)$$

<sup>13</sup>Im Rahmen der Literaturrecherche zur vorliegenden Arbeit wurde weder der Nachweis dieser Aussage noch ein Gegenbeispiel hierzu gefunden. Auf die Richtigkeit der in [DEMPSTER ET AL. 1968] und [HORNEGGER 1996] getroffenen Aussage, daß als Maß für die Konvergenzgeschwindigkeit der größte Eigenwert der Matrix  $\mathbf{DA}(\beta^*)$  wird dennoch vertraut.

Ist ein Eigenwert der Matrix  $\mathbf{DA}(\boldsymbol{\beta}^*)$  größer als eins, so kann es sich nach den Ausführungen in [DEMPSTER ET AL. 1968] bei  $\boldsymbol{\beta}^*$  um einen Sattelpunkt der logarithmierten Likelihoodfunktion handeln: Wird im Rahmen der EM-Iterationen für einen Startwert  $\hat{\boldsymbol{\beta}}^{(0)}$  als Schätzung ein Sattelpunkt  $\hat{\boldsymbol{\beta}}^*$  der logarithmierten Likelihoodfunktion erhalten so bleibt die Schätzung  $\hat{\boldsymbol{\beta}}$  in nachfolgenden Iterationen in diesem Sattelpunkt. Wird jedoch eine Schätzung  $\boldsymbol{\beta}'$  erhalten, die auch nur minimal von  $\boldsymbol{\beta}^*$  abweicht, so bewegen sich nachfolgende Schätzungen wegen (324.11) vom Sattelpunkt fort.

Hat  $\mathbf{DA}(\boldsymbol{\beta}^*)$  einen mit eins identischen Eigenwert, so weist dies darauf hin, daß  $\boldsymbol{\beta}^*$  auf einer Kurve im  $\mathbb{R}^u$  liegt, deren Punkte allesamt Maxima der logarithmierten Likelihoodfunktion darstellen (vgl. [DEMPSTER ET AL. 1968]).

Der nachfolgende Satz gibt Aufschluß darüber, wie die Matrix  $\mathbf{DA}(\boldsymbol{\beta}^*)|_{\boldsymbol{\beta}=\boldsymbol{\beta}^*}$  praktisch berechnet werden kann:

**Satz:** Es sei  $\{\hat{\boldsymbol{\beta}}^{(i)}\}_{i \geq 0}$  die Parameterfolge eines (G)EM-Algorithmus mit den Eigenschaften

- 1.)  $\{\hat{\boldsymbol{\beta}}^{(i)}\}_{i \geq 0}$  konvergiert gegen einen Punkt und  $\boldsymbol{\beta}^*$  im Innern des Parameterraums  $\mathcal{B}$
- 2.)  $\mathbf{D}^{10}Q(\hat{\boldsymbol{\beta}}^{(i+1)} | \hat{\boldsymbol{\beta}}^{(i)}) = \mathbf{0}$

Dann gilt

$$\mathbf{DL}(\boldsymbol{\beta}^*) = \mathbf{0}, \quad (324.13)$$

und

$$\mathbf{DA}(\hat{\boldsymbol{\beta}}^*) = [\mathbf{D}^{20}Q(\boldsymbol{\beta}^* | \boldsymbol{\beta}^*)]^{-1} \mathbf{D}^{20}H(\boldsymbol{\beta}^* | \boldsymbol{\beta}^*). \quad (324.14)$$

BEWEIS: Nach (231.7) gilt  $L(\hat{\boldsymbol{\beta}}^{(i+1)}) = Q(\hat{\boldsymbol{\beta}}^{(i+1)} | \hat{\boldsymbol{\beta}}^{(i)}) - H(\hat{\boldsymbol{\beta}}^{(i+1)} | \hat{\boldsymbol{\beta}}^{(i)})$ . Somit gilt für den Gradienten  $\mathbf{DL}(\hat{\boldsymbol{\beta}}^{(i+1)})$  der logarithmierten Likelihoodfunktion in Abhängigkeit von  $\hat{\boldsymbol{\beta}}^{(i+1)}$

$$\begin{aligned} \mathbf{DL}(\hat{\boldsymbol{\beta}}^{(i+1)}) &= \frac{\partial L(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}} \Big|_{\hat{\boldsymbol{\beta}}^{(i+1)}} = \frac{\partial Q(\boldsymbol{\beta} | \hat{\boldsymbol{\beta}}^{(i)})}{\partial \boldsymbol{\beta}} \Big|_{\hat{\boldsymbol{\beta}}^{(i+1)}} - \frac{\partial H(\boldsymbol{\beta} | \hat{\boldsymbol{\beta}}^{(i)})}{\partial \boldsymbol{\beta}} \\ &= \underbrace{\mathbf{D}^{10}Q(\hat{\boldsymbol{\beta}}^{(i+1)} | \hat{\boldsymbol{\beta}}^{(i)})}_{=0 \text{ nach Vorauss.2.})} - \mathbf{D}^{10}H(\hat{\boldsymbol{\beta}}^{(i+1)} | \hat{\boldsymbol{\beta}}^{(i)}) \\ &= -\mathbf{D}^{10}H(\hat{\boldsymbol{\beta}}^{(i+1)} | \hat{\boldsymbol{\beta}}^{(i)}) \end{aligned} \quad (324.15)$$

und nach Übergang auf den Grenzwert für  $i \rightarrow \infty$ , also für unendlich viele Iterationen folgt

$$\mathbf{DL}(\boldsymbol{\beta}^*) = \lim_{i \rightarrow \infty} \mathbf{DL}(\hat{\boldsymbol{\beta}}^{(i+1)}) = - \lim_{i \rightarrow \infty} \mathbf{D}^{10}H(\hat{\boldsymbol{\beta}}^{(i+1)} | \hat{\boldsymbol{\beta}}^{(i)}) = \mathbf{0}, \quad (324.16)$$

da nach (323.12)  $\lim_{i \rightarrow \infty} \mathbf{D}^{10}H(\hat{\boldsymbol{\beta}}^{(i+1)} | \hat{\boldsymbol{\beta}}^{(i)}) = \mathbf{0}$  gilt. Somit folgt die erste Aussage  $\mathbf{DL}(\boldsymbol{\beta}^*) = \mathbf{0}$ .

Zum Beweis der zweiten Aussage wird die von den beiden Vektoren  $\hat{\boldsymbol{\beta}}^{(i)}$  und  $\hat{\boldsymbol{\beta}}^{(i+1)}$  abhängige Größe  $\mathbf{D}^{10}Q(\hat{\boldsymbol{\beta}}^{(i+1)} | \hat{\boldsymbol{\beta}}^{(i)})$  in einer Taylorreihe um  $\boldsymbol{\beta}^*$  entwickelt. Bekanntlich gilt für eine vom Vektor  $\mathbf{u}$  abhängige skalarwertige Funktion  $b(\mathbf{u})$  (vgl. [HELFRICH II 1996])

$$b(\mathbf{u}_0 + \mathbf{h}) = b(\mathbf{u}_0) + \mathbf{h}^T \nabla b(\mathbf{u})|_{\mathbf{u}_0} + \frac{1}{2} \mathbf{h}^T H(\mathbf{u})|_{\mathbf{u}_0 + \theta \mathbf{h}} \mathbf{h} \quad \text{mit } 0 \leq \theta \leq 1, \quad (324.17)$$

wobei  $\nabla b(\mathbf{u})|_{\mathbf{u}_0}$  den Gradienten der Funktion  $b$  an der Stelle  $\mathbf{u}_0$  und  $H(\mathbf{u})|_{\mathbf{u}_0 + \theta \mathbf{h}}$  die Hesse-Matrix der Funktion  $b$  für einen Punkt  $\mathbf{u}_0 + \theta \mathbf{h}$  auf der geradlinigen Verbindung zwischen

$\mathbf{u}_0$  und  $\mathbf{u}_0 + \mathbf{h}$  bedeuten.

Setzt sich der Vektor  $\mathbf{u}$  aus den beiden Vektoren  $\mathbf{u}_1$  und  $\mathbf{u}_2$  zusammen, also  $\mathbf{u} = (\mathbf{u}_1^T, \mathbf{u}_2^T)^T$  und entsprechend  $\mathbf{u}_0 = (\mathbf{u}_{10}^T, \mathbf{u}_{20}^T)^T$  und  $\mathbf{h} = (\mathbf{h}_1, \mathbf{h}_2)^T$ , so folgt hieraus

$$b\left(\begin{bmatrix} \mathbf{u}_{10} \\ \mathbf{u}_{20} \end{bmatrix} + \begin{bmatrix} \mathbf{h}_1 \\ \mathbf{h}_2 \end{bmatrix}\right) = b\left(\begin{bmatrix} \mathbf{u}_{10} \\ \mathbf{u}_{20} \end{bmatrix}\right) + \begin{bmatrix} \mathbf{h}_1 \\ \mathbf{h}_2 \end{bmatrix}^T \nabla b\left(\begin{bmatrix} \mathbf{u}_1 \\ \mathbf{u}_2 \end{bmatrix}\right) \Big|_{\begin{bmatrix} \mathbf{u}_{10} \\ \mathbf{u}_{20} \end{bmatrix}} \\ + \frac{1}{2} \begin{bmatrix} \mathbf{h}_1 \\ \mathbf{h}_2 \end{bmatrix}^T H\left(\begin{bmatrix} \mathbf{u}_1 \\ \mathbf{u}_2 \end{bmatrix}\right) \Big|_{\begin{bmatrix} \mathbf{u}_{10} \\ \mathbf{u}_{20} \end{bmatrix} + \theta \begin{bmatrix} \mathbf{h}_1 \\ \mathbf{h}_2 \end{bmatrix}} \begin{bmatrix} \mathbf{h}_1 \\ \mathbf{h}_2 \end{bmatrix}. \quad (324.18)$$

Angewendet auf die  $k$ -te Komponente  $q_k^{10}(\hat{\beta}^{(i+1)} | \hat{\beta}^{(i)})$  des Vektors  $\mathbf{D}^{10}Q(\hat{\beta}^{(i+1)} | \hat{\beta}^{(i)}) = (q_1^{10}(\hat{\beta}^{(i+1)} | \hat{\beta}^{(i)}), q_2^{10}(\hat{\beta}^{(i+1)} | \hat{\beta}^{(i)}), \dots, q_u^{10}(\hat{\beta}^{(i+1)} | \hat{\beta}^{(i)}))^T$  liefert diese Gleichung

$$q_k^{10}(\hat{\beta}^{(i+1)} | \hat{\beta}^{(i)}) = q_k^{10}(\beta^* | \beta^*) + \begin{bmatrix} \hat{\beta}^{(i+1)} - \beta^* \\ \hat{\beta}^{(i)} - \beta^* \end{bmatrix}^T \begin{bmatrix} \frac{\partial q_k^{10}(\hat{\beta}^{(i+1)} | \hat{\beta}^{(i)})}{\partial \hat{\beta}^{(i+1)}} \Big|_{(\beta^* | \beta^*)} \\ \frac{\partial q_k^{10}(\hat{\beta}^{(i+1)} | \hat{\beta}^{(i)})}{\partial \hat{\beta}^{(i)}} \Big|_{(\beta^* | \beta^*)} \end{bmatrix} \\ + \frac{1}{2} \begin{bmatrix} \hat{\beta}^{(i+1)} - \beta^* \\ \hat{\beta}^{(i)} - \beta^* \end{bmatrix}^T \mathbf{H}_k \begin{bmatrix} \hat{\beta}^{(i+1)} - \beta^* \\ \hat{\beta}^{(i)} - \beta^* \end{bmatrix} \quad (324.19)$$

wobei die konkrete Gestalt der Hessematrix  $\mathbf{H}_k$  für die weiteren Überlegungen irrelevant ist. Nach Voraussetzung 1.) verschwindet in dieser Gleichung die linke Seite und wegen  $\mathbf{D}^{10}Q(\beta^* | \beta^*) = \lim_{i \rightarrow \infty} \mathbf{D}^{10}Q(\hat{\beta}^{(i+1)} | \hat{\beta}^{(i)}) = \mathbf{0}$  auch der erste Term der rechten Seite. Mit den Abkürzungen

$$\frac{\partial q_k^{10}(\hat{\beta}^{(i+1)} | \hat{\beta}^{(i)})}{\partial \hat{\beta}^{(i+1)}} =: \mathbf{q}_k^{20}(\hat{\beta}^{(i+1)} | \hat{\beta}^{(i)}) \quad \text{und} \quad \frac{\partial q_k^{10}(\hat{\beta}^{(i+1)} | \hat{\beta}^{(i)})}{\partial \hat{\beta}^{(i)}} =: \mathbf{q}_k^{11}(\hat{\beta}^{(i+1)} | \hat{\beta}^{(i)}),$$

folgt

$$0 = \begin{bmatrix} \hat{\beta}^{(i+1)} - \beta^* \\ \hat{\beta}^{(i)} - \beta^* \end{bmatrix}^T \begin{bmatrix} \mathbf{q}_k^{20}(\beta^* | \beta^*) \\ \mathbf{q}_k^{11}(\beta^* | \beta^*) \end{bmatrix} + \frac{1}{2} \begin{bmatrix} \hat{\beta}^{(i+1)} - \beta^* \\ \hat{\beta}^{(i)} - \beta^* \end{bmatrix}^T \mathbf{H}_k \begin{bmatrix} \hat{\beta}^{(i+1)} - \beta^* \\ \hat{\beta}^{(i)} - \beta^* \end{bmatrix} \quad (324.20)$$

Wegen  $\hat{\beta}^{(i+1)} = \mathbf{A}(\hat{\beta}^{(i)})$  und  $\beta^* = \mathbf{A}(\beta^*)$  gilt

$$\hat{\beta}^{(i+1)} - \beta^* = \mathbf{A}(\hat{\beta}^{(i)}) - \mathbf{A}(\beta^*)$$

und mit der Taylor-Entwicklung der Funktion  $\mathbf{A}$

$$\mathbf{A}(\hat{\beta}^{(i)}) = \mathbf{A}(\beta^*) + \underbrace{\frac{\partial \mathbf{A}(\beta)}{\partial \beta} \Big|_{\beta^*}}_{=\mathbf{DA}(\beta^*)} (\hat{\beta}^{(i)} - \beta^*)$$

unter Vernachlässigung von Termen der Ordnung 2 und größer gilt

$$\begin{bmatrix} \hat{\beta}^{(i+1)} - \beta^* \\ \hat{\beta}^{(i)} - \beta^* \end{bmatrix} = \begin{bmatrix} \mathbf{DA}(\beta^*) (\hat{\beta}^{(i)} - \beta^*) \\ (\hat{\beta}^{(i)} - \beta^*) \end{bmatrix} = \begin{bmatrix} \mathbf{DA}(\beta^*) \\ \mathbf{I} \end{bmatrix} (\hat{\beta}^{(i)} - \beta^*).$$

Einsetzen in (324.20) liefert

$$0 = (\hat{\beta}^{(i)} - \beta^*)^T \begin{bmatrix} \mathbf{DA}(\beta^*) \\ \mathbf{I} \end{bmatrix}^T \begin{bmatrix} \mathbf{q}_k^{20}(\beta^* | \beta^*) \\ \mathbf{q}_k^{11}(\beta^* | \beta^*) \end{bmatrix} \\ + \frac{1}{2} (\hat{\beta}^{(i)} - \beta^*)^T \underbrace{\begin{bmatrix} \mathbf{DA}(\beta^*) \\ \mathbf{I} \end{bmatrix}^T \mathbf{H}_k \begin{bmatrix} \mathbf{DA}(\beta^*) \\ \mathbf{I} \end{bmatrix}}_{:=\mathbf{C}} (\hat{\beta}^{(i)} - \beta^*) \quad (324.21)$$

und nach Ausklammern von  $(\hat{\beta}^{(i)} - \beta^*)^T$  gilt

$$0 = (\hat{\beta}^{(i)} - \beta^*)^T \left\{ \begin{bmatrix} \mathbf{DA}(\beta^*) \\ \mathbf{I} \end{bmatrix}^T \begin{bmatrix} \mathbf{q}_k^{20}(\beta^* | \beta^*) \\ \mathbf{q}_k^{11}(\beta^* | \beta^*) \end{bmatrix} + \frac{1}{2} \mathbf{C}(\hat{\beta}^{(i)} - \beta^*) \right\} \quad (324.22)$$

Da  $\hat{\beta}^{(i+1)} - \beta^*$  i.a. nicht der Nullvektor ist, muß der Ausdruck innerhalb der Klammer  $\{\}$  der Nullvektor sein. Es folgt also

$$\mathbf{0} = \begin{bmatrix} \mathbf{DA}(\beta^*) \\ \mathbf{I} \end{bmatrix}^T \begin{bmatrix} \mathbf{q}_k^{20}(\beta^* | \beta^*) \\ \mathbf{q}_k^{11}(\beta^* | \beta^*) \end{bmatrix} + \frac{1}{2} \mathbf{C}(\hat{\beta}^{(i)} - \beta^*) \quad (324.23)$$

und durch Übergang auf den Grenzwert für  $i \rightarrow \infty$ , bei dem der zweite Term wegen  $\lim_{i \rightarrow \infty} (\hat{\beta}^{(i)} - \beta^*) = 0$  verschwindet

$$\mathbf{0} = \begin{bmatrix} \mathbf{DA}(\beta^*) \\ \mathbf{I} \end{bmatrix}^T \begin{bmatrix} \mathbf{q}_k^{20}(\beta^* | \beta^*) \\ \mathbf{q}_k^{11}(\beta^* | \beta^*) \end{bmatrix} = \mathbf{DA}(\beta^*)^T \mathbf{q}_k^{20}(\beta^* | \beta^*) + \mathbf{q}_k^{11}(\beta^* | \beta^*). \quad (324.24)$$

Hieraus ergibt sich

$$\mathbf{0} = \mathbf{DA}(\beta^*)^T [\mathbf{q}_1^{20}(\beta^* | \beta^*), \dots, \mathbf{q}_u^{20}(\beta^* | \beta^*)] + [\mathbf{q}_1^{11}(\beta^* | \beta^*), \dots, \mathbf{q}_u^{11}(\beta^* | \beta^*)],$$

wobei  $\mathbf{0}$  die Nullmatrix bedeutet. Mit den Jacobi-Matrizen

$$\mathbf{D}^{20}Q(\hat{\beta}^{(i+1)} | \hat{\beta}^{(i)}) = \begin{bmatrix} \mathbf{q}_1^{20}(\hat{\beta}^{(i+1)} | \hat{\beta}^{(i)})^T \\ \vdots \\ \mathbf{q}_u^{20}(\hat{\beta}^{(i+1)} | \hat{\beta}^{(i)})^T \end{bmatrix} \quad (324.25)$$

bzw.

$$\mathbf{D}^{11}Q(\hat{\beta}^{(i+1)} | \hat{\beta}^{(i)}) = \begin{bmatrix} \mathbf{q}_1^{11}(\hat{\beta}^{(i+1)} | \hat{\beta}^{(i)})^T \\ \vdots \\ \mathbf{q}_u^{11}(\hat{\beta}^{(i+1)} | \hat{\beta}^{(i)})^T \end{bmatrix}, \quad (324.26)$$

die die 1. Ableitungen des Gradienten  $\mathbf{D}^{10}Q(\hat{\beta}^{(i+1)} | \hat{\beta}^{(i)})$  der Kullback-Leibler-Statistik nach den Elementen des ersten Argumentes  $\hat{\beta}^{(i+1)}$  bzw. nach den Elementen des 2. Argumentes  $\hat{\beta}^{(i)}$  enthalten, ergibt sich

$$\mathbf{0} = \mathbf{DA}(\beta^*)^T \mathbf{D}^{20}Q(\beta^* | \beta^*)^T + \mathbf{D}^{11}Q(\beta^* | \beta^*)^T \quad (324.27)$$

und daraus

$$\mathbf{0} = \mathbf{D}^{20}Q(\beta^* | \beta^*) \mathbf{DA}(\beta^*) + \mathbf{D}^{11}Q(\beta^* | \beta^*). \quad (324.28)$$

Schließlich folgt

$$\mathbf{DA}(\beta^*) = - [\mathbf{D}^{20}Q(\beta^* | \beta^*)]^{-1} \mathbf{D}^{11}Q(\beta^* | \beta^*) \quad (324.29)$$

Nach (321.5) gilt  $L(\hat{\beta}^{(i)}) = Q(\hat{\beta}^{(i)} | \hat{\beta}^{(i)}) - H(\hat{\beta}^{(i)} | \hat{\beta}^{(i)}) = Q(\hat{\beta}^{(i)} | \hat{\beta}^{(i-1)}) - H(\hat{\beta}^{(i)} | \hat{\beta}^{(i-1)})$ , so daß in der  $i$ -ten Iteration  $L(\hat{\beta}^{(i)})$  unabhängig von der vorherigen Schätzung  $\beta^{(i-1)}$  ist. Es gilt also  $\mathbf{D}^{11}L(\hat{\beta}^{(i)}) = \mathbf{0}$  und aus  $L(\beta^*) = Q(\beta^* | \beta^*) - H(\beta^* | \beta^*)$  folgt

$$\mathbf{0} = \mathbf{D}^{11}L(\beta^*) = \mathbf{D}^{11}Q(\beta^* | \beta^*) - \mathbf{D}^{11}H(\beta^* | \beta^*), \quad (324.30)$$

so daß gilt

$$\mathbf{D}^{11}Q(\beta^* | \beta^*) = \mathbf{D}^{11}H(\beta^* | \beta^*). \quad (324.31)$$

Für das Element  $[\mathbf{D}^{11}H(\boldsymbol{\psi} | \boldsymbol{\phi})]_{kl}$  der  $k$ -ten Zeile und der  $l$ -ten Spalte der Matrix  $\mathbf{D}^{11}H(\boldsymbol{\psi} | \boldsymbol{\phi})$  gilt mit  $\boldsymbol{\psi} = (\psi_1, \dots, \psi_u)^T$  und  $\boldsymbol{\phi} = (\phi_1, \dots, \phi_u)^T$

$$\begin{aligned} [\mathbf{D}^{11}H(\boldsymbol{\psi} | \boldsymbol{\phi})]_{kl} &= \frac{\partial}{\partial \psi_k} \left( \frac{\partial}{\partial \phi_l} \int \dots \int_{\mathcal{Z}(\mathbf{y})} \log f(z | \mathbf{y}, \boldsymbol{\psi}) f(z | \mathbf{y}, \boldsymbol{\phi}) dz \right) \\ &= \frac{\partial}{\partial \psi_k} \int \dots \int_{\mathcal{Z}(\mathbf{y})} \log f(z | \mathbf{y}, \boldsymbol{\psi}) \frac{\partial f(z | \mathbf{y}, \boldsymbol{\phi})}{\partial \phi_l} dz \\ &= \int \dots \int_{\mathcal{Z}(\mathbf{y})} \frac{1}{f(z | \mathbf{y}, \boldsymbol{\psi})} \frac{\partial f(z | \mathbf{y}, \boldsymbol{\psi})}{\partial \psi_l} \frac{\partial f(z | \mathbf{y}, \boldsymbol{\phi})}{\partial \phi_k} dz. \end{aligned} \quad (324.32)$$

Entsprechend gilt für das Element  $[\mathbf{D}^{20}H(\boldsymbol{\psi} | \boldsymbol{\phi})]_{kl}$  der  $k$ -ten Zeile und der  $l$ -ten Spalte der Matrix  $\mathbf{D}^{20}H(\boldsymbol{\psi} | \boldsymbol{\phi})$

$$\begin{aligned} [\mathbf{D}^{20}H(\boldsymbol{\psi} | \boldsymbol{\phi})]_{kl} &= \frac{\partial}{\partial \psi_k} \left( \frac{\partial}{\partial \psi_l} \int \dots \int_{\mathcal{Z}(\mathbf{y})} \log f(z | \mathbf{y}, \boldsymbol{\psi}) f(z | \mathbf{y}, \boldsymbol{\phi}) dz \right) \\ &= \frac{\partial}{\partial \psi_k} \int \dots \int_{\mathcal{Z}(\mathbf{y})} \frac{1}{f(z | \mathbf{y}, \boldsymbol{\psi})} \frac{\partial f(z | \mathbf{y}, \boldsymbol{\psi})}{\partial \psi_l} f(z | \mathbf{y}, \boldsymbol{\phi}) dz \\ &= \int \dots \int_{\mathcal{Z}(\mathbf{y})} - \frac{1}{f(z | \mathbf{y}, \boldsymbol{\psi})^2} \frac{\partial f(z | \mathbf{y}, \boldsymbol{\psi})}{\partial \psi_k} \frac{\partial f(z | \mathbf{y}, \boldsymbol{\psi})}{\partial \psi_k} f(z | \mathbf{y}, \boldsymbol{\phi}) dz \\ &\quad + \underbrace{\int \dots \int_{\mathcal{Z}(\mathbf{y})} \frac{1}{f(z | \mathbf{y}, \boldsymbol{\psi})} \frac{\partial^2 f(z | \mathbf{y}, \boldsymbol{\psi})}{\partial \psi_k \partial \psi_l} f(z | \mathbf{y}, \boldsymbol{\phi}) dz}_{= \frac{\partial^2}{\partial \psi_k \partial \psi_l} \int \dots \int_{\mathcal{Z}(\mathbf{y})} f(z | \mathbf{y}, \boldsymbol{\psi}) dz = \frac{\partial^2}{\partial \psi_k \partial \psi_l} 1 = 0} \\ &= - \int \dots \int_{\mathcal{Z}(\mathbf{y})} \frac{1}{f(z | \mathbf{y}, \boldsymbol{\psi})} \frac{\partial f(z | \mathbf{y}, \boldsymbol{\psi})}{\partial \psi_k} \frac{\partial f(z | \mathbf{y}, \boldsymbol{\psi})}{\partial \psi_k} dz. \end{aligned} \quad (324.33)$$

Aus (324.32) und (324.33) folgt

$$\mathbf{D}^{11}H(\boldsymbol{\beta}^* | \boldsymbol{\beta}^*) = -\mathbf{D}^{20}H(\boldsymbol{\beta}^* | \boldsymbol{\beta}^*) \quad (324.34)$$

Durch Einsetzen von (324.31) und (324.34) in (324.29) erhält man schließlich

$$\mathbf{D}\mathbf{A}(\boldsymbol{\beta}^*) = [\mathbf{D}^{20}Q(\boldsymbol{\beta}^* | \boldsymbol{\beta}^*)]^{-1} \mathbf{D}^{20}H(\boldsymbol{\beta}^* | \boldsymbol{\beta}^*)$$

und damit die zweite Aussage.  $\square$

In der Literatur zum (G)EM-Algorithmus wird dem EM-Algorithmus eine relativ langsame Konvergenz zugeschrieben (vgl. [DEMPSTER ET AL. 1968], [HORNEGGER 1996], [REDNER 1984]. Wie in [DEMPSTER ET AL. 1968] angedeutet, kann gezeigt werden, dass der EM-Algorithmus um so schneller konvergiert, je weniger fehlende Beobachtungen in die Ausgleichung eingeführt werden. Der GEM-Algorithmus konvergiert wegen der weniger strengen Anforderungen an die Schätzung  $\hat{\boldsymbol{\beta}}^{(i)}$  im M-Schritt noch langsamer als der EM-Algorithmus. Besonders dann, wenn sich die Kullback-Leibler-Statistik eines Parameterschätzproblems nicht analytisch angeben lässt, so dass sie im Rahmen der EM-Iterationen mittels numerischer Verfahren berechnet werden muß, ergibt sich ein enormer Rechenaufwand. Insbesondere in einem solchen Fall lohnt es sich vor dem Einsatz des



EM-Algorithmus sorgfältig zu prüfen, ob das gewünschte Ergebnis nicht schneller mit einem anderen Algorithmus erhalten werden kann.

Oft kann die Kullback-Leibler-Statistik einer Parameterschätzung aus unvollständigen Daten jedoch analytisch angegeben werden, so daß sich der im Zusammenhang mit dem EM-Algorithmus anfallende Rechenaufwand stark reduziert. Es bleibt allerdings dann immer noch die Tatsache, daß der EM-Algorithmus i.d.R. aufgrund seiner langsamen Konvergenz sehr viele EM-Iterationen erforderlich macht.

### **33 Zusammenfassung**

In diesem Kapitel wurden die allgemeinen Eigenschaften des EM-Algorithmus diskutiert. Ausgehend von der Betrachtung des Verhaltens der logarithmierten Likelihoodfunktion im Rahmen der EM-Iterationen und dem Nachweis ihrer Konvergenz wurde erläutert, unter welchen Bedingungen die logarithmierte Likelihoodfunktion im Rahmen der (G)EM-Iterationen wie gewünscht maximiert wird bzw. unter welchen Bedingungen Konvergenz gegen einen stationären Punkt erfolgt. Im Anschluß daran wurden Bedingungen genannt, unter denen die Parameterfolge gegen einen stationären Punkt bzw. gegen ein lokales Maximum der logarithmierten Likelihoodfunktion konvergiert, sofern die Folge der Likelihoodwerte konvergiert. Zum Abschluß des Kapitels erfolgten einige Bemerkungen zur Konvergenzgeschwindigkeit des (G)EM-Algorithmus, wobei insbesondere auf die Langsamkeit des Verfahrens hingewiesen wurde.

Im folgenden Kapitel wird der EM-Algorithmus nun im Zusammenhang mit einer konkreten Aufgabe angewandt, um zu zeigen, wie der Algorithmus zur Parameterschätzung und Klassifikation eingesetzt werden kann.

## Übersicht 3.1: Zur Konvergenz der Parameterfolge:

1. Konvergiert die Folge  $\{L(\hat{\beta}^{(i)})\}_{i \geq 0}$  beim (G)EM-Algorithmus gemäß Satz (322.6) für  $i \rightarrow \infty$  gegen  $L^*$ , den Funktionswert  $L^* = L(\beta)$  eines stationären Punktes  $\beta \in \Gamma^S$  von  $L(\beta)$ , dann gilt für das Konvergenzverhalten der Parameterfolge  $\{\hat{\beta}^{(i)}\}_{i \geq 0}$  für  $i \rightarrow \infty$ :
  - (a) Besteht die Menge  $\Gamma^S(L^*) = \{\beta \in \Gamma^S : L(\beta) = L^*\}$  der stationären Punkte  $\beta$  der logarithmierten Likelihoodfunktion mit  $L(\beta) = L^*$  aus genau einem Punkt  $\beta^*$ , dann konvergiert  $\{\hat{\beta}^{(i)}\}_{i \geq 0}$  gegen diesen Punkt,  $\lim_{i \rightarrow \infty} \hat{\beta}^{(i)} = \hat{\beta}^*$ .
  - (b) Falls  $\|\hat{\beta}^{(i+1)} - \hat{\beta}^{(i)}\| \rightarrow \mathbf{0}$  für  $i \rightarrow \infty$  gilt und  $\Gamma^S(L^*)$  eine diskrete Menge ist, so konvergiert  $\hat{\beta}^{(i)}$  gegen ein  $\beta^*$  in  $\Gamma^S(L^*)$ . Ist  $\Gamma^S(L^*)$  nicht diskret, dann liegen zumindest alle Konvergenzpunkte  $\beta^*$  von  $\{\hat{\beta}^{(i)}\}_{i \geq 0}$  in einer zusammenhängenden und kompakten Teilmenge von  $\Gamma^S(L^*)$ .
2. Konvergiert die Folge  $\{L(\hat{\beta}^{(i)})\}_{i \geq 0}$  beim (G)EM-Algorithmus gemäß Satz (322.7) für  $i \rightarrow \infty$  gegen  $L^*$ , den Funktionswert  $L^* = L(\beta)$  eines lokalen Maximums  $\beta \in \Gamma^L$  von  $L(\beta)$ , dann konvergiert die Parameterfolge  $\{\hat{\beta}^{(i)}\}_{i \geq 0}$  für  $i \rightarrow \infty$  gegen ein lokales Maximum  $\beta^*$  von  $L(\beta)$ , falls eine der folgenden Bedingungen erfüllt ist:
  - (a) Die Menge  $\Gamma^L(L^*) = \{\beta \in \Gamma^S : L(\beta) = L^*\}$  der lokalen Maxima  $\beta$  mit  $L(\beta) = L^*$  besteht nur aus einem einzigen Punkt  $\beta^*$ .  $\{\hat{\beta}^{(i)}\}_{i \geq 0}$  konvergiert dann gegen diesen Punkt,  $\lim_{i \rightarrow \infty} \hat{\beta}^{(i)} = \hat{\beta}^*$ .
  - (b) Falls  $\|\hat{\beta}^{(i+1)} - \hat{\beta}^{(i)}\| \rightarrow \mathbf{0}$  für  $i \rightarrow \infty$  gilt und  $\Gamma^L(L^*)$  eine diskrete Menge ist, so konvergiert  $\hat{\beta}^{(i)}$  gegen ein  $\beta^*$  in  $\Gamma^L(L^*)$ . Ist  $\Gamma^L(L^*)$  nicht diskret, dann liegen zumindest alle Konvergenzpunkte  $\beta^*$  von  $\{\hat{\beta}^{(i)}\}_{i \geq 0}$  in einer zusammenhängenden und kompakten Teilmenge von  $\Gamma^L(L^*)$ .
3. Die Parameterfolge  $\{\hat{\beta}^{(i)}\}_{i \geq 0}$  des (G)EM-Algorithmus konvergiert gegen einen stationären Punkt  $\beta^*$  der logarithmierten Likelihoodfunktion, wenn in jeder Iteration  $\mathbf{D}^{10}Q(\hat{\beta}^{(i+1)} | \hat{\beta}^{(i)}) = \mathbf{0}$  gilt und  $\mathbf{D}^{10}Q(\psi | \phi)$  in  $\psi$  und  $\phi$  stetig ist und zugleich entweder
  - (a) die Menge  $\mathcal{L}(L^*)$  aus einem einzigen Element besteht, also  $\mathcal{L}(L^*) = \{\beta^*\}$  oder
  - (b)  $\lim_{i \rightarrow \infty} \|\hat{\beta}^{(i+1)} - \hat{\beta}^{(i)}\| = \mathbf{0}$  gilt und  $\mathcal{L}(L^*)$  diskret ist.

**Speziell für den EM-Algorithmus gilt:**

4. Die Parameterfolge  $\{\hat{\beta}^{(i)}\}_{i \geq 0}$  des EM-Algorithmus konvergiert gegen einen stationären Punkt  $\beta^*$ , falls die Kullback-Leibler-Statistik  $Q(\psi | \phi)$  in  $\psi$  und  $\phi$  stetig ist,  $\lim_{i \rightarrow \infty} \|\hat{\beta}^{(i+1)} - \hat{\beta}^{(i)}\| = \mathbf{0}$  gilt und  $\Gamma^S(L^*)$  eine diskrete Menge ist. Unabhängig von der Diskretheit von  $\Gamma^S(L^*)$  liegen alle Konvergenzpunkte  $\beta^*$  von  $\{\hat{\beta}^{(i)}\}_{i \geq 0}$  in einer zusammenhängenden und kompakten Teilmenge von  $\Gamma^S(L^*)$ .
5. Ist  $L(\beta)$  im Innern des Parameterraums  $\mathcal{B}$  eine unimodale Funktion mit dem lokalen Maximum  $\beta^*$  und ist  $\mathbf{D}^{10}Q(\psi | \phi)$  stetig in  $\psi$  und  $\phi$ , so konvergiert die Parameterfolge gegen das lokale Maximum.

## Kapitel 4

# Beispiel: Parameterschätzung und Klassifikation

Im folgenden soll anhand des Bereits in der Einleitung skizzierten, sehr einfachen Beispiels aufgezeigt werden, wie der EM-Algorithmus zur Parameterschätzung bei gleichzeitiger Lösung eines Zuordnungsproblems eingesetzt werden kann.

### 41 Aufgabe: Bestimmung der Parameter zweier ausgleichender Kurven in einem Bild

Gegeben sei ein CCD-Bild (d.i. ein Bild im Rasterformat), in dem zwei linienhafte Objekte  $\mathcal{O}_1$  und  $\mathcal{O}_2$  – hier eine Hochspannungsleitung und eine Straße – abgebildet sind (vgl. Abb. 41a). Zur Überführung des Bildes in ein Vektorformat seien die Bildkoordinaten  $u_i$  und  $y_i$  einzelner Punkte auf den Mittellinien<sup>1</sup> der beiden Objekte bestimmt und in der Datei `Linien.dat` abgespeichert worden. Als Beobachtungsmaterial liegt also eine Meßreihe der Form

$$\begin{array}{c|c|c|c|c} u_i & u_1 & u_2 & \dots & u_N \\ \hline y_i & y_1 & y_2 & \dots & y_N \end{array} \quad (410.1)$$

vor. In Abb. 41b sind die einzelnen Objektpunkte<sup>2</sup> in einem Koordinatensystem dargestellt.

Bei den Koordinatenmessungen seien die  $N$  Objektpunkte in willkürlicher Reihenfolge angemessen worden, so daß in der Datei `Linien.dat` die zu Objekt  $\mathcal{O}_1$  und  $\mathcal{O}_2$  gehörenden Koordinatenpaare zufällig aufeinander folgen. (Das bedeutet, daß es nicht möglich ist, aus der Reihenfolge der einzelnen Koordinatenpaare in der Koordinatenliste auf deren Zugehörigkeit zu einem der beiden Objekte zu schließen.) Bei den Messungen sei es weiter versäumt worden, festzuhalten, welches gemessene Koordinatenpaar zu welchem Objekt gehört. Für die weitere Auswertung ist demnach zwar bekannt, daß jedes gemessene Koordinatenpaar zu einem der beiden Objekte  $\mathcal{O}_1$  (Hochspannungsleitung) oder  $\mathcal{O}_2$  (Straße) gehört, es geht aus dem Datenmaterial jedoch nicht hervor, zu welchem von beiden (vgl. Anhang A2).

<sup>1</sup>Im Falle der Straße sind dies Punkte auf dem Mittelstreifen und im Falle der Hochspannungsleitung Punkte auf dem mittleren Leitungskabel.

<sup>2</sup>Wenn im folgenden von Objektpunkten oder Punkten von Objekten die Rede ist, so sind damit jeweils die Punkte auf den Mittellinien der beiden Objekte gemeint.

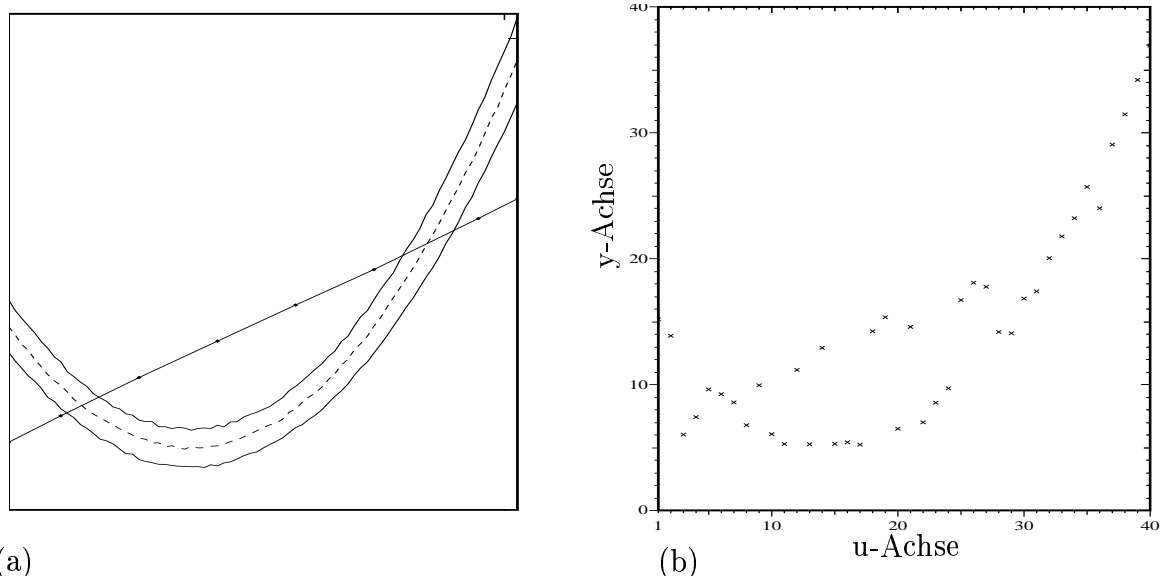


Abbildung 1: (a) Bild mit zwei linienhaften Objekten (Straße und Hochspannungsleitung). (b) Gemessene Punkte auf den Mittellinien der beiden Objekte. Zur Visualisierung der Messunsicherheiten wurde die Messwertstreuung in der Abbildung extrem groß gewählt, was insbesondere aus den zu der Hochspannungsleitung gehörenden Punkten ersichtlich ist.

Die *Aufgabe* besteht nun darin, eine Software zu entwickeln, die nach Vorgabe zweier Beobachtungsmodelle,<sup>3</sup> mit denen die Mittellinien der beiden Objekte  $\mathcal{O}_1$  und  $\mathcal{O}_2$  im Bild parametrisiert werden können, dazu in der Lage ist, nur aufgrund der in der Datei `Linien.dat` abgelegten Daten automatisch diejenigen Parameter der beiden Modelle zu schätzen, mit denen jeweils die optimale Approximation der Mittellinie der Hochspannungsleitung bzw. der Straße erreicht wird. Hierbei muß das Problem der Zuordnung der einzelnen Koordinatenpaare zu einem der beiden Objekte gelöst werden.

Um die nachfolgenden Betrachtungen möglichst einfach zu halten, wird hier davon ausgegangen, daß die Mittellinie der Hochspannungsleitung (Objekt  $\mathcal{O}_1$ ) durch eine Gerade und die Mittellinie der Straße (Objekt  $\mathcal{O}_2$ ) durch ein Polynom 2ten Grades approximiert werden können. (Prinzipiell können hier beliebige Funktionen gewählt werden, die zur Parametrisierung der beiden Mittellinien im Bild geeignet sind.) Weiter werden der Einfachheit halber die Abszissenwerte  $u_i$  der beobachteten Punkte als fest vorgegeben angesehen und nur die Ordinaten  $y_i$  als mit Meßunsicherheiten behaftete Beobachtungen betrachtet.

### Zur Eignung herkömmlicher Schätzverfahren für die Aufgabe

An dieser Stelle sei darauf hingewiesen, daß die oben formulierte Aufgabe mit herkömmlichen<sup>4</sup> Methoden der Parameterschätzung im Gauß Markoff Modell *nicht* ohne weiteres gelöst werden kann. Mit herkömmlichen Schätzmethoden können aus einem Satz von Be-

<sup>3</sup>Gemeint ist hier eine Modellierung der Beobachtungen im Sinne eines Gauß Markoff Modells bzw. eines Gauß Helmert Modells. (vgl. [KOCH 1998])

<sup>4</sup>Unter den herkömmlichen Schätzverfahren wird hier die Parameterschätzung im Gauß Markoff Modell nach der Methode der kleinsten Quadrate, der Maximum - Likelihood -Methode und nach der Methode der besten linearen erwartungstreuen Schätzung verstanden (vgl. [KOCH 1998])

obachtungsdaten jeweils nur Schätzwerte für die unbekannt Parameter eines einzigen Beobachtungsmodells abgeleitet werden; es ist nicht möglich, aus einem Beobachtungssatz gleichzeitig die Parameter mehrerer Modelle zu schätzen, von denen für jede Beobachtung jeweils nur eines zutrifft<sup>5</sup>. Übertragen auf die vorliegende Aufgabe bedeutet dies, daß mit herkömmlichen Schätzmethode aus den Beobachtungsdaten lediglich *entweder* Schätzwerte für die unbekannt Parameter der Gerade *oder* Schätzwerte für die unbekannt Parameter des Polynoms 2ten Grades bestimmt werden können. Sollen die Parameter *beider* Modelle geschätzt werden, so werden zwei getrennte Ausgleichungen erforderlich. Diese Ausgleichungen führen allerdings nur dann zu optimalen Approximationen der Mittellinien von  $\mathcal{O}_1$  und  $\mathcal{O}_2$ , wenn vorab bekannt ist, welche Beobachtungen in welche der beiden Ausgleichungen einzuführen sind, d.h. welche Koordinatenpaare zu welchem Objekt gehören; denn wird beispielsweise in die Schätzung der Geradenparameter (d.h. der Parameter zur Beschreibung der Mittellinie der Hochspannungsleitung  $\mathcal{O}_1$ ) eine Beobachtung (d.h. ein Koordinatenpaar) eingeführt, die tatsächlich zu der Straße (Objekt  $\mathcal{O}_2$ ) gehört, so stellt diese Beobachtung in der Ausgleichung zur Schätzung der Geradenparameter einen Ausreißer dar, der das Ergebnis der Ausgleichung verfälscht.

Es ist demnach wegen der Unkenntnis um die Zuordnung der gemessenen Koordinaten zu den beiden Objekten  $\mathcal{O}_1$  und  $\mathcal{O}_2$  nicht ohne weiteres möglich, die oben gestellte Aufgabe mit herkömmlichen Schätzmethode zu lösen, da das Zuordnungsproblem mit herkömmlichen Schätzverfahren zumindest nicht gleichzeitig mit der Parameterschätzung gelöst werden kann.

## 42 Lösung mittels des EM-Algorithmus:

Im folgenden wird aufgezeigt, wie die oben gestellte Aufgabe durch Anwendung des EM-Algorithmus gelöst werden kann. Dazu wird die Aufgabe zunächst als Schätzproblem mit unvollständigen Beobachtungsdaten formuliert. Anschließend wird auf dieses Schätzproblem dann der EM-Algorithmus angewandt.

### 421 Formulierung der Aufgabe als Schätzproblem aus unvollständigen Beobachtungsdaten

Zur Formulierung der Aufgabe als Schätzproblem aus unvollständigen Beobachtungsdaten werden zunächst im Vektor  $\mathbf{y}$  der zugänglichen Beobachtungen die  $y$ -Werte der Meßreihe (410.1) zusammengefasst:

$$\mathbf{y} = [y_1, y_2, \dots, y_N]^T \quad (421.1)$$

Um die folgenden Betrachtungen möglichst einfach zu halten, wird angenommen, daß die Beobachtungen  $y_i$  paarweise statistisch unabhängig voneinander sind, d.h. für die Kovarianz  $C(y_i, y_j)$  zweier Beobachtungen  $y_i$  und  $y_j$  gilt

$$C(y_i, y_j) = 0 \text{ für } i \neq j.$$

Weiter wird angenommen, daß die zur Hochspannungsleitung gehörenden Beobachtungen untereinander gleich genau sind und die Varianz  $\sigma_1^2$  besitzen. Ebenso werden die zur

<sup>5</sup>Man beachte, daß es sich bei der multivariaten Parameterschätzung (vgl. [KOCH 1998]) im Sinne einer Deformationsanalyse um eine mehrfache Bestimmung *derselben* Parameter handelt und nicht um die Bestimmung von Parametern mehrerer Modelle

Straße gehörenden Beobachtungen als untereinander gleich genau angenommenen, ihre Varianz beträgt  $\sigma_2^2$ .<sup>6</sup> Somit ergibt sich die Kovarianzmatrix  $\mathbf{D}(\mathbf{y})$  der Beobachtungen als Diagonalmatrix

$$\mathbf{D}(\mathbf{y}) = \text{diag} [V(y_1), V(y_2), \dots, V(y_N)]$$

$$\text{mit } V(y_i) = \begin{cases} \sigma_1^2 & \text{falls } y_i \text{ zu } \mathcal{O}_1 \text{ (Hochspannungsleitung) gehört,} \\ \sigma_2^2 & \text{falls } y_i \text{ zu } \mathcal{O}_2 \text{ (Straße) gehört.} \end{cases} \quad (421.2)$$

Die Beobachtungen  $y_i$  gehören jeweils entweder zum Modell „Gerade“ oder zum Modell „Polynom 2. Grades“. Zur Charakterisierung der Beobachtungen kommen also folgende Modellansätze im Sinne eines Gauß-Markoff-Modells in Frage:

**Modell 1 (Objekt  $\mathcal{O}_1$ , „Hochspannungsleitung“),** Linearen Approximation:

$$y_i + r_{i1} = a_0 + a_1 u_i = \underbrace{\begin{bmatrix} 1 \\ u_i \end{bmatrix}}_{:= \mathbf{m}_i^T} \underbrace{\begin{bmatrix} a_0 \\ a_1 \end{bmatrix}}_{:= \boldsymbol{\beta}_1} = \mathbf{m}_i^T \boldsymbol{\beta}_1 \quad \text{mit } V(y_i) = \sigma_1^2 \text{ und } C(y_i, y_j) = 0$$

$$\text{für } i, j \in \{1, \dots, N\} \text{ } j \neq i \quad (421.3)$$

**Modell 2 (Objekt  $\mathcal{O}_2$ , „Straße“),** Approximation durch ein Polynom 2ten Grades:

$$y_i + r_{i2} = b_0 + b_1 u_i + b_2 u_i^2 = \underbrace{\begin{bmatrix} 1 \\ u_i \\ u_i^2 \end{bmatrix}}_{:= \mathbf{m}_i^T} \underbrace{\begin{bmatrix} b_0 \\ b_1 \\ b_2 \end{bmatrix}}_{:= \boldsymbol{\beta}_2} = \mathbf{m}_i^T \boldsymbol{\beta}_2 \quad \text{mit } V(y_i) = \sigma_2^2 \text{ und } C(y_i, y_j) = 0$$

$$\text{für } i, j \in \{1, \dots, N\} \text{ } j \neq i \quad (421.4)$$

Hierin stellen die Elemente von  $\boldsymbol{\beta}_1$  und  $\boldsymbol{\beta}_2$  sowie  $\sigma_1$  und  $\sigma_2$  die unbekannt Parameter des Schätzproblems dar. Sie werden im Vektor

$$\boldsymbol{\beta} = [\boldsymbol{\beta}_1^T, \sigma_1, \boldsymbol{\beta}_2^T, \sigma_2]^T \quad (421.5)$$

zusammengefasst. Mit  $r_{i1}$  und  $r_{i2}$  sind die Beobachtungsresiduen bezeichnet, d.h. die Abweichungen der Beobachtungen  $y_i$  von ihren Erwartungswerten  $\hat{y}_{i1}$  bzw.  $\hat{y}_{i2}$  in den Modellen 1 und 2.

Wie in Abschnitt 41 angedeutet, fehlen im Beobachtungsmaterial Informationen über die Zugehörigkeit der einzelnen Beobachtungen  $y_i$  zu Modell 1 oder Modell 2, so daß diese Informationen im folgenden als unzugängliche Beobachtungen aufgefaßt werden. Formal wird für jede tatsächliche Beobachtung  $y_i$  ein Vektor

$$\mathbf{x}_i = \begin{bmatrix} x_{i1} \\ x_{i2} \end{bmatrix} \in \{\mathbf{e}_1, \mathbf{e}_2\} \quad \text{mit } \mathbf{e}_1 = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \mathbf{e}_2 = \begin{bmatrix} 0 \\ 1 \end{bmatrix} \quad (421.6)$$

nicht zugänglicher Beobachtungen  $x_{i1}$  und  $x_{i2}$  eingeführt. In Abhängigkeit davon, ob die Beobachtung  $y_i$  zur Hochspannungsleitung (Modell 1) oder zur Straße (Modell 2) gehört, gilt  $\mathbf{x}_i = \mathbf{e}_1$  oder  $\mathbf{x}_i = \mathbf{e}_2$ . Gehört die Beobachtung  $y_i$  zu Modell 1, so gilt  $\mathbf{x}_i = \mathbf{e}_1$ ; im

<sup>6</sup>U.U. sind die Mittellinien der beiden Objekte im Bild nicht gleich gut lokalisierbar. (Beispielsweise ist die Mittellinie einer Hochspannungsleitung durch das mittlere Kabel schärfer festgelegt als die Mittellinie einer Straße durch den doch (relativ breiten) Mittelstreifen.) Daher ist zwischen der Varianz der zu Objekt 1 und Objekt 2 gehörenden Beobachtungen zu unterscheiden.

anderen Fall, in dem die Beobachtung  $y_i$  zu Modell 2 gehört, gilt  $\mathbf{x}_i = \mathbf{e}_2$ . Daher wird  $\mathbf{x}_i$  auch als *Indikatorvektor* mit den *Indikatoren*  $x_{i1}$  und  $x_{i2}$  bezeichnet.

Die Vektoren unzugänglicher Beobachtungen  $\mathbf{x}_i$ ,  $i \in \{1, \dots, N\}$  werden in dem Vektor  $\mathbf{x}$

$$\mathbf{x} = \begin{bmatrix} \mathbf{x}_1 \\ \dots \\ \mathbf{x}_N \end{bmatrix} = [x_{11}, x_{12}, \dots, x_{N1}, x_{N2}]^T \quad (421.7)$$

zusammengefasst. Dies ist der unzugängliche Beobachtungsvektor des zu lösenden Schätzproblems aus unvollständigen Beobachtungsdaten.

Der zur Beobachtung  $y_i$  gehörende vollständige Beobachtungsvektor ergibt sich unter Einbeziehung der unzugänglichen Beobachtungen  $\mathbf{x}_i$  zu

$$\mathbf{z}_i = \begin{bmatrix} y_i \\ \mathbf{x}_i \end{bmatrix} = [y_i, x_{i1}, x_{i2}]^T. \quad (421.8)$$

Für den vollständigen Beobachtungsvektor des Schätzproblems aus unvollständigen Daten gilt damit

$$\mathbf{z} = \begin{bmatrix} \mathbf{z}_1 \\ \dots \\ \mathbf{z}_N \end{bmatrix} = [y_1, x_{11}, x_{12}, \dots, y_N, x_{N1}, x_{N2}]^T. \quad (421.9)$$

Die bisherigen Ausführungen zusammenfassend, kann die im Rahmen der Parameterschätzung zu lösende Aufgabe wie folgt als Schätzproblem aus unvollständigen Beobachtungsdaten formuliert werden:

**Schätzproblem aus unvollständigen Beobachtungsdaten:**

Von den in (421.9) spezifizierten vollständigen Beobachtungsdaten  $\mathbf{z}$  sind nur die gemäß (421.7) in dem Vektor  $\mathbf{y}$  enthaltenen Beobachtungen tatsächlich gemessen, es fehlen die Beobachtungen  $x_{i1}$  und  $x_{i2}$  für  $i \in \{1, \dots, N\}$  aus (421.6).

Ausgehend hiervon sind die unbekannt Parameter  $\beta_1$  und  $\beta_2$  der Modelle 1 und 2 (vgl. (421.3) bzw. (421.4)) sowie die zugehörigen Varianzen  $\sigma_1^2$  und  $\sigma_2^2$  der Beobachtungen zu schätzen.

Dieses Schätzproblem läßt sich mittels des EM-Algorithmus lösen. Im nächsten Abschnitt wird aufgezeigt, welche Berechnungen hierfür im E-Schritt und im M-Schritt erforderlich sind.

## 422 Ableitung der EM-Iterationsschritte

Wie in Kapitel 2 ausführlich erläutert, wird im Rahmen des EM-Algorithmus die Kullback-Leibler-Statistik iterativ berechnet und maximiert (vgl. 324.20). Die Berechnung der Kullback-Leibler-Statistik setzt voraus, daß die von den unbekannt Parametern  $\beta$  abhängige Dichtefunktion  $f(\mathbf{z} | \beta)$  des vollständigen Beobachtungsvektors  $\mathbf{z}$  bekannt ist. Daher wird zunächst diese Dichtefunktion angegeben.

### Herleitung der Dichte $f(\mathbf{z} | \beta)$ des vollständigen Beobachtungsvektors

Zur Ableitung der Dichtefunktion  $f(\mathbf{z} | \beta)$  wird zunächst mit  $p(\mathbf{x}_i | \beta)$  Wahrscheinlichkeitsverteilung der diskreten Zufallsvariablen  $\mathbf{x}_i \in \{\mathbf{e}_1, \mathbf{e}_2\}$  bezeichnet, die a priori, d.h.

ohne Kenntnis von der numerischen Realisierung des Beobachtungsvektors  $\mathbf{y}$ , jeder Beobachtung  $y_i$  die Wahrscheinlichkeiten  $P(\mathbf{x}_i = \mathbf{e}_1 | \boldsymbol{\beta}) = p(\mathbf{e}_1 | \boldsymbol{\beta})$  bzw.  $P(\mathbf{x}_i = \mathbf{e}_2 | \boldsymbol{\beta}) = p(\mathbf{e}_2 | \boldsymbol{\beta})$  zuordnet, daß sie zu Modell 1 (d.h. zur Hochspannungsleitung) bzw. zu Modell 2 (d.h. zur Straße) gehört, sofern mit  $\boldsymbol{\beta}$  die Parameter der beiden Modelle gegeben sind. Es wird hier angenommen, daß diese Dichtefunktion für alle Beobachtungen  $y_i$  dieselbe ist.

Ginge aus der Aufgabenstellung ein Zusammenhang zwischen den Modellparametern  $\boldsymbol{\beta}$  und den Wahrscheinlichkeiten  $P(\mathbf{x}_i = \mathbf{e}_1)$  bzw.  $P(\mathbf{x}_i = \mathbf{e}_2)$  hervor, so würde sich hieraus für  $p(\mathbf{x}_i | \boldsymbol{\beta})$  auch tatsächlich eine von den Parametern  $\boldsymbol{\beta}$  abhängige Funktion ergeben. Da aus der hier vorliegenden Aufgabenstellung ein solcher Zusammenhang jedoch nicht ersichtlich ist, werden die Wahrscheinlichkeiten  $P(\mathbf{x}_i = \mathbf{e}_1)$  und  $P(\mathbf{x}_i = \mathbf{e}_2)$  nach dem erwarteten Anteil der zur Hochspannungsleitung (Modell 1) bzw. zur Straße (Modell 2) gehörenden Beobachtungen an der Gesamtzahl  $N$  der Beobachtungen bemessen. Es wird angenommen, daß 1/3 der Beobachtungen zur Hochspannungsleitung und 2/3 der Beobachtungen zur Straße gehören. Daher gilt

$$\begin{aligned} P(\text{Beobachtung } y_i \text{ gehört zu Objekt } \mathcal{O}_1 \text{ (HS-Leitung)}) \\ = P(\mathbf{x}_i = \mathbf{e}_1 | \boldsymbol{\beta}) = p(\mathbf{e}_1 | \boldsymbol{\beta}) = \frac{1}{3} \end{aligned} \quad (422.1)$$

und

$$\begin{aligned} P(\text{Beobachtung } y_i \text{ gehört zu Objekt } \mathcal{O}_2 \text{ (Straße)}) \\ = P(\mathbf{x}_i = \mathbf{e}_2 | \boldsymbol{\beta}) = p(\mathbf{e}_2 | \boldsymbol{\beta}) = \frac{2}{3}. \end{aligned} \quad (422.2)$$

Die Beobachtungen  $y_i$  werden unabhängig davon, zu welchem Objekt sie gehören (d.h. unabhängig davon, welches Beobachtungsmodell ihnen zugrunde liegt), als normalverteilt vorausgesetzt. Je nachdem, ob eine Beobachtung  $y_i$  zu Objekt  $\mathcal{O}_1$  oder zu Objekt  $\mathcal{O}_2$  gehört, liegen ihrer Normalverteilung allerdings unterschiedliche Parameter zugrunde (vgl. Gl. (421.4)). Gehört die Beobachtung  $y_i$  zu der Hochspannungsleitung (Modell 1), so gilt für die Wahrscheinlichkeitsdichte  $f(y_i | \mathbf{x}_i = \mathbf{e}_1, \boldsymbol{\beta})$  in Abhängigkeit von den Parametern  $\boldsymbol{\beta}_1$  der linearen Approximation und der Standardabweichung  $\sigma_1$  der zur Hochspannungsleitung gehörenden Beobachtungen

**Modell 1** (lineare Approximation):

$$f(y_i | \mathbf{x}_i = \mathbf{e}_1, \boldsymbol{\beta}) = f(y_i | \mathbf{e}_1, \boldsymbol{\beta}_1, \sigma_1) = \frac{1}{\sqrt{2\pi}\sigma_1} \exp\left(-\frac{1}{2\sigma_1^2}[y_i - \mathbf{m}_i^T \boldsymbol{\beta}_1]^2\right) \quad (422.3)$$

Gehört die Beobachtung  $y_i$  hingegen zur Straße (Modell 2), so gilt für  $f(y_i | \mathbf{x}_i = \mathbf{e}_2, \boldsymbol{\beta})$  in Abhängigkeit von den Parametern  $\boldsymbol{\beta}_2$  der Approximation mittels eines Polynoms 2. Grades und der Standardabweichung  $\sigma_2$  der Beobachtungen, die zur Straße gehören:

**Modell 2** (Approximation durch Polynom 2ten Grades):

$$f(y_i | \mathbf{x}_i = \mathbf{e}_2, \boldsymbol{\beta}) = f(y_i | \mathbf{e}_2, \boldsymbol{\beta}_2, \sigma_2) = \frac{1}{\sqrt{2\pi}\sigma_2} \exp\left(-\frac{1}{2\sigma_2^2}[y_i - \mathbf{m}_i^T \boldsymbol{\beta}_2]^2\right) \quad (422.4)$$

Mit (422.1), (422.2), (422.3) und (422.4) läßt sich die Dichtefunktion  $f(\mathbf{z}_i | \boldsymbol{\beta}) = f(\mathbf{z}_i | \boldsymbol{\beta}_1, \sigma_1, \boldsymbol{\beta}_2, \sigma_2)$  der vollständigen Beobachtung  $\mathbf{z}_i$  ableiten. Aufgrund der Bayes - Formel



gilt

$$f(y_i, \mathbf{x}_i = \mathbf{e}_1 \mid \boldsymbol{\beta}) = p(\mathbf{e}_1 \mid \boldsymbol{\beta}) \cdot f(y_i \mid \mathbf{e}_1, \boldsymbol{\beta}) = \frac{1}{3\sqrt{2\pi}\sigma_1} e^{-\frac{1}{2\sigma_1^2}[y_i - \mathbf{m}_i^T \boldsymbol{\beta}_1]^2} \quad (422.5)$$

und

$$f(y_i, \mathbf{x}_i = \mathbf{e}_2 \mid \boldsymbol{\beta}) = p(\mathbf{e}_2 \mid \boldsymbol{\beta}) \cdot f(y_i \mid \mathbf{e}_2, \boldsymbol{\beta}) = \frac{2}{3\sqrt{2\pi}\sigma_2} e^{-\frac{1}{2\sigma_2^2}[y_i - \mathbf{m}_i^T \boldsymbol{\beta}_2]^2}. \quad (422.6)$$

Hiermit erhält man unter Ausnutzung der Tatsache, daß  $\mathbf{x}_i$  in Abhängigkeit von dem für die Beobachtung  $y_i$  zutreffenden Modell entweder die Elemente  $x_{i1} = 1$  und  $x_{i2} = 0$  oder  $x_{i1} = 0$  und  $x_{i2} = 1$  enthält, die Dichte

$$\begin{aligned} f(\mathbf{z}_i \mid \boldsymbol{\beta}) &= f(y_i, \mathbf{x}_i \mid \boldsymbol{\beta}) \\ &= x_{i1} \cdot p(\mathbf{e}_1 \mid \boldsymbol{\beta}) \cdot f(y_i \mid \mathbf{e}_1, \boldsymbol{\beta}) + x_{i2} \cdot p(\mathbf{e}_2 \mid \boldsymbol{\beta}) \cdot f(y_i \mid \mathbf{e}_2, \boldsymbol{\beta}) \\ &= [p(\mathbf{e}_1 \mid \boldsymbol{\beta}) \cdot f(y_i \mid \mathbf{e}_1, \boldsymbol{\beta})]^{x_{i1}} \cdot [p(\mathbf{e}_2 \mid \boldsymbol{\beta}) \cdot f(y_i \mid \mathbf{e}_2, \boldsymbol{\beta})]^{x_{i2}}, \end{aligned} \quad (422.7)$$

so daß folgt

$$f(\mathbf{z}_i \mid \boldsymbol{\beta}) = \left( \frac{1}{3\sqrt{2\pi}\sigma_1} e^{-\frac{1}{2\sigma_1^2}[y_i - \mathbf{m}_i^T \boldsymbol{\beta}_1]^2} \right)^{x_{i1}} \cdot \left( \frac{2}{3\sqrt{2\pi}\sigma_2} e^{-\frac{1}{2\sigma_2^2}[y_i - \mathbf{m}_i^T \boldsymbol{\beta}_2]^2} \right)^{x_{i2}}. \quad (422.8)$$

Die tatsächlichen Beobachtungen  $y_i$  und die nicht zugänglichen Beobachtungsvektoren  $\mathbf{x}_i$  werden als paarweise statistisch voneinander unabhängig aufgefaßt. Daher erhält man nach (vgl. [KOCH 1998], Gl. (215.1)) die Dichtefunktion des vollständigen Beobachtungsvektors  $\mathbf{z}$  aus dem Produkt der Einzeldichten  $f(\mathbf{z}_i \mid \boldsymbol{\beta})$ , es ergibt sich also die **Dichtefunktion des vollständigen Beobachtungsvektors**

$$\begin{aligned} f(\mathbf{z} \mid \boldsymbol{\beta}) &= f(\mathbf{z} \mid \boldsymbol{\beta}_1, \sigma_1, \boldsymbol{\beta}_2, \sigma_2) = \prod_{i=1}^N f(\mathbf{z}_i \mid \boldsymbol{\beta}) \\ &= \prod_{i=1}^N \left( \frac{1}{3\sqrt{2\pi}\sigma_1} e^{-\frac{1}{2\sigma_1^2}[y_i - \mathbf{m}_i^T \boldsymbol{\beta}_1]^2} \right)^{x_{i1}} \cdot \left( \frac{2}{3\sqrt{2\pi}\sigma_2} e^{-\frac{1}{2\sigma_2^2}[y_i - \mathbf{m}_i^T \boldsymbol{\beta}_2]^2} \right)^{x_{i2}}. \end{aligned} \quad (422.9)$$

Durch Logarithmieren dieser Dichtefunktion erhält man den **Logarithmus der Dichte des vollständigen Beobachtungsvektors**

$$\begin{aligned} \log f(\mathbf{z} \mid \boldsymbol{\beta}) &= \sum_{i=1}^N \log f(\mathbf{z}_i \mid \boldsymbol{\beta}_1, \sigma_1, \boldsymbol{\beta}_2, \sigma_2) \\ &= \sum_{i=1}^N \left\{ x_{i1} \log \left( \frac{1}{3\sqrt{2\pi}\sigma_1} e^{-\frac{1}{2\sigma_1^2}[y_i - \mathbf{m}_i^T \boldsymbol{\beta}_1]^2} \right) + x_{i2} \log \left( \frac{2}{3\sqrt{2\pi}\sigma_2} e^{-\frac{1}{2\sigma_2^2}[y_i - \mathbf{m}_i^T \boldsymbol{\beta}_2]^2} \right) \right\}. \end{aligned} \quad (422.10)$$

Auf Grundlage der Gleichungen (422.9) und (422.10) soll nun die beim EM-Algorithmus in jedem Iterationsschritt zu maximierende Kullback-Leibler-Statistik abgeleitet werden.

### E-Schritt: Berechnung der Kullback-Leibler-Statistik

Sind mit  $\hat{\boldsymbol{\beta}}^{(j)} = [\hat{\boldsymbol{\beta}}_1^{(j)T}, \hat{\sigma}_1^{(j)}, \hat{\boldsymbol{\beta}}_2^{(j)T}, \hat{\sigma}_2^{(j)}]^T$  vorläufige Schätzwerte für die unbekanten Parameter bekannt, so ergibt sich mit  $\boldsymbol{\beta} = [\boldsymbol{\beta}_1^T, \sigma_1, \boldsymbol{\beta}_2^T, \sigma_2]^T$  die Kullback-Leibler-Statistik

$Q(\boldsymbol{\beta} \mid \hat{\boldsymbol{\beta}}^{(j)}) = E[\log f(\mathbf{z} \mid \boldsymbol{\beta}) \mid \mathbf{y}, \hat{\boldsymbol{\beta}}^{(j)}]$  als Erwartungswert der logarithmierten Likelihoodfunktion (422.10). Im Falle einer überabzählbaren Menge  $\mathcal{Z}(\mathbf{y})$  wäre  $Q(\boldsymbol{\beta} \mid \hat{\boldsymbol{\beta}}^{(j)})$  nach (231.8) durch Integration zu berechnen. In unserem Beispiel handelt es sich bei  $\mathcal{Z}(\mathbf{y})$  wegen

$$\begin{aligned} \mathcal{Z}(\mathbf{y}) &= \mathcal{Z}([y_1, y_2, \dots, y_N]^T) \\ &= \left\{ [y_1, x_{11}, x_{12}, \dots, y_N, x_{N1}, x_{N2}]^T \mid \begin{pmatrix} x_{i1} \\ x_{i2} \end{pmatrix} \in \left\{ \begin{pmatrix} 1 \\ 0 \end{pmatrix}; \begin{pmatrix} 0 \\ 1 \end{pmatrix} \right\} \forall i \in \{1, \dots, N\} \right\} \end{aligned} \quad (422.11)$$

allerdings um eine diskrete Menge, so daß anstelle der Integration (231.8) eine Summation über alle Elemente  $\mathbf{z} \in \mathcal{Z}(\mathbf{y})$  durchzuführen ist. Mit (422.10) ergibt sich

$$\begin{aligned} Q(\boldsymbol{\beta} \mid \hat{\boldsymbol{\beta}}^{(j)}) &= E[\log f(\mathbf{z} \mid \boldsymbol{\beta}) \mid \mathbf{y}, \hat{\boldsymbol{\beta}}^{(j)}] = \sum_{\mathbf{z} \in \mathcal{Z}(\mathbf{y})} \log f(\mathbf{z} \mid \boldsymbol{\beta}) \cdot f(\mathbf{z} \mid \mathbf{y}, \hat{\boldsymbol{\beta}}^{(j)}) \\ &= \sum_{x_{N1}=0}^1 \cdots \sum_{x_{21}=0}^1 \sum_{x_{11}=0}^1 \sum_{i=1}^N \left\{ x_{i1} \log \left( \frac{1}{3\sqrt{2\pi}\sigma_1} e^{-\frac{1}{2\sigma_1^2}[y_i - \mathbf{m}_i^T \boldsymbol{\beta}_1]^2} \right) \right. \\ &\quad \left. + (1 - x_{i1}) \log \left( \frac{2}{3\sqrt{2\pi}\sigma_2} e^{-\frac{1}{2\sigma_2^2}[y_i - \mathbf{m}_i^T \boldsymbol{\beta}_2]^2} \right) \right\} \cdot f(\mathbf{z} \mid \mathbf{y}, \hat{\boldsymbol{\beta}}^{(j)}). \end{aligned} \quad (422.12)$$

Man erkennt in dieser Gleichung, daß der Erwartungswert  $Q(\boldsymbol{\beta} \mid \hat{\boldsymbol{\beta}}^{(j)})$  linear in den Indikatoren  $x_{i1}$  und  $x_{i2} = 1 - x_{i1}$  ist. Wegen der Linearität des Erwartungswertoperators (vgl. [KOCH 1998], Gl. 231.5) gilt daher

$$\begin{aligned} Q(\boldsymbol{\beta} \mid \hat{\boldsymbol{\beta}}^{(j)}) &= E[\log f(\mathbf{z} \mid \boldsymbol{\beta}) \mid \mathbf{y}, \hat{\boldsymbol{\beta}}^{(j)}] \\ &= \sum_{i=1}^N \left\{ E[x_{i1} \mid \mathbf{y}, \hat{\boldsymbol{\beta}}^{(j)}] \log \left( \frac{1}{3\sqrt{2\pi}\sigma_1} e^{-\frac{1}{2\sigma_1^2}[y_i - \mathbf{m}_i^T \boldsymbol{\beta}_1]^2} \right) \right. \\ &\quad \left. + E[x_{i2} \mid \mathbf{y}, \hat{\boldsymbol{\beta}}^{(j)}] \log \left( \frac{2}{3\sqrt{2\pi}\sigma_2} e^{-\frac{1}{2\sigma_2^2}[y_i - \mathbf{m}_i^T \boldsymbol{\beta}_2]^2} \right) \right\}. \end{aligned} \quad (422.13)$$

Der Expectation - Schritt des hier vorliegenden EM-Algorithmus (d.h. die Berechnung der Kullback-Leibler-Statistik) besteht also im wesentlichen in der Berechnung der Erwartungswerte  $E[x_{i1} \mid \mathbf{y}, \hat{\boldsymbol{\beta}}^{(j)}]$  und  $E[x_{i2} \mid \mathbf{y}, \hat{\boldsymbol{\beta}}^{(j)}]$ . Daher werden nun Gleichungen zur Berechnung dieser Erwartungswerte hergeleitet.

Die Indikatoren  $x_{i1}$  und  $x_{i2}$  für  $i \in \{1, \dots, N\}$  sind diskrete Zufallsvariablen, die nur die Werte 0 und 1 annehmen können. Der Erwartungswert einer solchen Zufallsvariable  $X$  ist nach ([KOCH 1998], Gl. (231.1)) definiert als die Summe

$$E(X) = \sum_{x=0}^1 x \cdot P(X = x) = 1 \cdot p(X = 1) + 0 \cdot p(X = 0) = P(X = 1). \quad (422.14)$$

Daher gilt für die Erwartungswerte von  $x_{i1}$  und  $x_{i2}$  mit  $i \in \{1, \dots, n\}$

$$E[x_{i1} \mid \mathbf{y}, \hat{\boldsymbol{\beta}}^{(j)}] = P(x_{i1} = 1 \mid \mathbf{y}, \hat{\boldsymbol{\beta}}^{(j)}) \quad \text{und} \quad (422.15)$$

$$E[x_{i2} \mid \mathbf{y}, \hat{\boldsymbol{\beta}}^{(j)}] = P(x_{i2} = 1 \mid \mathbf{y}, \hat{\boldsymbol{\beta}}^{(j)}), \quad (422.16)$$

d.h. der Erwartungswert  $E[x_{i1} | \mathbf{y}, \hat{\boldsymbol{\beta}}^{(j)}]$  gibt die Wahrscheinlichkeit dafür an, daß der Beobachtung  $y_i$  das Modell 1 (Hochspannungsleitung) zugrundeliegt, sofern insgesamt die Beobachtungen  $\mathbf{y}$  vorliegen und die vorläufigen Parameter  $\hat{\boldsymbol{\beta}}^{(j)}$  der beiden Modelle 1 und 2 gegeben sind. Entsprechend gibt der Erwartungswert  $E[x_{i2} | \mathbf{y}, \hat{\boldsymbol{\beta}}^{(j)}]$  die Wahrscheinlichkeit dafür an, daß der Beobachtung  $y_i$  das Modell 2 (Straße) zugrundeliegt, falls  $\mathbf{y}$  und  $\hat{\boldsymbol{\beta}}^{(j)}$  gegeben sind. Da die Zugehörigkeit der Beobachtung  $y_i$  zur Hochspannungsleitung (Modell 1) bzw. zur Straße (Modell 2) zwei sich gegenseitig ausschließende Ereignisse darstellen, ergibt sich nach ([KOCH 1998], Gl. (213.3))

$$E[x_{i1} | \mathbf{y}, \hat{\boldsymbol{\beta}}^{(j)}] + E[x_{i2} | \mathbf{y}, \hat{\boldsymbol{\beta}}^{(j)}] = 1. \quad (422.17)$$

Nach der Bayes-Formel gilt zunächst

$$P(x_{i1} = 1 | \mathbf{y}, \hat{\boldsymbol{\beta}}^{(j)}) = P(\mathbf{x}_i = \mathbf{e}_1 | \mathbf{y}, \hat{\boldsymbol{\beta}}^{(j)}) = \frac{f(\mathbf{y}, \mathbf{x}_i = \mathbf{e}_1 | \hat{\boldsymbol{\beta}}^{(j)})}{f(\mathbf{y} | \hat{\boldsymbol{\beta}}^{(j)})}. \quad (422.18)$$

Hierin gilt mit (422.1) wiederum aufgrund der Bayes-Formel

$$f(\mathbf{y}, \mathbf{x}_i = \mathbf{e}_1 | \hat{\boldsymbol{\beta}}^{(j)}) = p(\mathbf{e}_1 | \hat{\boldsymbol{\beta}}^{(j)}) \cdot f(\mathbf{y} | \mathbf{x}_i = \mathbf{e}_1, \hat{\boldsymbol{\beta}}^{(j)}), \quad (422.19)$$

woraus wegen der statistischen Unabhängigkeit der Beobachtungen  $y_i$  folgt

$$\begin{aligned} f(\mathbf{y}, \mathbf{x}_i = \mathbf{e}_1 | \hat{\boldsymbol{\beta}}^{(j)}) &= p(\mathbf{e}_1 | \hat{\boldsymbol{\beta}}^{(j)}) \cdot \prod_{k=1}^N f(y_k | \mathbf{x}_i = \mathbf{e}_1, \hat{\boldsymbol{\beta}}^{(j)}) \\ &= p(\mathbf{e}_1 | \hat{\boldsymbol{\beta}}^{(j)}) \cdot f(y_i | \mathbf{x}_i = \mathbf{e}_1, \hat{\boldsymbol{\beta}}^{(j)}) \cdot \prod_{\substack{k=1 \\ k \neq i}}^N f(y_k | \mathbf{x}_i = \mathbf{e}_1, \hat{\boldsymbol{\beta}}^{(j)}). \end{aligned} \quad (422.20)$$

Nach (422.3) und (422.4) sind die Dichtefunktionen  $f(y_k | \mathbf{x}_k = \mathbf{e}_1, \hat{\boldsymbol{\beta}}^{(j)})$  und  $f(y_k | \mathbf{x}_k = \mathbf{e}_2, \hat{\boldsymbol{\beta}}^{(j)})$  für  $k \neq i$  unabhängig vom Indikator  $\mathbf{x}_i$ . Daher gilt für alle  $k \neq i$  die Beziehung  $f(y_k | \mathbf{x}_i = \mathbf{e}_1, \hat{\boldsymbol{\beta}}^{(j)}) = f(y_k | \hat{\boldsymbol{\beta}}^{(j)})$  und es folgt

$$\begin{aligned} f(\mathbf{y}, \mathbf{x}_i = \mathbf{e}_1 | \hat{\boldsymbol{\beta}}^{(j)}) &= p(\mathbf{e}_1 | \hat{\boldsymbol{\beta}}^{(j)}) \cdot f(y_i | \mathbf{x}_i = \mathbf{e}_1, \hat{\boldsymbol{\beta}}^{(j)}) \cdot \prod_{\substack{k=1 \\ k \neq i}}^N f(y_k | \hat{\boldsymbol{\beta}}^{(j)}) \\ &= f(y_i, \mathbf{x}_i = \mathbf{e}_1 | \hat{\boldsymbol{\beta}}^{(j)}) \cdot \prod_{\substack{k=1 \\ k \neq i}}^N f(y_k | \hat{\boldsymbol{\beta}}^{(j)}). \end{aligned} \quad (422.21)$$

Wegen der statistischen Unabhängigkeit der Beobachtungen  $y_i$  gilt außerdem

$$f(\mathbf{y} | \hat{\boldsymbol{\beta}}^{(j)}) = \prod_{k=1}^N f(y_k | \hat{\boldsymbol{\beta}}^{(j)}). \quad (422.22)$$

Durch Einsetzen von (422.21) und (422.22) in (422.18) erhält man schließlich

$$\begin{aligned} P(x_{i1} = 1 | \mathbf{y}, \hat{\boldsymbol{\beta}}^{(j)}) &= \frac{f(y_i, \mathbf{x}_i = \mathbf{e}_1 | \hat{\boldsymbol{\beta}}^{(j)}) \cdot \prod_{\substack{k=1 \\ k \neq i}}^N f(y_k | \hat{\boldsymbol{\beta}}^{(j)})}{\prod_{k=1}^N f(y_k | \hat{\boldsymbol{\beta}}^{(j)})} \\ &= \frac{f(y_i, \mathbf{x}_i = \mathbf{e}_1 | \hat{\boldsymbol{\beta}}^{(j)})}{f(y_i | \hat{\boldsymbol{\beta}}^{(j)})}. \end{aligned} \quad (422.23)$$

Bei der Zugehörigkeit einer Beobachtung  $y_i$  zur Hochspannungsleitung bzw. zur Straße handelt es sich um zwei sich gegenseitig ausschließende Ereignisse. Daher gilt (vgl. [KOCH 1998], Gl.(213.3))

$$f(y_i | \hat{\boldsymbol{\beta}}^{(j)}) = f(y_i, \mathbf{x}_i = \mathbf{e}_1 | \hat{\boldsymbol{\beta}}^{(j)}) + f(y_i, \mathbf{x}_i = \mathbf{e}_2 | \hat{\boldsymbol{\beta}}^{(j)}), \quad (422.24)$$

Hiermit ergibt sich

$$P(x_{i1} = 1 | \mathbf{y}, \hat{\boldsymbol{\beta}}^{(j)}) = \frac{f(y_i, \mathbf{x}_i = \mathbf{e}_1 | \hat{\boldsymbol{\beta}}^{(j)})}{f(y_i, \mathbf{x}_i = \mathbf{e}_1 | \hat{\boldsymbol{\beta}}^{(j)}) + f(y_i, \mathbf{x}_i = \mathbf{e}_2 | \hat{\boldsymbol{\beta}}^{(j)})} = E[x_{i1} | \mathbf{y}, \hat{\boldsymbol{\beta}}^{(j)}], \quad (422.25)$$

worin die Dichten  $f(y_i, \mathbf{x}_i = \mathbf{e}_1 | \hat{\boldsymbol{\beta}}^{(j)})$  und  $f(y_i, \mathbf{x}_i = \mathbf{e}_2 | \hat{\boldsymbol{\beta}}^{(j)})$  nach (422.3) und (422.4) berechnet werden. Ganz analog läßt sich zeigen, daß gilt

$$P(x_{i2} = 1 | \mathbf{y}, \hat{\boldsymbol{\beta}}^{(j)}) = \frac{f(y_i, \mathbf{x}_i = \mathbf{e}_2 | \hat{\boldsymbol{\beta}}^{(j)})}{f(y_i, \mathbf{x}_i = \mathbf{e}_1 | \hat{\boldsymbol{\beta}}^{(j)}) + f(y_i, \mathbf{x}_i = \mathbf{e}_2 | \hat{\boldsymbol{\beta}}^{(j)})} = E[x_{i2} | \mathbf{y}, \hat{\boldsymbol{\beta}}^{(j)}]. \quad (422.26)$$

Mit (422.5) und (422.6) ergeben sich schließlich die **Erwartungswerte der fehlenden Beobachtungen**:

$$\mu_{i1}^{(j)} := E[x_{i1} | \mathbf{y}, \hat{\boldsymbol{\beta}}^{(j)}] = \frac{\frac{1}{3\sqrt{2\pi}\hat{\sigma}_1^{(j)}} e^{-\frac{1}{2(\hat{\sigma}_1^{(j)})^2} [y_j - \mathbf{m}_j^T \hat{\boldsymbol{\beta}}_1^{(j)}]^2}}{\frac{1}{3\sqrt{2\pi}\hat{\sigma}_1^{(j)}} e^{-\frac{1}{2(\hat{\sigma}_1^{(j)})^2} [y_j - \mathbf{m}_j^T \hat{\boldsymbol{\beta}}_1^{(j)}]^2} + \frac{2}{3\sqrt{2\pi}\hat{\sigma}_2^{(j)}} e^{-\frac{1}{2(\hat{\sigma}_2^{(j)})^2} [y_j - \mathbf{m}_j^T \hat{\boldsymbol{\beta}}_2^{(j)}]^2}} \quad (422.27)$$

und

$$\mu_{i2}^{(j)} := E[x_{i2} | \mathbf{y}, \hat{\boldsymbol{\beta}}^{(j)}] = \frac{\frac{2}{3\sqrt{2\pi}\hat{\sigma}_2^{(j)}} e^{-\frac{1}{2(\hat{\sigma}_2^{(j)})^2} [y_j - \mathbf{m}_j^T \hat{\boldsymbol{\beta}}_2^{(j)}]^2}}{\frac{1}{3\sqrt{2\pi}\hat{\sigma}_1^{(j)}} e^{-\frac{1}{2(\hat{\sigma}_1^{(j)})^2} [y_j - \mathbf{m}_j^T \hat{\boldsymbol{\beta}}_1^{(j)}]^2} + \frac{2}{3\sqrt{2\pi}\hat{\sigma}_2^{(j)}} e^{-\frac{1}{2(\hat{\sigma}_2^{(j)})^2} [y_j - \mathbf{m}_j^T \hat{\boldsymbol{\beta}}_2^{(j)}]^2}} \quad (422.28)$$

Der E-Schritt des zur Lösung der Aufgabe anzuwendenden EM-Algorithmus besteht in der Berechnung der Erwartungswerte  $\mu_{i1}^{(j)}$  und  $\mu_{i2}^{(j)}$  für  $i \in \{1, \dots, N\}$  nach diesen beiden Gleichungen. Im  $j$ -ten Iterationsschritt ergibt sich dann nach (422.13) die **Kullback-Leibler-Statistik** in Abhängigkeit von dem Parametervektor  $\boldsymbol{\beta}$  zu

$$Q(\boldsymbol{\beta} | \hat{\boldsymbol{\beta}}^{(j)}) = \sum_{i=1}^N \left\{ \mu_{i1}^{(j)} \log \left( \frac{1}{3\sqrt{2\pi}\sigma_1} e^{-\frac{1}{2\sigma_1^2} [y_i - \mathbf{m}_i^T \boldsymbol{\beta}_1]^2} \right) + \mu_{i2}^{(j)} \log \left( \frac{2}{3\sqrt{2\pi}\sigma_2} e^{-\frac{1}{2\sigma_2^2} [y_i - \mathbf{m}_i^T \boldsymbol{\beta}_2]^2} \right) \right\}. \quad (422.29)$$

Kontrolle der Gleichungen (422.27) und (422.28):

Zum gleichen Ergebnis wie in (422.27) und (422.28) kommt man *ohne* die Betrachtung der Erwartungswerte als Wahrscheinlichkeiten, wenn man ausgehend von (422.13) beispielsweise den Erwartungswert  $E[x_{i1} | \mathbf{y}, \hat{\boldsymbol{\beta}}^{(j)}]$  nach

$$E[x_{i1} | \mathbf{y}, \hat{\boldsymbol{\beta}}^{(j)}] = \sum_{\mathbf{z} \in \mathcal{Z}(\mathbf{y})} x_{i1} \cdot f(\mathbf{z} | \mathbf{y}, \hat{\boldsymbol{\beta}}^{(j)}),$$

berechnet. Nach (231.1) gilt  $f(\mathbf{z} | \mathbf{y}, \hat{\boldsymbol{\beta}}^{(j)}) = f(\mathbf{z} | \hat{\boldsymbol{\beta}}^{(j)}) [g(\mathbf{y} | \hat{\boldsymbol{\beta}}^{(j)})]^{-1}$ . Hierin ist die Dichte  $f(\mathbf{z} | \hat{\boldsymbol{\beta}}^{(j)})$  mit (422.9) gegeben, die Dichte  $g(\mathbf{y} | \hat{\boldsymbol{\beta}}^{(j)})$  ergibt sich aus (222.1), wobei an die Stelle der Integration wieder eine Summation über alle  $\mathbf{x} \in \mathcal{Z}(\mathbf{y})$  tritt:

$$\begin{aligned} g(\mathbf{y} | \hat{\boldsymbol{\beta}}^{(j)}) &= \sum_{\mathbf{z} \in \mathcal{Z}(\mathbf{y})} f(\mathbf{z} | \hat{\boldsymbol{\beta}}^{(j)}) \quad \text{mit } f(\mathbf{z} | \hat{\boldsymbol{\beta}}^{(j)}) \stackrel{(422.9)}{=} \prod_{i=1}^N f(z_i | \hat{\boldsymbol{\beta}}^{(j)}) \\ &= \sum_{\mathbf{z} \in \mathcal{Z}(\mathbf{y})} \prod_{i=1}^N f(z_i | \hat{\boldsymbol{\beta}}_1^{(j)}, \hat{\sigma}_1^{(j)}, \hat{\boldsymbol{\beta}}_2^{(j)}, \hat{\sigma}_2^{(j)}) \\ &= \sum_{\mathbf{z} \in \mathcal{Z}(\mathbf{y})} \prod_{i=1}^N \underbrace{\left( \frac{1}{3\sqrt{2\pi}\hat{\sigma}_1^{(j)}} e^{-\frac{1}{2(\hat{\sigma}_1^{(j)})^2} [y_i - \mathbf{m}_i^T \hat{\boldsymbol{\beta}}_1^{(j)}]^2} \right)^{x_{i1}}}_{=: W(y_i, \hat{\boldsymbol{\beta}}_1^{(j)}, \hat{\sigma}_1^{(j)})^{x_{i1}}} \cdot \underbrace{\left( \frac{2}{3\sqrt{2\pi}\hat{\sigma}_2^{(j)}} e^{-\frac{1}{2(\hat{\sigma}_2^{(j)})^2} [y_i - \mathbf{m}_i^T \hat{\boldsymbol{\beta}}_2^{(j)}]^2} \right)^{x_{i2}}}_{=: W(y_i, \hat{\boldsymbol{\beta}}_2^{(j)}, \hat{\sigma}_2^{(j)})^{x_{i2}}} \\ &= \sum_{x_{N1}=0}^1 \cdots \sum_{x_{21}=0}^1 \sum_{x_{11}=0}^1 \prod_{i=1}^N W(y_i, \hat{\boldsymbol{\beta}}_1^{(j)}, \hat{\sigma}_1^{(j)})^{x_{i1}} \cdot W(y_i, \hat{\boldsymbol{\beta}}_2^{(j)}, \hat{\sigma}_2^{(j)})^{1-x_{i1}} \\ &= \sum_{x_{N1}=0}^1 \cdots \sum_{x_{21}=0}^1 \prod_{i=2}^N W(y_i, \hat{\boldsymbol{\beta}}_1^{(j)}, \hat{\sigma}_1^{(j)})^{x_{i1}} \cdot W(y_i, \hat{\boldsymbol{\beta}}_2^{(j)}, \hat{\sigma}_2^{(j)})^{1-x_{i1}} \\ &\quad \cdot \left( W(y_1, \hat{\boldsymbol{\beta}}_1^{(j)}, \hat{\sigma}_1^{(j)}) + W(y_1, \hat{\boldsymbol{\beta}}_2^{(j)}, \hat{\sigma}_2^{(j)}) \right) \\ &= \sum_{x_{N1}=0}^1 \cdots \sum_{x_{31}=0}^1 \prod_{i=3}^N W(y_i, \hat{\boldsymbol{\beta}}_1^{(j)}, \hat{\sigma}_1^{(j)})^{x_{i1}} \cdot W(y_i, \hat{\boldsymbol{\beta}}_2^{(j)}, \hat{\sigma}_2^{(j)})^{1-x_{i1}} \\ &\quad \cdot \left( W(y_1, \hat{\boldsymbol{\beta}}_1^{(j)}, \hat{\sigma}_1^{(j)}) + W(y_2, \hat{\boldsymbol{\beta}}_2^{(j)}, \hat{\sigma}_2^{(j)}) \right) \cdot \left( W(y_2, \hat{\boldsymbol{\beta}}_1^{(j)}, \hat{\sigma}_1^{(j)}) + W(y_2, \hat{\boldsymbol{\beta}}_2^{(j)}, \hat{\sigma}_2^{(j)}) \right) \\ &= \dots \\ g(\mathbf{y} | \hat{\boldsymbol{\beta}}^{(j)}) &= \prod_{j=1}^N \left\{ \frac{1}{3\sqrt{2\pi}\hat{\sigma}_1^{(j)}} e^{-\frac{1}{2(\hat{\sigma}_1^{(j)})^2} [y_j - \mathbf{m}_j^T \hat{\boldsymbol{\beta}}_1^{(j)}]^2} + \frac{2}{3\sqrt{2\pi}\hat{\sigma}_2^{(j)}} e^{-\frac{1}{2(\hat{\sigma}_2^{(j)})^2} [y_j - \mathbf{m}_j^T \hat{\boldsymbol{\beta}}_2^{(j)}]^2} \right\} \quad (422.30) \end{aligned}$$

Hiermit erhält man schließlich

$$f(\mathbf{z} | \mathbf{y}, \hat{\boldsymbol{\beta}}^{(j)}) = \prod_{j=1}^N \frac{\left( \frac{1}{3\sqrt{2\pi}\hat{\sigma}_1^{(j)}} e^{-\frac{1}{2(\hat{\sigma}_1^{(j)})^2} [y_j - \mathbf{m}_j^T \hat{\boldsymbol{\beta}}_1^{(j)}]^2} \right)^{z_{j1}} \cdot \left( \frac{2}{3\sqrt{2\pi}\hat{\sigma}_2^{(j)}} e^{-\frac{1}{2(\hat{\sigma}_2^{(j)})^2} [y_j - \mathbf{m}_j^T \hat{\boldsymbol{\beta}}_2^{(j)}]^2} \right)^{z_{j2}}}{\frac{1}{3\sqrt{2\pi}\hat{\sigma}_1^{(j)}} e^{-\frac{1}{2(\hat{\sigma}_1^{(j)})^2} [y_j - \mathbf{m}_j^T \hat{\boldsymbol{\beta}}_1^{(j)}]^2} + \frac{2}{3\sqrt{2\pi}\hat{\sigma}_2^{(j)}} e^{-\frac{1}{2(\hat{\sigma}_2^{(j)})^2} [y_j - \mathbf{m}_j^T \hat{\boldsymbol{\beta}}_2^{(j)}]^2}}.$$

Diese Dichtefunktion wird zur Berechnung des Erwartungswertes  $E[x_{i1} | \mathbf{y}, \hat{\boldsymbol{\beta}}^{(j)}]$  verwendet. Man erhält

$$\begin{aligned} E[x_{i1} | \mathbf{y}, \hat{\boldsymbol{\beta}}^{(j)}] &= \sum_{\mathbf{z} \in \mathcal{Z}(\mathbf{y})} x_{i1} \cdot f(\mathbf{z} | \mathbf{y}, \hat{\boldsymbol{\beta}}^{(j)}) \\ &= \sum_{x_{N1}=0}^1 \cdots \sum_{x_{11}=0}^1 x_{i1} \cdot \prod_{j=1}^N \frac{W(y_j, \hat{\boldsymbol{\beta}}_1^{(j)}, \hat{\sigma}_1^{(j)})^{x_{j1}} \cdot W(y_j, \hat{\boldsymbol{\beta}}_2^{(j)}, \hat{\sigma}_2^{(j)})^{1-x_{j1}}}{W(y_j, \hat{\boldsymbol{\beta}}_1^{(j)}, \hat{\sigma}_1^{(j)}) + W(y_j, \hat{\boldsymbol{\beta}}_2^{(j)}, \hat{\sigma}_2^{(j)})} \quad \text{für } i = 1, 2, \dots, N \\ &= \sum_{x_{N1}=0}^1 \cdots \sum_{x_{11}=0}^1 \left\{ x_{i1} \cdot \frac{W(y_i, \hat{\boldsymbol{\beta}}_1^{(j)}, \hat{\sigma}_1^{(j)})^{x_{i1}} \cdot W(y_i, \hat{\boldsymbol{\beta}}_2^{(j)}, \hat{\sigma}_2^{(j)})^{1-x_{i1}}}{W(y_i, \hat{\boldsymbol{\beta}}_1^{(j)}, \hat{\sigma}_1^{(j)}) + W(y_i, \hat{\boldsymbol{\beta}}_2^{(j)}, \hat{\sigma}_2^{(j)})} \right. \\ &\quad \left. \cdot \prod_{\substack{j=1 \\ j \neq i}}^N \frac{W(y_j, \hat{\boldsymbol{\beta}}_1^{(j)}, \hat{\sigma}_1^{(j)})^{x_{j1}} \cdot W(y_j, \hat{\boldsymbol{\beta}}_2^{(j)}, \hat{\sigma}_2^{(j)})^{1-x_{j1}}}{W(y_j, \hat{\boldsymbol{\beta}}_1^{(j)}, \hat{\sigma}_1^{(j)}) + W(y_j, \hat{\boldsymbol{\beta}}_2^{(j)}, \hat{\sigma}_2^{(j)})} \right\} \end{aligned}$$

$$\begin{aligned}
&= \sum_{x_{N1}=0}^1 \cdots \sum_{x_{(i-1),1}=0}^1 \sum_{x_{(i+1),1}=0}^1 \cdots \sum_{x_{11}=1}^N \left\{ \frac{W(y_i, \hat{\beta}_1^{(j)}, \hat{\sigma}_1^{(j)})}{W(y_i, \hat{\beta}_1^{(j)}, \hat{\sigma}_1^{(j)}) + W(y_i, \hat{\beta}_2^{(j)}, \hat{\sigma}_2^{(j)})} \right. \\
&\quad \left. \cdot \prod_{\substack{j=1 \\ j \neq i}}^N \frac{W(y_j, \hat{\beta}_1^{(j)}, \hat{\sigma}_1^{(j)})^{x_{j1}} \cdot W(y_j, \hat{\beta}_2^{(j)}, \hat{\sigma}_2^{(j)})^{1-x_{j1}}}{W(y_j, \hat{\beta}_1^{(j)}, \hat{\sigma}_1^{(j)}) + W(y_j, \hat{\beta}_2^{(j)}, \hat{\sigma}_2^{(j)})} \right\} \\
&= \frac{W(y_i, \hat{\beta}_1^{(j)}, \hat{\sigma}_1^{(j)})}{W(y_i, \hat{\beta}_1^{(j)}, \hat{\sigma}_1^{(j)}) + W(y_i, \hat{\beta}_2^{(j)}, \hat{\sigma}_2^{(j)})} \cdot \frac{1}{\prod_{j \neq i}^N [W(y_j, \hat{\beta}_1^{(j)}, \hat{\sigma}_1^{(j)}) + W(y_j, \hat{\beta}_2^{(j)}, \hat{\sigma}_2^{(j)})]} \\
&\quad \cdot \underbrace{\sum_{x_{N1}=0}^1 \cdots \sum_{x_{(i-1),1}=0}^1 \sum_{x_{(i+1),1}=0}^1 \cdots \sum_{x_{11}=1}^N \prod_{\substack{j=1 \\ j \neq i}}^N [W(y_j, \hat{\beta}_1^{(j)}, \hat{\sigma}_1^{(j)})^{x_{j1}} \cdot W(y_j, \hat{\beta}_2^{(j)}, \hat{\sigma}_2^{(j)})^{1-x_{j1}}]}_{= \prod_{j \neq i}^N [W(y_j, \hat{\beta}_1^{(j)}, \hat{\sigma}_1^{(j)}) + W(y_j, \hat{\beta}_2^{(j)}, \hat{\sigma}_2^{(j)})] \quad \text{vgl. (422.30)}} \\
&= \frac{W(y_i, \hat{\beta}_1^{(j)}, \hat{\sigma}_1^{(j)})}{W(y_i, \hat{\beta}_1^{(j)}, \hat{\sigma}_1^{(j)}) + W(y_i, \hat{\beta}_2^{(j)}, \hat{\sigma}_2^{(j)})} \\
\Rightarrow E[x_{i1} | \mathbf{y}, \hat{\beta}^{(j)}] &= \frac{\frac{1}{3\sqrt{2\pi}\hat{\sigma}_1^{(j)}} e^{-\frac{1}{2(\hat{\sigma}_1^{(j)})^2} [y_i - \hat{\mathbf{m}}_1^T \hat{\beta}_1^{(j)}]^2}}{\frac{1}{3\sqrt{2\pi}\hat{\sigma}_1^{(j)}} e^{-\frac{1}{2(\hat{\sigma}_1^{(j)})^2} [y_i - \hat{\mathbf{m}}_1^T \hat{\beta}_1^{(j)}]^2} + \frac{2}{3\sqrt{2\pi}\hat{\sigma}_2^{(j)}} e^{-\frac{1}{2(\hat{\sigma}_2^{(j)})^2} [y_i - \hat{\mathbf{m}}_2^T \hat{\beta}_2^{(j)}]^2}}.
\end{aligned}$$

Analog gilt

$$E[x_{i2} | \mathbf{y}, \hat{\beta}^{(j)}] = \frac{\frac{2}{3\sqrt{2\pi}\hat{\sigma}_2^{(j)}} e^{-\frac{1}{2(\hat{\sigma}_2^{(j)})^2} [y_i - \hat{\mathbf{m}}_2^T \hat{\beta}_2^{(j)}]^2}}{\frac{1}{3\sqrt{2\pi}\hat{\sigma}_1^{(j)}} e^{-\frac{1}{2(\hat{\sigma}_1^{(j)})^2} [y_i - \hat{\mathbf{m}}_1^T \hat{\beta}_1^{(j)}]^2} + \frac{2}{3\sqrt{2\pi}\hat{\sigma}_2^{(j)}} e^{-\frac{1}{2(\hat{\sigma}_2^{(j)})^2} [y_i - \hat{\mathbf{m}}_2^T \hat{\beta}_2^{(j)}]^2}},$$

es werden also auf diese Weise die gleichen Ergebnisse erhalten wie in (422.27) und (422.28), womit diese Ergebnisse kontrolliert sind.

### M-Schritt: Maximierung der Kullback-Leibler-Statistik

Nachdem im E-Schritt des EM-Algorithmus die Erwartungswerte  $\mu_{ik}^{(j)}$  der Indikatoren  $x_{ik}$  für  $i \in \{1, \dots, N\}$  und  $k \in \{1, 2\}$  nach (422.27) und (422.28) berechnet worden sind, ist im M-Schritt muß die Kullback-Leibler-Statistik (422.29) durch Variation der unbekannt Parameter  $\beta = [\beta_1^T, \sigma_1, \beta_2^T, \sigma_2]^T$  zu maximieren. Maximiert werden muß also die Funktion

$$\begin{aligned}
Q(\beta | \hat{\beta}^{(j)}) &= \sum_{i=1}^N \left\{ \mu_{i1}^{(j)} \log \left( \frac{1}{3\sqrt{2\pi}\sigma_1} e^{-\frac{1}{2(\sigma_1)^2} [y_i - \hat{\mathbf{m}}_1^T \beta_1]^2} \right) \right. \\
&\quad \left. + \mu_{i2}^{(j)} \log \left( \frac{2}{3\sqrt{2\pi}\sigma_2} e^{-\frac{1}{2(\sigma_2)^2} [y_i - \hat{\mathbf{m}}_2^T \beta_2]^2} \right) \right\} \\
&= - \sum_{i=1}^N \left[ \mu_{i1}^{(j)} \log(3\sqrt{2\pi}\sigma_1) + \mu_{i2}^{(j)} \log(3/2\sqrt{2\pi}\sigma_2) \right] \\
&\quad - \sum_{i=1}^N \left( \frac{1}{2(\sigma_1)^2} \mu_{i1}^{(j)} [y_i - \hat{\mathbf{m}}_1^T \beta_1]^2 + \frac{1}{2(\sigma_2)^2} \mu_{i2}^{(j)} [y_i - \hat{\mathbf{m}}_2^T \beta_2]^2 \right).
\end{aligned}$$

Wie man sieht, hängt hierin nur die zweite Summe von den Modellparametern  $\beta_1$  und  $\beta_2$  ab, so daß es zur Bestimmung neuer Schätzwerte  $\hat{\beta}_1^{(j+1)}$  und  $\hat{\beta}_2^{(j+1)}$  für diese Parameter

zunächst ausreicht, die zweite Summe separat zu maximieren. Gesucht werden also jetzt Schätzwerte  $\hat{\beta}_1^{(j+1)}$  und  $\hat{\beta}_2^{(j+1)}$ , die die Summe

$$-\sum_{i=1}^N \left( \frac{1}{2\sigma_1^2} \mu_{i1}^{(j)} [y_i - \mathbf{m}_i^T \beta_1]^2 + \frac{1}{2\sigma_2^2} \mu_{i2}^{(j)} [y_i - \mathbf{m}_i^T \beta_2]^2 \right)$$

für  $\beta_1 = \hat{\beta}_1^{(j+1)}$  und  $\beta_2 = \hat{\beta}_2^{(j+1)}$  maximieren. Diese Summe kann wieder in zwei Teilsammen zerlegt werden, die getrennt voneinander zu maximieren sind:

$$-\frac{1}{2\sigma_1^2} \sum_{i=1}^N \mu_{i1}^{(j)} [y_i - \mathbf{m}_i^T \beta_1]^2 \rightarrow \max \quad \text{und} \quad -\frac{1}{2\sigma_2^2} \sum_{i=1}^N \mu_{i2}^{(j)} [y_i - \mathbf{m}_i^T \beta_2]^2 \rightarrow \max \quad (422.31)$$

Fasst man die Erwartungswerte  $\mu_{i1}^{(j)}$  bzw.  $\mu_{i2}^{(j)}$  der unzugänglichen Beobachtungen in den Gewichtsmatrizen  $\mathbf{P}_1^{(j)}$  bzw.  $\mathbf{P}_2^{(j)}$  zusammen, also

$$\mathbf{P}_1^{(j)} = \text{diag} \left( \mu_{11}^{(j)}, \mu_{21}^{(j)}, \dots, \mu_{N1}^{(j)} \right), \quad \mathbf{P}_2 = \text{diag} \left( \mu_{12}^{(j)}, \mu_{22}^{(j)}, \dots, \mu_{N2}^{(j)} \right) \quad (422.32)$$

und entsprechend die Koeffizientenvektoren  $\mathbf{m}_i^T$  bzw.  $\mathbf{m}_i^T$  der Beobachtungsgleichungen des Modells 1 (421.3) bzw. des Modells 2 (421.4) in den Matrizen  $\mathbf{M}^1$  und  $\mathbf{M}^2$ ,

$$\mathbf{M}^1 = \begin{bmatrix} \mathbf{m}_1^T \\ \mathbf{m}_2^T \\ \dots \\ \mathbf{m}_N^T \end{bmatrix}, \quad \mathbf{M}^2 = \begin{bmatrix} \mathbf{m}_1^T \\ \mathbf{m}_2^T \\ \dots \\ \mathbf{m}_N^T \end{bmatrix}, \quad (422.33)$$

so ergeben sich die Schätzwerte  $\hat{\beta}_1^{(j+1)}$  und  $\hat{\beta}_2^{(j+1)}$  aufgrund der Beziehungen

$$\hat{\beta}_1^{(j+1)} = \arg \max_{\beta_1} -\frac{1}{2\sigma^2} \left( \mathbf{y} - \mathbf{M}^1 \beta_1 \right)^T \mathbf{P}_1 \left( \mathbf{y} - \mathbf{M}^1 \beta_1 \right) \quad \text{und} \quad (422.34)$$

$$\hat{\beta}_2^{(j+1)} = \arg \max_{\beta_2} -\frac{1}{2\sigma^2} \left( \mathbf{y} - \mathbf{M}^2 \beta_2 \right)^T \mathbf{P}_2 \left( \mathbf{y} - \mathbf{M}^2 \beta_2 \right). \quad (422.35)$$

Aus den Gleichungen (422.31) und (422.35) ist ersichtlich, daß sich die neuen Schätzwerte  $\hat{\beta}_1^{(j+1)}$  und  $\hat{\beta}_2^{(j+1)}$  für die unbekannt Parameter  $\beta_1$  und  $\beta_2$  als Schätzwerte zweier Parameterschätzungen in den Modellen 1 und 2 nach der Methode der kleinsten Quadrate ergeben, wenn die Beobachtungen  $y_i$  bei der Ausgleichung im Modell 1 mit den Gewichten  $\mu_{i1}^{(j)}$  und bei der Ausgleichung im Modell 2 mit den Gewichten  $\mu_{i2}^{(j)}$  versehen werden. Gemäß [KOCH 1998], Satz (323.3) ergeben sich die Schätzungen

$$\hat{\beta}_1^{(j+1)} = (\mathbf{M}^{1T} \mathbf{P}_1^{(j)} \mathbf{M}^1)^{-1} \mathbf{M}^{1T} \mathbf{P}_1^{(j)} \mathbf{y} \quad \text{und} \quad \hat{\beta}_2^{(j+1)} = (\mathbf{M}^{2T} \mathbf{P}_2^{(j)} \mathbf{M}^2)^{-1} \mathbf{M}^{2T} \mathbf{P}_2^{(j)} \mathbf{y}. \quad (422.36)$$

Für die unbekannt Standardabweichungen  $\sigma_1$  und  $\sigma_2$  der zu den Modellen 1 und 2 gehörenden Beobachtungen sollen nun erwartungstreue Schätzungen  $\hat{\sigma}_1^{(j+1)}$  und  $\hat{\sigma}_2^{(j+1)}$  abgeleitet werden. Nach [KOCH 1998], Gl. (325.6) gilt

$$(\hat{\sigma}_1^{(j+1)})^2 = \frac{\Omega_1^{(j)}}{N - u_1} \quad \text{und} \quad (\hat{\sigma}_2^{(j+1)})^2 = \frac{\Omega_2^{(j)}}{N - u_2}, \quad (422.37)$$

worin  $\Omega_1^{(j)}$  und  $\Omega_2^{(j)}$  die Quadratsummen der geschätzten Residuen in den Modellen 1 und 2 bezeichnen, die sich im  $j$ -ten Iterationsschritt mit den Schätzungen  $\hat{\beta}_1^{(j+1)}$  bzw.  $\hat{\beta}_2^{(j+1)}$  ergeben.  $u_1$  und  $u_2$  bezeichnen jeweils die Anzahl der unbekannt Parameter in Modell 1 und Modell 2; hier gilt  $u_1 = 2$  und  $u_2 = 3$ .

Die Vektoren  $\hat{\mathbf{r}}_1^{(j)}$  und  $\hat{\mathbf{r}}_2^{(j)}$  der geschätzten Beobachtungsresiduen in den Modellen 1 und 2 ergeben sich in der  $j$ -ten EM-Iteration zu

$$\hat{\mathbf{r}}_1^{(j)} = \mathbf{M}^1 \hat{\beta}_1^{(j+1)} - \mathbf{y} \quad \text{und} \quad \hat{\mathbf{r}}_2^{(j)} = \mathbf{M}^2 \hat{\beta}_2^{(j+1)} - \mathbf{y}. \quad (422.38)$$

Hiermit erhält man in der  $j$ -ten EM-Iteration die Quadratsummen  $\Omega_1$  und  $\Omega_2$  der Beobachtungsresiduen

$$\Omega_1^{(j)} = (\hat{\mathbf{r}}_1^{(j)})^T \cdot \mathbf{P}_1^{(j)} \cdot \hat{\mathbf{r}}_1^{(j)} \quad \text{und} \quad \Omega_2^{(j)} = (\hat{\mathbf{r}}_2^{(j)})^T \cdot \mathbf{P}_2^{(j)} \cdot \hat{\mathbf{r}}_2^{(j)}. \quad (422.39)$$

Durch Einsetzen von (422.39) und (422.38) in (422.21) ergeben sich die Schätzwerte  $\hat{\sigma}_1^{(j+1)}$  und  $\hat{\sigma}_2^{(j+1)}$  für die Standardabweichungen der zu den Modellen 1 (Hochspannungsleitung) und 2 (Straße) gehörenden Beobachtungen

$$\hat{\sigma}_1^{(j+1)} = \left( \frac{(\mathbf{M}^1 \hat{\beta}_1^{(j+1)} - \mathbf{y})^T \cdot \mathbf{P}_1^{(j)} \cdot (\mathbf{M}^1 \hat{\beta}_1^{(j+1)} - \mathbf{y})}{N - 2} \right)^{\frac{1}{2}} \quad (422.40)$$

und

$$\hat{\sigma}_2^{(j+1)} = \left( \frac{(\mathbf{M}^2 \hat{\beta}_2^{(j+1)} - \mathbf{y})^T \cdot \mathbf{P}_2^{(j)} \cdot (\mathbf{M}^2 \hat{\beta}_2^{(j+1)} - \mathbf{y})}{N - 3} \right)^{\frac{1}{2}} \quad (422.41)$$

Mit (422.36), (422.40) und (422.41) sind alle im M-Schritt des EM-Algorithmus zu berechnenden Größen bestimmt.

### Abbruchkriterium, Kovarianzmatrizen der Schätzwerte und Klassifikation der Beobachtungen

**Abbruchkriterium:** Nachdem die im E-Schritt und im M-Schritt des zur Lösung der in Abschnitt 41 Schätz- und Klassifikationsaufgabe zu implementierenden EM-Algorithmus abgeleitet wurden, muß noch ein Kriterium eingeführt werden, aufgrund dessen am Ende jeder Iteration entschieden wird, ob die EM-Sequenz nochmals durchlaufen wird oder nicht. Zur Definition eines solchen Abbruchkriteriums wird hier die Quadratwurzel

$$\delta = \left( \|\hat{\beta}_1^{(j+1)} - \hat{\beta}_1^{(j)}\|^2 + \|\hat{\beta}_2^{(j+1)} - \hat{\beta}_2^{(j)}\|^2 \right)^{1/2} \quad (422.42)$$

aus der Quadratsumme der euklidischen Abstände der Schätzungen  $\hat{\beta}_1^{(j+1)}$  und  $\hat{\beta}_2^{(j+1)}$  der  $j$ -ten Iteration von den entsprechenden Schätzungen  $\hat{\beta}_1^{(j)}$  und  $\hat{\beta}_2^{(j)}$  der  $j - 1$ ten Iteration betrachtet. Bleibt diese Quadratsumme kleiner als ein vorzugebender Schwellwert  $\epsilon$ , was gleichbedeutend damit ist, daß sich die Schätzungen der unbekannt Parameter



zwischen zwei Iterationen um weniger als  $\epsilon$  ändern, so sollen die EM-Iterationen abgebrochen werden. Sonst wird mit der nächsten Iteration begonnen. Es ergibt sich also das **Abbruchkriterium**:

$$\text{Falls } \delta > \epsilon : \text{Übergang zur nächsten EM-Iteration} \quad (422.43)$$

$$\text{sonst : Abbruch der EM-Iterationen} \quad (422.44)$$

Für  $\epsilon$  wird im folgenden der Wert  $1 \cdot 10^{-9}$  gewählt. Die sich in der letzten EM-Iteration ergebende Schätzung  $\hat{\beta}^{(j+1)} = [(\hat{\beta}_1^{(j+1)})^T, \hat{\sigma}_1^{(j+1)}, (\hat{\beta}_2^{(j+1)})^T, \hat{\sigma}_2^{(j)}]$  für die unbekannt Parameter  $\hat{\beta} = [\hat{\beta}_1^T, \sigma_1, \hat{\beta}_2^T, \sigma_2]$  wird im folgenden mit  $\hat{\beta}^* = [(\hat{\beta}_1^{(*)})^T, \hat{\sigma}_1^{(*)}, (\hat{\beta}_2^{(*)})^T, \hat{\sigma}_2^{(*)}]^T$  bezeichnet.

**Kovarianzmatrizen der Schätzungen:** Schätzwerte für unbekannte Parameter, die im Rahmen einer Ausgleichung bestimmt werden, sind für sich allein nur beschränkt aussagekräftig, da aus den Schätzwerten zunächst nicht hervorgeht, wie sicher bzw. wie unsicher die Schätzung ist. Daher werden zum Abschluß der EM-Iterationen noch Schätzungen  $\hat{D}(\hat{\beta}_1^{(*)})$  und  $\hat{D}(\hat{\beta}_2^{(*)})$  für die Kovarianzmatrizen  $D(\hat{\beta}_1^{(*)})$  und  $D(\hat{\beta}_2^{(*)})$  der geschätzten Parameter  $\hat{\beta}_1^{(*)}$  und  $\hat{\beta}_2^{(*)}$  angegeben. Aus den Quadratwurzeln der Diagonalelemente dieser Matrizen gehen die geschätzten Standardabweichungen der Schätzwerte hervor, die hier als Maße für die Genauigkeit der einzelnen Schätzwerte dienen. Aus den Nebendiagonalelementen lassen sich Korrelationen zwischen den einzelnen Parametern ableiten, also Maße für statistische Abhängigkeiten zwischen den Schätzwerten. Nach [KOCH 1998], Gl. (325.7) ergeben sich die Schätzungen

$$\hat{D}(\hat{\beta}_1^{(*)}) = (\hat{\sigma}_1^{(*)})^2 \cdot (\mathbf{M}^T \mathbf{P}_1^{(*)} \mathbf{M})^{-1} \quad \text{und} \quad (422.45)$$

$$\hat{D}(\hat{\beta}_2^{(*)}) = (\hat{\sigma}_2^{(*)})^2 \cdot (\mathbf{M}^T \mathbf{P}_2^{(*)} \mathbf{M})^{-1} \quad (422.46)$$

für die beiden Kovarianzmatrizen der in den Modellen 1 und 2 erhaltenen Schätzwerte  $\hat{\beta}_1^{(*)}$  und  $\hat{\beta}_2^{(*)}$ .

**Klassifikation der Beobachtungen:** Die angestrebte Klassifikation der Beobachtungen, also die Zuordnung der einzelnen Beobachtungen jeweils entweder zur Hochspannungsleitung (Objekt  $\mathcal{O}_1$ ) oder zur Straße (Objekt  $\mathcal{O}_2$ ) kann mittels der sich nach der letzten EM-Iteration im E-Schritt ergebenden Erwartungswerte  $\mu_{i1}^{(*)}$  und  $\mu_{i2}^{(*)}$ ,  $i = 1, \dots, N$  auf einfache Weise erfolgen: Nach (422.25) und (422.26) sind  $\mu_{i1}^{(*)}$  und  $\mu_{i2}^{(*)}$  die Wahrscheinlichkeiten, daß die Beobachtung  $y_i$  zu Modell 1 (d.h. zur Hochspannungsleitung) bzw. zu Modell 2 (d.h. zur Straße) gehört. Es liegt nun nahe, jede Beobachtung  $y_i$  dem Modell  $M_k$ ,  $k \in \{1, 2\}$  zuzuordnen, für das die Wahrscheinlichkeit  $\mu_{ik}^{(*)}$  am größten ist<sup>7</sup>. Ist also  $y_i$  eine Beobachtung und ist  $o_i$  die Bezeichnung des Objektes, zu dem  $y_i$  gehört ( $o_i = 1$ , falls  $y_i$  zur Hochspannungsleitung gehört,  $o_i = 2$ , falls  $y_i$  zur Straße gehört), so ergeben sich die geschätzte Objektbezeichnungen  $\hat{o}_i$  für  $i = 1, \dots, N$  mit

$$\hat{o}_i = \begin{cases} 1 & , \text{ falls } \mu_{i1}^{(*)} > \mu_{i2}^{(*)} \\ 2 & , \text{ falls } \mu_{i1}^{(*)} \leq \mu_{i2}^{(*)} \end{cases} \quad , \quad i = 1, \dots, N \quad (422.47)$$

<sup>7</sup>Diese Vorgehensweise entspricht einer Maximum-Likelihood-Schätzung der Objektbezeichnung  $o_i$ ,  $o_i \in \{1, 2\}$  des zur Beobachtung  $y_i$  gehörenden Objektes

Die in diesem Abschnitt abgeleiteten Rechenformeln, die zur Lösung der zu Beginn formulierten Aufgabe mittels des EM-Algorithmus benötigt werden, sollen nun zusammengefasst werden.

### Zusammenfassung der Ergebnisse

In Übersicht (4.1) ist der zur Lösung der in Abschnitt 4.1 gestellten Aufgabe zu implementierende Algorithmus mit allen erforderlichen Rechenschritten dargestellt. Aus dem dargestellten Ablaufschema wird klar, daß es sich bei dem hier vorliegenden EM-Algorithmus um eine Parameterschätzung mit *iterativer Regewichtung der Beobachtungen* handelt: Im E-Schritt der  $j$ -ten Iteration werden auf Grundlage der in der vorhergehenden Iteration berechneten Schätzwerte  $\hat{\beta}^{(j)}$  die Beobachtungen für die nachfolgenden Parameterschätzungen in den Modellen 1 und 2 neu gewichtet: Eine Beobachtung  $y_i$ , die im  $j$ -ten Iterationsschritt besser durch das sich mit den Schätzwerten  $\hat{\beta}^{(j)}$  ergebende Modell 1 als durch das entsprechende Modell 2 beschrieben wird, wird bei der nächsten Parameterschätzung in Modell 1 herauf- und in Modell 2 herabgewichtet. Analog dazu wird sie bei der Neuschätzung in Modell 2 herauf- und in Modell 1 herabgewichtet, wenn sie besser durch das sich aus der vorhergehenden Iteration ergebende Modell 2 als durch das entsprechende Modell 1 beschrieben wird.

Auf Dauer führt diese iterative Neugewichtung der Beobachtungen dazu, daß Beobachtungen, die nicht zu dem Modell 1 (d.h. zur Hochspannungsleitung) gehören, bei der Parameterschätzung in diesem Modell ein solch geringes Gewicht erhalten, daß sie praktisch keinen Einfluß mehr auf die Schätzwerte für die unbekannt Parameter dieses Modells nehmen. Dann bestimmen nur noch die tatsächlich zur Hochspannungsleitung gehörenden Beobachtungen die Schätzwerte für die unbekannt Parameter des Modells, das die zur Hochspannungsleitung gehörenden Beobachtungen beschreibt. Genau so erhalten die nicht zur Straße gehörenden Beobachtungen bei den Parameterschätzungen in Modell 2 ein solch geringes Gewicht, daß sie sich auf die Schätzungen in diesem Modell kaum mehr auswirken. Somit ist der sich hier ergebende EM-Algorithmus eng verwandt mit Methoden zur robusten Parameterschätzung<sup>8</sup> (vgl. [KOCH 1998],Kap.3.8), bei denen Ausreißer in den Beobachtungsdaten durch iterative Regewichtung der Beobachtungen sukzessive herabgewichtet werden, bis ihr Einfluß auf das Schätzergebnis vernachlässigbar gering ist.

Der in diesem Abschnitt beschriebene Algorithmus wurde in der in Übersicht 4.1 dargestellten Form mittels der Programmiersprache C in einem Programm mit dem Namen `EM.c` umgesetzt. (Der Quellcode dieses Programms befindet sich in Anhang A1.) Das Programm liest eine Messreihe der Form

40 (Anzahl der Koordinatenpaare)	(u-Wert)	(y-Wert)
	1	15.18094175
	2	13.87552839
	3	6.03249538
	4	7.420570315
	5	9.613662336
	⋮	⋮

aus der Datei `Linien.dat` ein und bestimmt aus dem Datenmaterial in der oben beschrie-

<sup>8</sup>Es sei darauf hingewiesen, daß dies nur für das hier vorliegende Beispiel gilt.

benen Weise Schätzwerte für die unbekannt Parameter einer ausgleichenden Gerade und eines ausgleichenden Polynoms 2ten Grades. Hierbei erfolgt eine Klassifikation der Beobachtungen nach der Zugehörigkeit zum linearen Modell oder zum Polynommodell 2ten Grades. Die Ergebnisse werden in die Datei `Erg.dat` ausgegeben.

## 423 Ergebnisse

Das Programm `EM.c` wurde auf den Datensatz `Linien.dat` (vgl. Anhang A2) angewendet. Abbildung (2) visualisiert die eingelesenen Daten:

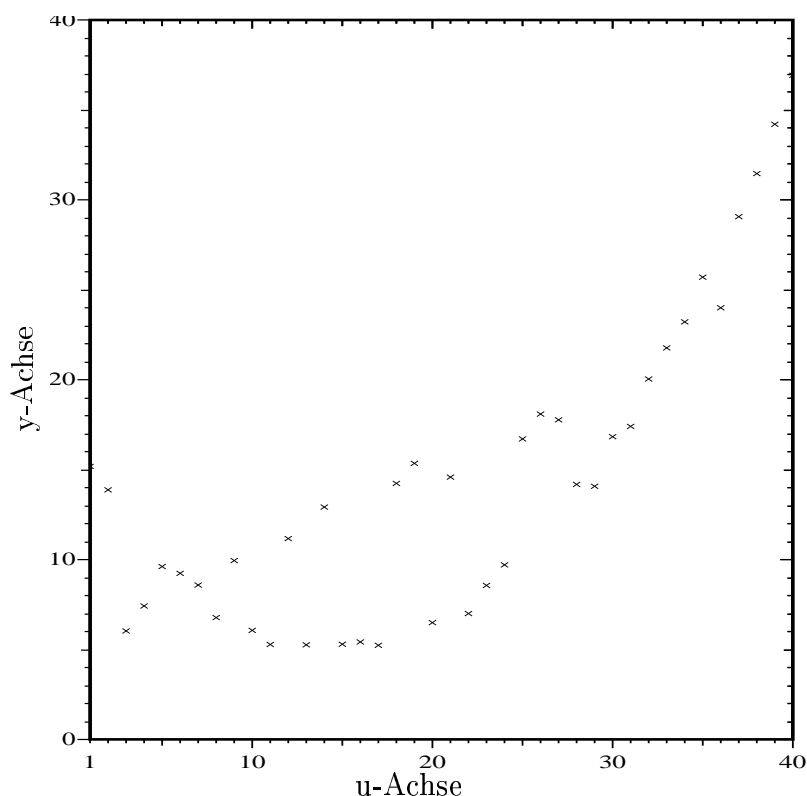


Abbildung 2: Visualisierung der Koordinatenpaare aus dem Datensatz `Linien.dat`

Mit diesen Eingangsdaten wurden mit dem Programm `EM.c` folgende Ergebnisse erhalten, wobei für das Abbruchkriterium der Schwellwert  $\epsilon = 1 \cdot 10^{-9}$  gewählt wurde:

### Schätzwerte des Modells 1 (Hochspannungsleitung):

Als *Schätzwerte für die Geradenparameter* wurden erhalten:

$$\hat{\beta}_1^{(*)} = \begin{bmatrix} \hat{a}_0 \\ \hat{a}_1 \end{bmatrix} = \begin{bmatrix} 5.1330 \\ 0.4985 \end{bmatrix} \quad \text{mit} \quad \begin{matrix} \hat{\sigma}_{a_0} = 0.214983 \\ \hat{\sigma}_{a_1} = 0.010525 \end{matrix} \quad (423.1)$$

Hierbei wurden die Standardabweichungen mittels der unten aufgeführten Kovarianzmatrix berechnet. Aufgrund der kleinen Standardabweichungen können die Schätzwerte als relativ sicher und die Modellparameter  $a_0$  und  $a_1$  als signifikant gelten<sup>9</sup>. Mit den Schätzwerten ergibt sich die *ausgleichende Gerade* im Bild, d.h. die Gerade, die die Mittellinie der Hochspannungsleitung im Bild approximiert, zu

$$\hat{y}_1(u) = 5.1330 + 0.4985 \cdot u \quad (423.2)$$

Als geschätzte Kovarianzmatrix des Vektors  $\hat{\beta}_1^{(*)}$  wurde die Matrix

$$\hat{D}(\hat{\beta}_1^{(*)}) = \begin{bmatrix} \hat{\sigma}_{a_0}^2 & \hat{\sigma}_{a_0, a_1} \\ \hat{\sigma}_{a_1, a_0} & \hat{\sigma}_{a_1}^2 \end{bmatrix} = \begin{bmatrix} 0.046218 & -0.001987 \\ -0.001987 & 0.000111 \end{bmatrix} \quad (423.3)$$

erhalten. Die geschätzte Standardabweichung der Gewichtseinheit der zur Hochspannungsleitung gehörenden Beobachtungen ergibt sich zu  $\hat{\sigma}_1^{(*)} = 0.3687$

### Schätzwerte des Modells 2 (Straße):

Es ergaben sich folgende Schätzwerte für die Koeffizienten des Polynoms 2ten Grades:

$$\hat{\beta}_2^{(*)} = \begin{bmatrix} \hat{b}_0 \\ \hat{b}_1 \\ \hat{b}_2 \end{bmatrix} = \begin{bmatrix} 16.3266 \\ -1.5097 \\ 0.0505 \end{bmatrix} \quad \text{mit} \quad \begin{array}{l} \hat{\sigma}_{b_0} = 0.214657 \\ \hat{\sigma}_{b_1} = 0.024097 \\ \hat{\sigma}_{b_2} = 0.000560 \end{array} \quad (423.4)$$

Auch hier sind die erhaltenen Standardabweichungen der Schätzwerte im Vergleich zu den Schätzwerten selbst klein, was ein Beleg für die Signifikanz der Modellparameter ist.<sup>9</sup> Weiter können aufgrund dieser kleinen Standardabweichungen die Schätzungen als relativ sicher gelten. Mit den Schätzwerten erhält man das *ausgleichende Polynom 2ten Grades* im Bild, d.h. das Polynom 2ten Grades, das die Mittellinie der Straße im Bild approximiert, zu

$$\hat{y}_2(u) = 16.3266 - 1.5097 \cdot u + 0.0505 \cdot u^2 \quad (423.5)$$

Als Geschätzte Kovarianzmatrix des Vektors  $\hat{\beta}_2^{(*)}$  wurde die Matrix

$$\hat{D}(\hat{\beta}_2^{(*)}) = \begin{bmatrix} \hat{\sigma}_{b_0}^2 & \hat{\sigma}_{b_0, b_1} & \hat{\sigma}_{b_0, b_2} \\ \hat{\sigma}_{b_1, b_0} & \hat{\sigma}_{b_1}^2 & \hat{\sigma}_{b_1, b_2} \\ \hat{\sigma}_{b_2, b_0} & \hat{\sigma}_{b_2, b_1} & \hat{\sigma}_{b_2}^2 \end{bmatrix} = \begin{bmatrix} 0.046078 & -0.004504 & 0.000091 \\ -0.004504 & 0.000581 & -0.000013 \\ 0.000091 & -0.000013 & 0.000003 \end{bmatrix} \quad (423.6)$$

erhalten. Die geschätzte Standardabweichung der Gewichtseinheit der zur Straße gehörenden Beobachtungen beträgt ebenfalls  $\hat{\sigma}_2^{(*)} = 0.3687$ . Dies deutet auf eine einheitliche Meßgenauigkeit der zur Hochspannungsleitung und zur Straße gehörenden Beobachtungen hin.

Die Ergebnisse der Klassifikation der Beobachtungen gehen aus Übersicht 423 hervor. Hierin gelten folgende Bezeichnungen:

$u$  : (feste) Ordinatenwerte der Koordinatenpaare des Eingabedatensatzes

$y$  : gemessene Abszissenwerte des Eingabedatensatzes

$\mu_{i1}^{(*)}, \mu_{i2}^{(*)}$  : Wahrscheinlichkeiten der Zugehörigkeit eines Koordinatenpaares zur HS-Leitung bzw. zur Straße

$\hat{o}_i$  : Klassifikationsergebnis: Geschätzte Objektbezeichnung des Objektes, zu der eine Beobachtung  $y_i$  gehört. (HS-Leitung:  $o_i = 1$ , Straße:  $o_i = 2$ )

$\hat{r}_{i1}, \hat{r}_{i2}$  : Beobachtungsresiduen im Modell 1 bzw. im Modell 2. Für jede Beobachtung ist das Residuum nur in dem Modell angegeben, dem die Beobachtung bei der Klassifikation zugeordnet wird.

$\hat{y}_{i1}, \hat{y}_{i2}$  : geschätzte Erwartungswerte der Beobachtungen. Der Erwartungswert wird nur in dem Modell angegeben, dem die Beobachtung bei der Klassifikation zugeordnet wird.

Es wird deutlich, daß alle Beobachtungsresiduen klein bleiben, was darauf hindeutet, daß keine Fehlklassifikationen von Beobachtungen vorliegen.

Die Ergebnisse von Parameterschätzung und Klassifikation sind in Abbildung 3 visualisiert. Hieraus wird offensichtlich, daß alle Koordinatenpaare des Eingangsdatensatzes den

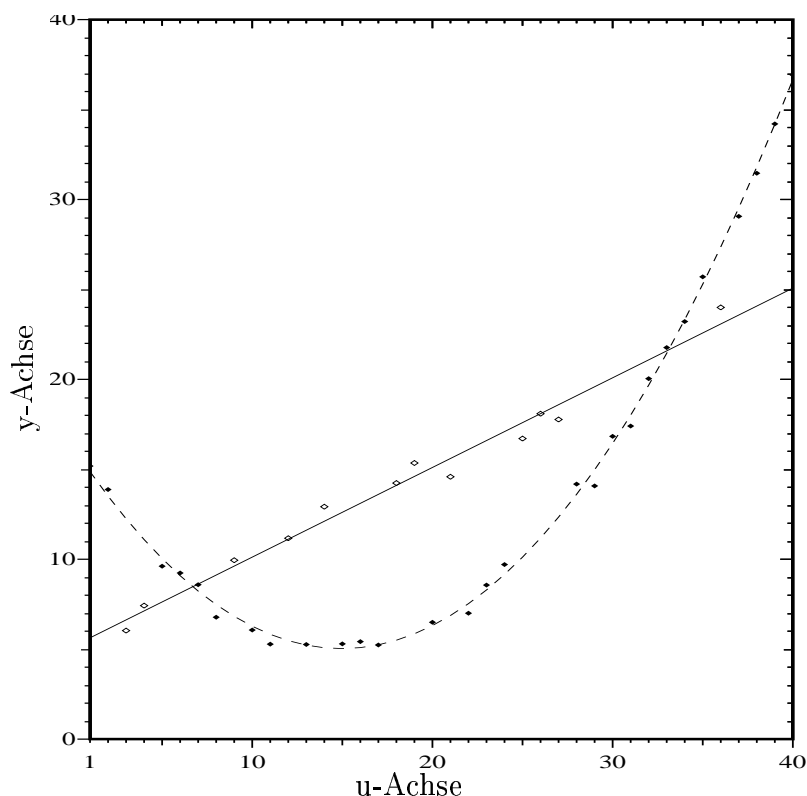


Abbildung 3: Visualisierung der Ergebnisse von Parameterschätzung und Klassifikation: Die durchgezogene Linie stellt die mittels des EM-Algorithmus bestimmte ausgleichende Gerade, die gestrichelte Linie das ausgleichende Polynom 2ten Grades dar. Die voll gezeichneten Punkte stellen Datenpunkte des Eingangsdatensatzes dar, die vom Programm als zur Straße (also zum Modell 2) gehörig klassifiziert wurden; die nicht ausgefüllten Punkte stellen Datenpunkte dar, die als zur Hochspannungsleitung (Objekt 1, Gerade) gehörig klassifiziert wurden.

Objekten 1 und 2 korrekt zugeordnet werden. Die beiden ausgleichenden Kurven approximieren das Datenmaterial augenscheinlich optimal.

Dieses optimale Ergebnis der Parameterschätzung und Klassifikation erklärt sich z.T. dadurch, daß sich die Objekte 1 und 2 in ihrer Form klar unterscheiden und daß die Messwertstreuung im Beispiel relativ klein gewählt wurde, so daß in Abbildung 2 die Beobachtungen klar dem einen oder anderen Objekt zugeschrieben werden können. Es wäre also noch zu testen, wie der Algorithmus auf stark verrauschten Daten arbeitet, bzw. was passiert, wenn die Objekte im Bild nicht so deutlich voneinander unterschieden werden können. (Z.B. wenn sich die zwei gerade Objektmittellinien im Bild schleifend schneiden.) Diese Untersuchungen sollen jedoch nicht zum weiteren Inhalt dieser Arbeit gemacht werden, da es hier nur darum geht, eine mögliche Anwendung des EM-Algorithmus *vorzustellen*.

<sup>9</sup>Die 3fache Standardabweichung jedes Schätzwertes ist kleiner als der Schätzwert selbst, so daß mit einer Irrtumswahrscheinlichkeit von weniger als 0,001 die Parameter des Modells von 0 verschieden sind.

Abschließend soll die iterative Regewichtung der Beobachtungen im Rahmen des hier angewandten EM-Algorithmus noch veranschaulicht werden: In Abbildung 423 sind beispielhaft die Gewichte  $\mu_{31}^{(j)}$  und  $\mu_{32}^{(j)}$  der Beobachtung  $y_3$  bei den Parameterschätzungen in den Modellen 1 und 2 gegen die laufenden Nummern  $j$  der Iterationen aufgetragen. Man er-

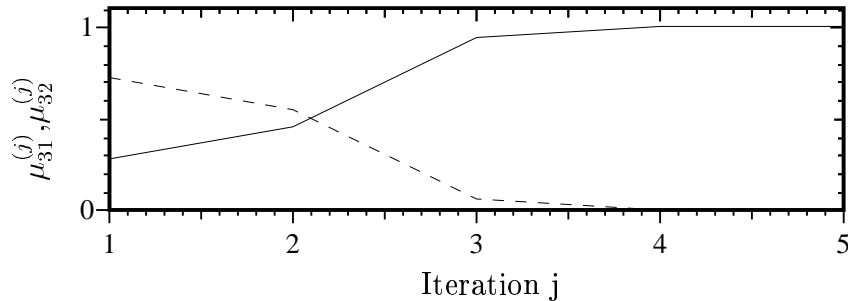


Abbildung 4: Visualisierung der iterativen Regewichtung der Beobachtungen im Rahmen des hier angewandten EM-Algorithmus. Die durchgezogene Linie stellt den Verlauf der Gewichtung der Beobachtung  $y_3$  bei den Parameterschätzungen im Modell 1 dar, die gestrichelte Linie zeigt den Verlauf der Gewichtung dieser Beobachtung bei den Schätzungen in Modell 2.

kennt, daß die Beobachtung  $y_3$  in Modell 1 zunächst ein geringes Gewicht hat. Im Rahmen des hier angewandten EM-Algorithmus wird sie in diesem Modell jedoch von Iteration zu Iteration z.T. drastisch heraufgewichtet und entsprechend im Modell 2 herabgewichtet, so daß sie sich am Ende nur noch auf die Schätzwerte der Parameterschätzungen in Modell 1 auswirkt. In Modell 2 konvergiert ihr Gewicht gegen Null. Die Gewichte konvergieren sehr schnell (nach 5 Iterationsschritten ist keine Veränderung mehr erkennbar), so daß der hier implementierte EM-Algorithmus relativ schnell konvergiert.

## Übersicht 4.1: Zusammenfassung der Rechenschritte

**EM-Algorithmus zur Lösung der Aufgabe:**

1. **Bestimmung von Näherungswerten:** Näherungswerte  $\hat{\beta}^{(0)}$  für die unbekannt Parameter  $\beta$  der beiden Modelle 1 und 2 werden durch zwei initiale Parameterschätzungen in diesen Modellen nach der Methode der kleinsten Quadrate erhalten. Als Gewichtsmatrix der Beobachtungen  $\mathbf{y}$  wird dabei jeweils die Einheitsmatrix  $\mathbf{P}_1^{(0)} = \mathbf{P}_2^{(0)} = \mathbf{I}$  gewählt. Es ergeben sich also für  $k = 1, 2$  die Näherungswerte

$$\hat{\beta}_k^{(0)} = (\mathbf{M}^T \mathbf{M})^{-1} \mathbf{M}^T \mathbf{y} \quad , \quad \sigma_k^{(0)} = [(\mathbf{M} \hat{\beta}_k^{(0)} - \mathbf{y})^T \cdot (\mathbf{M} \hat{\beta}_k^{(0)} - \mathbf{y}) / (N - k - 1)]^{1/2}$$

2. **EM-Iterationen:**

- i.) **E-Schritt:** Berechnung der Erwartungswerte  $\mu_{ik}^{(j)}$  für  $k = 1, 2$  und  $i = 1, \dots, N$  mittels

$$\mu_{ik}^{(j)} = \frac{\frac{k}{3\sqrt{2\pi}\hat{\sigma}_k^{(j)}} e^{-\frac{1}{2(\hat{\sigma}_k^{(j)})^2} [y_j - \mathbf{m}_j^T \hat{\beta}_k^{(j)}]^2}}{\frac{1}{3\sqrt{2\pi}\hat{\sigma}_1^{(j)}} e^{-\frac{1}{2(\hat{\sigma}_1^{(j)})^2} [y_j - \mathbf{m}_j^T \hat{\beta}_1^{(j)}]^2} + \frac{2}{3\sqrt{2\pi}\hat{\sigma}_2^{(j)}} e^{-\frac{1}{2(\hat{\sigma}_2^{(j)})^2} [y_j - \mathbf{m}_j^T \hat{\beta}_2^{(j)}]^2}}$$

- ii.) **M-Schritt:** Aufstellen der Matrizen  $\mathbf{P}_k^{(j+1)} = \text{diag}(\mu_{1k}^{(j+1)}, \mu_{2k}^{(j+1)}, \dots, \mu_{Nk}^{(j+1)})$  für  $k = 1, 2$  und Berechnung der Schätzungen

$$\begin{aligned} \hat{\beta}_k^{(j+1)} &= (\mathbf{M}^T \mathbf{P}_k^{(j+1)} \mathbf{M})^{-1} \mathbf{M}^T \mathbf{P}_k^{(j+1)} \mathbf{y} \quad \text{und} \\ \hat{\sigma}_k^{(j+1)} &= [(\mathbf{M} \hat{\beta}_k^{(j+1)} - \mathbf{y})^T \mathbf{P}_k^{(j+1)} (\mathbf{M} \hat{\beta}_k^{(j+1)} - \mathbf{y}) / (N - k - 1)]^{1/2} \end{aligned}$$

für  $k = 1, 2$

- iii.) **Abbruchkriterium:** Falls  $(\hat{\beta}_1^{(k+1)} - \hat{\beta}_1^{(k)})^2 + (\hat{\beta}_2^{(k+1)} - \hat{\beta}_2^{(k)})^2 < \epsilon$  gilt, wird der Algorithmus bei 3. fortgesetzt, sonst ist eine weitere EM-Iteration erforderlich, so daß der Algorithmus bei i.) fortgesetzt wird.

3. **Berechnung der Kovarianzmatrizen und Klassifikation der Beobachtungen:**

Für  $k = 1, 2$  ergeben sich die Kovarianzmatrizen der Schätzungen  $\hat{\beta}_1^{(*)}$  und  $\hat{\beta}_2^{(*)}$  zu

$$\hat{\mathbf{D}}(\hat{\beta}_k^{(*)}) = (\hat{\sigma}_k^{(*)})^2 \cdot (\mathbf{M}^T \mathbf{P}_k^{(*)} \mathbf{M})^{-1}.$$

Die geschätzten Bezeichnungen  $\hat{o}_i$  der Objekte, zu denen die Beobachtungen  $y_i$  gehören, ergeben sich mit

$$\hat{o}_i = \begin{cases} 1 \text{ ( HS-Leitung) } & , \text{ falls } \mu_{i1}^{(*)} > \mu_{i2}^{(*)} \\ 2 \text{ ( Straße) } & , \text{ falls } \mu_{i1}^{(*)} \leq \mu_{i2}^{(*)} . \end{cases} \quad , \quad i = 1, \dots, N$$

Übersicht 4.2: Klassifikationsergebnisse und Beobachtungsresiduen

u	y	$\mu_{i1}^{(*)}$	$\mu_{i2}^{(*)}$	$\hat{\sigma}_i$	$\hat{r}_{i1}$	$\hat{r}_{i2}$	$\hat{y}_{i1}$	$\hat{y}_{i2}$
1.0	15.181	0.000	1.000	2	–	0.000	–	14.867
2.0	13.876	0.000	1.000	2	–	0.000	–	13.509
3.0	6.032	1.000	0.000	1	0.596	–	6.628	–
4.0	7.421	1.000	0.000	1	0.000	–	7.127	–
5.0	9.614	0.000	1.000	2	–	0.426	–	10.040
6.0	9.236	0.006	0.994	2	–	0.000	–	9.085
7.0	8.585	0.440	0.560	2	–	0.000	–	8.231
8.0	6.777	0.000	1.000	2	–	0.701	–	7.478
9.0	9.945	1.000	0.000	1	0.000	–	9.619	–
10.0	6.065	0.000	1.000	2	–	0.210	–	6.275
11.0	5.280	0.000	1.000	2	–	0.545	–	5.825
12.0	11.168	1.000	0.000	1	0.000	–	11.115	–
13.0	5.261	0.000	1.000	2	–	0.000	–	5.227
14.0	12.919	1.000	0.000	1	0.000	–	12.111	–
15.0	5.295	0.000	1.000	2	–	0.000	–	5.033
16.0	5.421	0.000	1.000	2	–	0.000	–	5.087
17.0	5.234	0.000	1.000	2	–	0.009	–	5.243
18.0	14.235	1.000	0.000	1	0.000	–	14.105	–
19.0	15.350	1.000	0.000	1	0.000	–	14.604	–
20.0	6.495	0.000	1.000	2	–	0.000	–	6.314
21.0	14.583	1.000	0.000	1	1.017	–	15.601	–
22.0	7.000	0.000	1.000	2	–	0.532	–	7.532
23.0	8.564	0.000	1.000	2	–	0.000	–	8.293
24.0	9.705	0.000	1.000	2	–	0.000	–	9.154
25.0	16.707	1.000	0.000	1	0.887	–	17.595	–
26.0	18.089	1.000	0.000	1	0.004	–	18.093	–
27.0	17.768	1.000	0.000	1	0.824	–	18.591	–
28.0	14.179	0.000	1.000	2	–	0.000	–	13.609
29.0	14.075	0.000	1.000	2	–	0.900	–	14.975
30.0	16.832	0.000	1.000	2	–	0.000	–	16.442
31.0	17.404	0.000	1.000	2	–	0.606	–	18.010
32.0	20.039	0.014	0.986	2	–	0.000	–	19.679
33.0	21.767	0.388	0.612	2	–	0.000	–	21.448
34.0	23.214	0.004	0.996	2	–	0.105	–	23.319
35.0	25.699	0.000	1.000	2	–	0.000	–	25.290
36.0	23.998	1.000	0.000	1	0.000	–	23.078	–
37.0	29.062	0.000	1.000	2	–	0.474	–	29.536
38.0	31.459	0.000	1.000	2	–	0.351	–	31.810
39.0	34.200	0.000	1.000	2	–	0.000	–	34.185
40.0	36.924	0.000	1.000	2	–	0.000	–	36.661



# Kapitel 5

## Zusammenfassung

Die Photogrammetrie ist untrennbar mit Methoden der Ausgleichsrechnung und Statistik verbunden; dies gilt sowohl für die klassische Photogrammetrie als auch für moderne Verfahren der digitalen Bildverarbeitung. Klassische Methoden der geodätischen Ausgleichsrechnung reichen allerdings zur Lösung bestimmter Probleme aus der Photogrammetrie nicht immer aus, wie in Abschnitt 122 anhand eines einfachen Beispiels gezeigt wurde. Zur Lösung solcher Probleme müssen diese Parameterschätzverfahren an die konkrete Aufgabe angepaßt, erweitert, ergänzt oder durch andere Parameterschätzverfahren ersetzt werden. Bestimmte Aufgaben, bei denen die klassischen Schätzverfahren versagen, können mittels Schätzverfahren aus unvollständigen Beobachtungsdaten - wenn auch z.T. unter erheblichem Rechenaufwand - gelöst werden. Ein Beispiel hierfür sind Parameterschätzungen, bei denen die zugrundeliegenden Beobachtungen hinsichtlich ihrer Zugehörigkeit zu einem von mehreren alternativ gültigen Modellen unbestimmt sind (vgl. Abschnitt 122 und Kapitel 4).

Der in dieser Arbeit diskutierte EM-Algorithmus stellt ein solches Parameterschätzverfahren dar, mit dem Parameter aus unvollständigen Beobachtungsdaten geschätzt werden können. Er wurde erstmals im Jahre 1968 von den amerikanischen Professoren A. Dempster, N. Laird und D. Rubin unter dem Namen „EM-Algorithmus“ vorgestellt (vgl. [DEMPSTER ET AL. 1968]) und ergab sich damals als verallgemeinerte Zusammenfassung mehrerer Einzellösungen von Parameterschätzproblemen, die von verschiedenen Autoren entwickelt worden waren. Zwar gestaltet sich die Umsetzung des EM-Algorithmus in bezug auf eine konkrete Aufgabe häufig einfach, eine allgemeine Diskussion der zugrundeliegenden Theorie stellt jedoch eine sehr komplexe Aufgabe dar<sup>1</sup>. Die überwiegend englischsprachige Literatur zu diesem Thema richtet sich überwiegend an Mathematiker und Informatiker, deren Fachtermini sich z.T. von der im Rahmen der geodätischen Ausgleichsrechnung verwendeten Terminologie unterscheiden.

Um dem aus dem Bereich der Geodäsie stammenden Leser den Zugang zum EM-Algorithmus zu erleichtern, wurde in Abschnitt 13 dieser Arbeit ausgehend von den klassischen geodätischen Parameterschätzverfahren das Prinzip der Parameterschätzung aus unvollständigen Beobachtungsdaten erläutert. Anschließend wurde in Kapitel 2 schrittweise der EM-Algorithmus abgeleitet: Ausgehend von der in Abschnitt 22 beschriebenen Maximum-Likelihood-Methode, die die Grundlage des EM-Algorithmus darstellt, wurde

---

<sup>1</sup>Dies zeigt sich u.a. daran, daß –wie [WU 1982] nachwies – der Aufsatz der oben genannten und offensichtlich mit dem Thema vertrauten Professoren Dempster, Laird und Rubin Fehler enthält, die dazu führen, daß die in [DEMPSTER ET AL. 1968] getroffenen Aussagen zur Konvergenz des EM-Algorithmus ungültig sind.

zunächst in Abschnitt 222 gezeigt, wie die Maximum-Likelihood-Methode zur Schätzung von Parametern aus unvollständigen Beobachtungsdaten genutzt werden kann. Dann wurde in Abschnitt 231 die Schlüsselgleichung des (G)EM-Algorithmus abgeleitet und daraus schließlich der in Abschnitt 232 definierte (G)EM-Algorithmus entwickelt.

In Kapitel 3 wurden die Eigenschaften des (G)EM-Algorithmus in allgemeiner Form diskutiert, wobei Betrachtungen zur Konvergenz der logarithmierten Likelihoodfunktion bzw. deren Argument sowie die Konvergenzgeschwindigkeit im Vordergrund standen. Da in der diesbezüglichen Literatur einige Autoren verschiedene Auffassungen vertreten, wurden hier im wesentlichen die allgemein als richtig anerkannten Ergebnisse von [WU 1982] wiedergegeben. Die Darstellung wurde dabei relativ ausführlich gehalten, um dem nicht mit der englischsprachigen Terminologie der mathematischen Statistik bzw. Informatik vertrauten Leser den Zugang hierzu zu erleichtern.

Zunächst wurde in Abschnitt 321 das allgemeine Verhalten der logarithmierten Likelihoodfunktion im Rahmen der EM-Iterationen erklärt. Anschließend wurden in Abschnitt 322 die Bedingungen erarbeitet, unter denen die logarithmierte Likelihoodfunktion im Rahmen der EM-Iterationen gegen einen stationären Punkt bzw. gegen ein lokales Maximum konvergiert. In Abschnitt 323 wurden schließlich die Bedingungen genannt, unter denen die Folge der (vorläufigen) Schätzwerte für die unbekannt Parameter beim (G)EM-Algorithmus konvergiert.

Von besonderer Bedeutung für die Einsetzbarkeit eines Algorithmus ist nach der Frage, ob er zum Ziel führt, die Geschwindigkeit, mit dem er dies tut. In der Literatur (vgl. [DEMPSTER ET AL. 1968],[WU 1982]) wird dem EM-Algorithmus i.a. eine nur langsame Konvergenz zugeschrieben. Die Konvergenzgeschwindigkeit hängt dabei vom Anteil der fehlenden Beobachtungen an der Gesamtzahl der Beobachtungen ab. In Abschnitt 324 wurde ein Maß für die Konvergenzgeschwindigkeit des EM-Algorithmus abgeleitet.

In Kapitel 4 wurde das in der Einleitung skizzierte Beispiel für ein Schätz- und Klassifikationsproblem, aufgegriffen. Zur Lösung dieses Problems wurde es als Parameterschätzproblem aus unvollständigen Beobachtungsdaten formuliert und anschließend wurden die im E-Schritt und im M-Schritt durchzuführenden Berechnungen erläutert. Dies führte auf ein Ablaufschema zur Lösung der Schätz- gegebenen und Klassifikationsaufgabe, das in einem Programm zur Umsetzung gelangte (Vgl. Anhang A). Das Programm wurde auf einen Testdatensatz angewendet, wobei sehr gute Ergebnisse erzielt wurden. Dies kann u.a. neben der Qualität des Algorithmus auch dadurch begründet werden, daß das zugrundeliegende Datenmaterial eine nur geringe Streuung aufwies und daß die alternativ zur Auswahl stehenden Beobachtungsmodelle zur Parameterschätzung voneinander gut unterscheidbar waren. Es wäre noch zu testen, wie der Algorithmus auf stark verrauschten Meßdaten und sich ähnelnden Schätzmodellen arbeitet.

Diese Arbeit zeigt u.a., daß der EM-Algorithmus eine interessante Alternative zu herkömmlichen Verfahren der geodätischen Ausgleichsrechnung darstellt, falls diese zur Lösung einer Aufgabe ungeeignet ist. Aufgrund der im Allgemeinen nur langsamen Konvergenz des EM-Algorithmus sollte vor dessen Einsatz allerdings geprüft werden, ob sich das gewünschte Ergebnis nicht auch mit einer weniger rechenintensiven Methode erreichen läßt.

# Literaturverzeichnis

- [DEMPSTER ET AL. 1968] A. Dempster, N. Laird, D. Rubin. *Maximum Likelihood from Incomplete Data via the EM Algorithm*. Journal of the Royal Statistical Society, Series B (Methodological), Bd. 39, Nr. 1, 1977, S. 1-38.
- [FÖRSTNER 1991] W. Förstner. *Statistische Kenngrößen*, Institut für Photogrammetrie der Universität Bonn, 1991.
- [FUCHS 1998] C. Fuchs. *Extraktion polymorpher Bildstrukturen und ihre topologische und geometrische Gruppierung*, Dissertation am Institut für Photogrammetrie der Universität Bonn, 1998.
- [HELFRICH I 1995] H.-P. Helfrich. *Mathematik I*. Manuskript zur Vorlesung Mathematik I für Geodäten, Mathematisches Seminar der Landwirtschaftlichen Fakultät der Universität Bonn, 1995.
- [HELFRICH II 1996] H.-P. Helfrich. *Mathematik II*. Manuskript zur Vorlesung Mathematik II für Geodäten, Mathematisches Seminar der Landwirtschaftlichen Fakultät der Universität Bonn, 1996.
- [HINTON 1997] ] G.E. Hinton, M. Revow. *Using Mixtures of Factor Analyzers for Segmentation and Pose Estimation*. Department of Computer Science, University of Toronto, Canada, 1997.
- [HORNEGGER 1996] J. Hornegger. *Statistische Modellierung, Klassifikation und Lokalisierung von Objekten*, Shaker Verlag, Aachen, 1996.
- [JORDAN ET AL. 1994] M.I. Jordan, R.A. Jacobs. *Hierarchical mixtures of experts and the EM-algorithm*, Neural Computation, 6, S. 181 - 214, 1994.
- [KERNER ET AL. 1995] O. Kerner, J. Maurer, J. Steffens, T. Thode, R. Voller. *Vieweg Mathematik Lexikon*, 3. Aufl., Vieweg, 1995.
- [KOCH 1998] K.R. Koch. *Parameterschätzung und Hypothesentests*, Dümmler Verlag, Bonn, 1998.
- [KRAUS 1994] K. Kraus. *Photogrammetrie, Band 1, Grundlagen und Standardverfahren*, Dümmler Verlag, Bonn, 1994.

- [REDNER 1984] R.A.Redner, H.F.Walker. *Mixture Densities, Maximum Likelihood And The EM-Algorithm*, Society for Industrial and Applied Mathematics Review, Bd. 26, Nr. 2, S.195-239, 1984.
- [RAO 1965] C.R. Rao. *Linear Statistical Inference and its Applications*. Wiley-Verlag, New York, 1968.
- [TANNER 1993] M.A.Tanner. *Tools for Statistical Inference: Methods for the Exploration of Posterior Distributions and Likelihood Functions*, Springer Series in Statistics, Springer, Heidelberg, 1993.
- [WU 1982] C.F.J. Wu. *On The Convergence Properties Of The EM Algorithm*. The Annals of Statistics, Vol. 11, No.1, S. 95-103, 1983.
- [ZANGWILL 1969] W.I. Zangwill. *Nonlinear Programming: A Unified Approach*. Prentice Hall, Englewood Cliffs, New Jersey.

# Anhang A

## Anhang

### A1 C-Quellcode des Programms EM.c (Auszug)

Im Folgenden wird ein Auszug aus dem C-Quelltext des Programms EM.c wiedergegeben. Von den Ein- und Ausgaberoutinen, die hier aus Platzgründen nicht aufgeführt sind, werden nur die Funktionsprototypen angegeben.

```

/* Programm EM.c, Marc Luxen, 11.2000 */

#include<stdio.h>
#include<math.h>

/*Konstanten */
#define MaxDim 50                /*Maximale Matrixdimension */
#define pi 3.1415926
#define Epsilon 0.000000001     /*Schwellwert fuer Abbruchkriterium */

/*Datentypen*/
typedef struct
{
    int N_r, N_c;                /*Datentyp Matrix */
    double Elem[MaxDim][MaxDim];
}
t_Matrix;

/* Funktionsprototypen */
void Matrix_ausgeben (FILE * Disk, t_Matrix A);
t_Matrix Inverse (t_Matrix A);
t_Matrix Transponierte (t_Matrix A);
t_Matrix Produkt (t_Matrix A, t_Matrix B);
t_Matrix Differenz (t_Matrix A, t_Matrix B);
t_Matrix Einheitsmatrix (int n);
double Euklidnorm (t_Matrix X);
t_Matrix cProdukt (double c, t_Matrix A);
t_Matrix Schaetzung (t_Matrix X, t_Matrix y, t_Matrix P);
t_Matrix Residuen (t_Matrix X, t_Matrix y, t_Matrix beta);
t_Matrix Koeffizientenmatrix (int c, t_Matrix x);
void Messreihe_einlesen (char Dateiname[], t_Matrix * A, t_Matrix * B);
double W (double x, double y, double sigma, t_Matrix beta);

```

```

double mu (int c,
           double x,
           double y,
           double sigma1,
           t_Matrix beta1,
           double sigma2,
           t_Matrix beta2);

t_Matrix P (int c,
            t_Matrix u,
            t_Matrix y,
            double sigma1,
            t_Matrix beta1,
            double sigma2,
            t_Matrix beta2);

double Omega (t_Matrix X,
              t_Matrix y,
              t_Matrix P,
              t_Matrix beta);

void Plotdateien_schreiben (t_Matrix P1,    t_Matrix P2,
                           t_Matrix u,    t_Matrix y,
                           t_Matrix beta1, t_Matrix beta2);

void Ergebnisdatei_schreiben (int j,        t_Matrix u,
                              t_Matrix y,   double sigma1,
                              t_Matrix beta1, t_Matrix M1,
                              t_Matrix P1,   double sigma2,
                              t_Matrix beta2, t_Matrix P2,
                              t_Matrix M2,   t_Matrix Cov1,
                              t_Matrix Cov2);

/* Routinen:*/

void
Messreihe_einlesen (char Dateiname[], t_Matrix * A, t_Matrix * B)

/* Einlesen der Messreihe aus der Datei <Dateiname>.
   Die u-Werte der Messreihe werden in die Matrix A und
   die y-Werte in die Matrix B (beides N_r*1-Matrizen)
   geschrieben */

{
  FILE *Disk;
  int r;

  Disk = fopen (Dateiname, "r");

  /* Anzahl der Messwerte einlesen */
  fscanf (Disk, "%d", &((*A).N_r));

  /* Matrizenformate festlegen */
  (*A).N_c = 1;
  (*B).N_r = (*A).N_r;
  (*B).N_c = 1;

  /*Messwerte einlesen: */
  for (r = 0; r <= ((*A).N_r - 1); r++)

```

```

    {
        fscanf (Disk, "%lf", &((*A).Elem[r][0]));
        fscanf (Disk, "%lf", &((*B).Elem[r][0]));
    }
fclose (Disk);
}

void
Matrix_ausgeben (FILE * Disk, t_Matrix A)

/* Gibt die Matrix A in den Ausgabestrom Disk aus */

{
    int i, j;
    for (i = 0; i <= (A.N_r - 1); i++)
    {
        for (j = 0; j <= (A.N_c - 1); j++)
            fprintf (Disk, "%10.6lf    ", A.Elem[i][j]);
        fprintf (Disk, "\n");
    }
}

t_Matrix
Inverse (t_Matrix A)

/* Berechnet die Inverse der Matrix A nach der Gauss - Jordan - Methode
   (vgl. [Koch 1997], Gl. (133.21) */

{
    int i, j, k, r, c;          /* Laufvariablen */
    t_Matrix Inv;              /* Return - Variable

/* Format der Ausgabematrix festlegen */
Inv.N_r = A.N_r;
Inv.N_c = A.N_c;

/*Eliminationen */
for (i = 0; i <= (A.N_r - 1); i++)
    {
        for (j = 0; j <= (A.N_r - 1); j++)
            for (k = 0; k <= (A.N_r - 1); k++)
                {
                    if (k != i)
                        {
                            if (j != i)
                                Inv.Elem[j][k] = A.Elem[j][k] - A.Elem[j][i]
                                    * A.Elem[i][k] / A.Elem[i][i];
                            else
                                {
                                    Inv.Elem[i][k] = A.Elem[i][k] / A.Elem[i][i];
                                    Inv.Elem[k][i] = -A.Elem[k][i] / A.Elem[i][i];
                                }
                        }
                    }
                Inv.Elem[i][i] = 1 / A.Elem[i][i];
            }

    if (i < A.N_r - 1)
        for (r = 0; r <= (A.N_r - 1); r++)

```

```

        for (c = 0; c <= (A.N_c - 1); c++)
            A.Elem[r][c] = Inv.Elem[r][c];
    };
return (Inv);          /* Rueckgabe an aufrufende Routine */
}

```

t\_Matrix  
 Transponierte (t\_Matrix A)

```

/*Transponiert die Matrix A */

{
t_Matrix Trans;          /*Return - Matrix */
int i, j;                /*Laufvariablen */

/*Format der Ausgabematrix festlegen */
Trans.N_r = A.N_c;
Trans.N_c = A.N_r;

/*Element [i,j] der Return - Matrix wird das
  Element [j,i] der Matrix A zugewiesen: */
for (i = 0; i <= (Trans.N_r - 1); i++)
    for (j = 0; j <= (Trans.N_c - 1); j++)
        Trans.Elem[i][j] = A.Elem[j][i];
return (Trans);        /* Rueckgabe an aufrufende Routine */
}

```

t\_Matrix  
 Produkt (t\_Matrix A, t\_Matrix B)

```

/*Berechnet das Matrizenprodukt A*B der beiden Matrizen
  A und B und gibt es an die aufrufende Routine zurueck */

{
int r, c, k;            /* Laufvariablen */
t_Matrix Prod;        /* Return - Matrix */

/* Festlegung des Formates der Return - Matrix: */
Prod.N_r = A.N_r;
Prod.N_c = B.N_c;

/* Berechnung der Elemente der Return - Matrix: */
for (r = 0; r <= (Prod.N_r - 1); r++)
    for (c = 0; c <= (Prod.N_c - 1); c++)
        {
            Prod.Elem[r][c] = 0;
            for (k = 0; k <= (A.N_c - 1); k++)
                Prod.Elem[r][c] = Prod.Elem[r][c] + A.Elem[r][k] * B.Elem[k][c];
        };
return (Prod);        /* Rueckgabe an aufrufende Routine */
}

```

t\_Matrix  
 cProdukt (double c, t\_Matrix A)

```

/* Multipliziert die Matrix A mit dem Skalar c und gibt das
  Produkt an die aufrufende Routine zurueck */

```



```

{
    int z, s;                /*Laufvariablen */
    t_Matrix Prod;         /* Return - Matrix */

    /* Format der Return - Matrix festlegen: */
    Prod.N_r = A.N_r;
    Prod.N_c = A.N_c;

    /* Multiplikation der Matricelemente mit c, Prod[i][j]=c*A[i][j] */
    for (z = 0; z <= (A.N_r - 1); z++)
        for (s = 0; s <= (A.N_r - 1); s++)
            Prod.Elem[z][s] = c * A.Elem[z][s];

    return (Prod);        /* Rueckgabe an aufrufende Routine */
}

double
W (double u, double y, double sigma, t_Matrix beta)

/* Berechnet die Dichtefunktionen f(y_i|x_i=e_i,beta) und gibt sie
als double-Wert an die aufrufende Routine zurueck
y: Beobachtung
u: zu y gehoerender Abszissenwert*/

{
    if (beta.N_r == 3)    /* Entweder Berechnung der
                           Dichte des Modells 2... */
        return (2 / (3 * sqrt (2 * pi) * sigma)
                * exp (-pow (y - beta.Elem[0][0]
                            - beta.Elem[1][0] * u
                            - beta.Elem[2][0] * pow (u, 2), 2)
                    / (2 * pow (sigma, 2))));
    else                  /* ... oder Berechnung der
                           Dichte des Modells 1 */
        return (1 / (3 * sqrt (2 * pi) * sigma)
                * exp (-pow (y - beta.Elem[0][0]
                            - beta.Elem[1][0] * u, 2)
                    / (2 * pow (sigma, 2))));
}

double
mu (int c,                double x,
    double y,            double sigma1,
    t_Matrix beta1,     double sigma2,
    t_Matrix beta2)

/* Berechnung des Erwartungswertes mu_c im E-Schritt */

{
    double W1, W2;
    W1 = W (x, y, sigma1, beta1);
    W2 = W (x, y, sigma2, beta2);
    if (c == 1)
        return (W1 / (W1 + W2));
    else
        return (W2 / (W1 + W2));
}

```

```

t_Matrix
P (int c,          t_Matrix u,
   t_Matrix y,    double sigma1,
   t_Matrix beta1,
   double sigma2,
   t_Matrix beta2)

/* Baut die Matrix P1 bzw P2 auf und gibt sie an die aufrufende
   Routine zurueck
   falls c=1: Aufbauen der Matrix P1
   sonst      : Aufbauen der Matrix P2*/

{
  int i, j;                /* Laufvariablen */
  t_Matrix PMat;          /*Return - Matrix */

  /* Festlegung der Dimension der Return - Matrix */
  PMat.N_r = y.N_r;
  PMat.N_c = y.N_r;

  for (i = 0; i <= (PMat.N_r - 1); i++)
    for (j = 0; j <= (PMat.N_r - 1); j++)
      if (i == j)          /*Nur die Diagonalelemente
                           sind von 0 verschieden! */

/*Falls P1 aufgebaut werden Soll... */
      if (c == 1)

          /*Zuweisung der Erwartungswerte: */
          PMat.Elem[i][j] = mu (1, u.Elem[i][0], y.Elem[i][0], sigma1,
                                beta1, sigma2, beta2);

/*... sonst wird P2 aufgebaut */
      else
          PMat.Elem[i][j] = mu (2, u.Elem[i][0], y.Elem[i][0], sigma1,
                                beta1, sigma2, beta2);

      else
          PMat.Elem[i][j] = 0;
  return (PMat);          /* Rueckgabe an aufrufende Routine */
}

```

```

t_Matrix
Schaetzung (t_Matrix X, t_Matrix y, t_Matrix P)

/*Berechnet die Schaetzwerte eines Gauss-Markoff-Modells mit der
   Koeffizientenmatrix X, dem Beobachtungsvektor y und der
   Gewichtsmatrix P */

{
  t_Matrix XtP;           /*Dummy-Matrix */
  t_Matrix Inv;           /*Return - Matrix */
  XtP = Produkt (Transponierte (X), P);      /*Bekannte */
  Inv = Inverse (Produkt (XtP, X));          /* Gleichung */
  return (Produkt (Produkt (Inv, XtP), y));  /* Schuetzungen */
}

```

```

t_Matrix
Differenz (t_Matrix A, t_Matrix B)

    /* Berechnet die Differenz A-B zwischen den Matrizen A und B */

{
    int i, j;                /* Laufvariablen */
    t_Matrix Diff;          /*Return - Matrix */

    /*Festlegung der Dimension der Return - Matrix */
    Diff.N_r = A.N_r;
    Diff.N_c = A.N_c;

    /*Berechnen der Differenzen A[i,j]-B[i,j] */
    for (i = 0; i <= (A.N_r - 1); i++)
        for (j = 0; j <= (A.N_r - 1); j++)
            Diff.Elem[i][j] = A.Elem[i][j] - B.Elem[i][j];
    return (Diff);          /* Rueckgabe an aufrufende Routine */
}

t_Matrix
Residuen (t_Matrix X, t_Matrix y, t_Matrix beta)

    /* Berechnet die Beobachtungsresiduen eines GMM mit der
       Koeffizientenmatrix X, dem Beobachtungsvektor y und
       der Schaetzung beta */

{
    return (Differenz (Produkt (X, beta), y));    /*bekannte Formel */
}

double
Omega (t_Matrix X, t_Matrix y, t_Matrix P, t_Matrix beta)

    /*Berechnet die Quadratsumme der Residuen eines GMM mit der
       Koeffizientenmatrix X, dem Beobachtungsvektor y, der Gewichtsmatrix P
       und den Schaetzwerten beta */

{
    t_Matrix e;                /*Return - Matrix */

    e = Residuen (X, y, beta);    /*Bekannte */
    e = Produkt (Produkt (Transponierte (e), P), e);    /*Gleichung */
    return (e.Elem[0][0]);        /* Rueckgabe an aufrufende Routine */
}

t_Matrix
Koeffizientenmatrix (int c, t_Matrix u)

    /*Baut die Koeffizientenmatrix 1 (fuer c=1) bzw. 2 (fuer c=2) auf */
    /*u: Abszissenwert */

{
    int i, j;                /* Laufvariable */
    t_Matrix Koeff;          /* Return - Matrix */

    /*Dimension der Return - Matrix festlegen */

```

```

Koeff.N_c = c + 1;
Koeff.N_r = u.N_r;

/* Berechnung der Matricelemente */
for (i = 0; i <= (Koeff.N_r - 1); i++)
{
    /* Die ersten beiden Koeffizienten sind in beiden
       Koeffizientenmatrizen dieselben.... */
    Koeff.Elem[i][0] = 1;
    Koeff.Elem[i][1] = u.Elem[i][0];

    /*Bei der Koeffizientenmatrix 2 kommt noch ein dritter
       Koeffizient hinzu */
    if (c == 2)
        Koeff.Elem[i][2] = pow (u.Elem[i][0], 2);
};
return (Koeff);          /* Rueckgabe an aufrufende Routine */
}

t_Matrix
Einheitsmatrix (int n)

/* Erzeugt eine n*n - Einheitsmatrix */

{
    int r, c;              /*Laufvariablen */
    t_Matrix I;           /* Return - Matrix */

    /* Dimensionen der Return - Matrix festlegen*/
    I.N_r = n;
    I.N_c = n;

    /*Berechnung der Matricelemente */
    for (r = 0; r < n; r++)
        for (c = 0; c < n; c++)
            if (r == c)
                I.Elem[r][c] = 1;
            else
                I.Elem[r][c] = 0;
    return (I);          /* Rueckgabe an aufrufende Routine */
}

double
Euklidnorm (t_Matrix X)

/*Berechnet die euklidische Vektornorm eines Vektors X */

{
    double N = 0.0;       /*Return - Variable wird initialisiert */
    int i;                /* Laufvariable */

    /* Quadratsumme der Vektorelemente bilden: */
    for (i = 0; i <= (X.N_r - 1); i++)
        N = N + pow (X.Elem[i][0], 2);
    return (sqrt (N));    /* Rueckgabe an aufrufende Routine */
}

```

```
int
Klasse (double a, double b)

/* Gibt aus, welche von den beiden Zahlen a und b die groessere ist
   (Zahl 1 oder Zahl 2)
   Von der aufrufenden Routine aus sind a und b hier
   Wahrscheinlichkeiten des Zutreffens des Modells 1 oder 2 fuer
   eine Beobachtung. Daher nimmt diese Funktion hier einer
   Klassifikationsentscheidung vor. */

{
  if (a > b)
    return (1);
  else
    return (2);
}

void
main (void)
{
  t_Matrix u; /*Vektor der Abszissenwerte */
  t_Matrix y; /*Vektor der Ordinatenwerte (Beobachtungsvektor) */
  t_Matrix M1; /*Koeffizientenmatrix des Modells 1 */
  t_Matrix M2; /*Koeffizientenmatrix des Modells 2 */
  t_Matrix P1; /*Gewichtsmatrix des Modells 1 */
  t_Matrix P2; /*Gewichtsmatrix des Modells 2 */
  t_Matrix beta1; /*Geschaetzter Parametervektor des Modells 1 */
  t_Matrix Hbeta1; /*Dummy-Parametervektor (Hilfsvektor) */
  t_Matrix beta2; /*Geschaetzter Parametervektor des Modells 2 */
  t_Matrix Hbeta2; /*Dummy-Parametervektor (Hilfsvektor) */
  t_Matrix Cov1; /*Geschaetzte Kovarianzmatrix der Parameter beta1 */
  t_Matrix Cov2; /*Geschaetzte Kovarianzmatrix der Parameter beta2 */

  int i, j; /*Laufvariablen */
  double sigma1, sigma2; /*Standardabweichungen der Gewichtseinheit */
  FILE *Disk; /*Ausgabestrom fuer Zwischenergebnisse */

  Messreihe_einlesen ("/home/Marc/Arbeit/Software/Linien.dat", &u, &y);
  Disk = fopen ("/home/Marc/Arbeit/Software/P1.dat", "w");

  /* Naehierungswerte */

  /* Als Gewichtsmatrizen werden zunaechst Einheitsmatrizen gewaehlt */
  P1 = Einheitsmatrix (u.N_r);
  P2 = P1;

  /* Aufbauen der Koeffizientenmatrizen : */
  M1 = Koeffizientenmatrix (1, u);
  M2 = Koeffizientenmatrix (2, u);

  /* Initiale Parameterschaetzungen in den Modellen 1 und 2 */
  beta1 = Schaetzung (M1, y, P1);
  beta2 = Schaetzung (M2, y, P2);
  /* Bestimmung der Varianzen der Gewichtseinheit */
  sigma1 = sqrt (Omega (M1, y, P1, beta1) / (y.N_r - 2));
  sigma2 = sqrt (Omega (M2, y, P2, beta2) / (y.N_r - 3));

  j = 0; /* Iterationszaehler j wird auf 0 gesetzt */
```

```

do
{
  j++;
  /*Um einen Vergleich der neuen Schaetzwerte mit den alten
    Schaetzwerten zu ermoeeglichen, muessen die alten
    Schaetzwerte zwischengespeichert werden */
  Hbeta1 = beta1;
  Hbeta2 = beta2;

  /* E-Schritt: Aufbau der Gewichtsmatrizen P1 und P2 */
  P1 = P (1, u, y, sigma1, beta1, sigma2, beta2);
  P2 = P (2, u, y, sigma1, beta1, sigma2, beta2);

  /* M-Schritt: Neuschaetzung der unbekanntnen Parameter mit
    der neuen Gewichtsmatrix */
  beta1 = Schaetzung (M1, y, P1);
  beta2 = Schaetzung (M2, y, P2);
  /* Neuschaetzungen der Varianzen der Gewichtseinheit */
  sigma1 = sqrt (Omega (M1, y, P1, beta1) / (y.N_r - 2));
  sigma2 = sqrt (Omega (M2, y, P2, beta2) / (y.N_r - 3));

  /*Ausgabe von Zwischenergebnissen fuer Visualisierung */
  fprintf (Disk, "%d ", j);
  for (i = 0; i <= (y.N_r - 1); i++)
    fprintf (Disk, "%8.6lf ", P1.Elem[i][i]);
  fprintf (Disk, "\n");
}

/* Abbruchbedingung: */
while (sqrt (pow (Euklidnorm (Differenz (beta1, Hbeta1)), 2)
  + pow (Euklidnorm (Differenz (beta2, Hbeta2)), 2))
  >= Epsilon);

fclose (Disk); /* Schliessen des Ausgabestroms fuer Zwischenergebnisse */

/*Berechnung der Kovarianzmatrizen der finalen Schaetzungen beta1
  und beta2: */
Cov1 = cProdukt (pow (sigma1, 2),
  Inverse (Produkt (Produkt (Transponierte (M1), P1), M1)));

Cov2 = cProdukt (pow (sigma2, 2),
  Inverse (Produkt (Produkt (Transponierte (M2), P2), M2)));

/*Ausgabe der Ergebnisse */
Ergebnisdatei_schreiben (j, u, y, sigma1, beta1, M1, P1,
  sigma2, beta2, P2, M2, Cov1, Cov2);

Plotdateien_schreiben (P1, P2, u, y, beta1, beta2);
}

```

## A2 Datensatz Linien.dat

Den Berechnungen in Abschnitt 423 liegen folgende Daten zugrunde:

```
40
1      15.18094175
2      13.87552839
3      6.03249538
4      7.420570315
5      9.613662336
6      9.235540291
7      8.584545608
8      6.776858258
9      9.945421688
10     6.064615363
11     5.280097939
12     11.16776634
13     5.261352972
14     12.91940046
15     5.29507723
16     5.420944694
17     5.233561958
18     14.2346887
19     15.34965322
20     6.49533126
21     14.58344748
22     6.999892587
23     8.564241033
24     9.705304727
25     16.70748849
26     18.0886264
27     17.76763562
28     14.17874028
29     14.07531348
30     16.8324473
31     17.40444817
32     20.03882833
33     21.76675137
34     23.21420071
35     25.69938612
36     23.9982072
37     29.06170133
38     31.45851914
39     34.19962462
40     36.9238597
```