

BestMasters

Michael Seitz

Intervalldaten und generalisierte lineare Modelle



Springer Spektrum

BestMasters

Mit „BestMasters“ zeichnet Springer die besten Masterarbeiten aus, die an renommierten Hochschulen in Deutschland, Österreich und der Schweiz entstanden sind. Die mit Höchstnote ausgezeichneten Arbeiten wurden durch Gutachter zur Veröffentlichung empfohlen und behandeln aktuelle Themen aus unterschiedlichen Fachgebieten der Naturwissenschaften, Psychologie, Technik und Wirtschaftswissenschaften.

Die Reihe wendet sich an Praktiker und Wissenschaftler gleichermaßen und soll insbesondere auch Nachwuchswissenschaftlern Orientierung geben.

Michael Seitz

Intervalldaten und generalisierte lineare Modelle

Mit einem Geleitwort von Prof. Dr. Thomas Augustin

 Springer Spektrum

Michael Seitz
München, Deutschland

BestMasters

ISBN 978-3-658-08745-6

ISBN 978-3-658-08746-3 (eBook)

DOI 10.1007/978-3-658-08746-3

Die Deutsche Nationalbibliothek verzeichnet diese Publikation in der Deutschen Nationalbibliografie; detaillierte bibliografische Daten sind im Internet über <http://dnb.d-nb.de> abrufbar.

Springer Spektrum

© Springer Fachmedien Wiesbaden 2015

Das Werk einschließlich aller seiner Teile ist urheberrechtlich geschützt. Jede Verwertung, die nicht ausdrücklich vom Urheberrechtsgesetz zugelassen ist, bedarf der vorherigen Zustimmung des Verlags. Das gilt insbesondere für Vervielfältigungen, Bearbeitungen, Übersetzungen, Mikroverfilmungen und die Einspeicherung und Verarbeitung in elektronischen Systemen.

Die Wiedergabe von Gebrauchsnamen, Handelsnamen, Warenbezeichnungen usw. in diesem Werk berechtigt auch ohne besondere Kennzeichnung nicht zu der Annahme, dass solche Namen im Sinne der Warenzeichen- und Markenschutz-Gesetzgebung als frei zu betrachten wären und daher von jedermann benutzt werden dürften.

Der Verlag, die Autoren und die Herausgeber gehen davon aus, dass die Angaben und Informationen in diesem Werk zum Zeitpunkt der Veröffentlichung vollständig und korrekt sind. Weder der Verlag noch die Autoren oder die Herausgeber übernehmen, ausdrücklich oder implizit, Gewähr für den Inhalt des Werkes, etwaige Fehler oder Äußerungen.

Gedruckt auf säurefreiem und chlorfrei gebleichtem Papier

Springer Fachmedien Wiesbaden ist Teil der Fachverlagsgruppe Springer Science+Business Media (www.springer.com)

Geleitwort

Seit längerem herrscht in der angewandten Statistik ein Unbehagen, dass viele gängige statistische Verfahren eigentlich voraussetzen, dass die Daten in einer Qualität und Präzision vorliegen, die oft nicht gewährleistet werden kann. Gerade in großen Studien zu komplexen Themen treten häufig nicht nur Antwortverweigerungen auf, sondern die erhobenen Daten sind zusätzlich durch Antworteffekte verzerrt und unpräzise. Ein charakteristisches Beispiel sind *Intervalldaten*: In Surveys runden die Befragten oft grob und bevorzugen bestimmte “attraktive Werte”. Ferner wird bei der Konzeption von Fragebögen häufig empfohlen, Angaben zu heiklen Fragen (wie etwa dem Einkommen) von vornherein nur in Intervallen zu erheben, um Verweigerungen zu vermeiden, oder es werden von primären Verweigerern wenigstens Intervallangaben erbeten. Bei ereignisanalytischen Modellen kommt zusätzlich oft noch das durch individuelle Beobachtungszeitpunkte induzierte Problem der Intervallzensur hinzu.

Sensibilisiert dafür, dass eine solche “Defizitäten” negierende, naive Anwendung von statistischen Verfahren unter Umständen zu gravierenden Verzerrungen der Ergebnisse – und damit zu folgenschweren inhaltlichen Fehlinterpretationen – führen kann, wurden in der Literatur einerseits genauere Kriterien gewonnen, wann mit starken Verzerrungen zu rechnen ist, sowie entsprechende Korrekturverfahren entwickelt, um die Verzerrungen durch Messfehler, Vergrößerungen und fehlende Daten zu kompensieren.

Diese theoretisch sehr leistungsfähigen Korrekturmethode stellen einen wichtigen Fortschritt dar; ihre praktische Relevanz bleibt aber mit dem Makel behaftet, dass die Korrekturen typischerweise auf weitreichenden Zusatzannahmen über den “Defizitäts”-prozess beruhen, die empirisch nicht überprüfbar und oft auch inhaltlich nicht begründbar sind. Deshalb findet in den letzten Jahren insbesondere im Bereich der systematisch fehlenden Daten eine prinzipiell andere Vorgehensweise immer mehr Anklang, die v.a. unter dem Begriff partielle Identifikation bekannt wurde: Man verzichtet bewusst auf die vermeintliche Präzision von unter Einbezug von artifiziellen Zusatzannahmen erzielten Modellschätzungen und lässt “mengenwertige Ergebnisse” zu, indem man die Menge aller mit den Daten (und inhaltlich unbezweifelten Zusatzannahmen) verträglichen Modelle betrachtet. Diese so erzielten

Ergebnisse sind zwar häufig unpräzise, aber durch die Art ihrer Gewinnung glaubwürdiger und zuverlässiger, und, wie sich mittlerweile in vielfältigen Studien gezeigt hat, oft immer noch präzise genug, um die eigentlichen substanzwissenschaftlichen Fragestellungen beantworten zu können.

Die Arbeit von Herrn Seitz ist genau diesem neuem Paradigma einer zuverlässigen Inferenz unter unvollständiger Information verpflichtet. In vielen Situationen erscheint es bei Intervalldaten äußerst bedenklich, anzunehmen, der Vergrößerungsmechanismus sei nichtinformativ zufällig, und so sind Intervalldaten ein natürlicher Anwendungsbereich für Methoden der partiellen Information. In der Tat gibt es eine Reihe von entsprechenden Ansätzen für lineare Regressionsmodelle unter Intervalldaten, aber überraschenderweise praktisch keine Literatur zu verallgemeinerten linearen Modellen i. e. S., wie sie sich mittlerweile in der statistischen Modellierung als Standardmethode durchgesetzt haben. Herr Seitz präsentiert hier einen sehr allgemeinen Ansatz, der sich im Prinzip auch noch allgemeiner auf beliebige unverzerrte Schätzgleichungen ausdehnen lässt. Nach einer eher als Kontrollsituation dienenden direkten Lösung im Fall eines eindimensionalen Parameters wird die allgemeine Aufgabe konsequent als Optimierungsproblem über die Parameter gesehen, wobei die Intervalldaten in natürlicher Weise lineare Restriktionen formulieren. Um genau die Menge aller Maximum-Likelihood-Schätzer als zulässige Lösungen zu erhalten, wird die Bedingung "jeweiliger Wert der Scorefunktion = 0" zunächst als nichtlineare Nebenbedingung eingebracht. Die Arbeit entwickelt sodann mehrere Methoden, wie dieser Ansatz konkret operational gemacht werden kann. Eine Idee unter anderen besteht darin, diese Nebenbedingungen mit einem Pönalisierungstrick in die Zielfunktion mitaufzunehmen, so dass nun ein nichtlineares Optimierungsproblem über einem konvexen Polyeder entsteht. Die verschiedenen Methoden werden v.a. für die Exponentialregression konkret implementiert, verglichen und sorgfältig evaluiert.

Die Ergebnisse der Arbeit von Herrn Seitz erlauben erstmalig die konkrete Implementation von Methoden der partiellen Identifikation in generalisierten linearen Modellen unter Intervalldaten und erweitern damit den möglichen Anwendungsbereich dieses neuen methodischen Paradigmas substantiell. Sie sind ein wichtiger Meilenstein auf dem Weg zu einem allgemeinen Rahmen für eine zuverlässige statistische Modellierung komplexer Daten.

Prof. Dr. Thomas Augustin

Vorwort

Mit der Theorie der partiellen Identifizierung wird Unsicherheit in den beobachteten Daten auf die Schätzung der Modellparameter übertragen. Im Kontext der generalisierten Regression mit Intervalldaten erhält man für die Parameterschätzer so keine skalaren Werte, sondern ebenfalls Intervalle. Im Allgemeinen kann dies als Optimierungsproblem mit Nebenbedingungen formuliert werden. Die Herausforderung besteht hier insbesondere darin, die globalen Extrema zu bestimmen. Für Spezialfälle sind bereits Lösungen aus der Literatur bekannt: Bei der linearen Regression mit skalarer abhängiger Variable kann das Optimierungsproblem analytisch gelöst werden. Für eindimensionale Parameter in generalisierten linearen Modellen wird in der vorliegenden Arbeit ein neuer Ansatz realisiert und untersucht: Optimiert man die Score-Funktion mit festem Parameter über die beobachteten Intervalldaten, so lassen sich dadurch die Extrema des Parameterschätzers bestimmen. Dieser Ansatz liefert zuverlässig die richtige Lösung. Als allgemeine Herangehensweise kann der Parameterschätzer direkt numerisch optimiert werden. Daneben wird ein alternativer Ansatz verwendet: Um das Intervall der zulässigen Schätzer zu erhalten, kann die Score-Funktion als Strafterm in die Zielfunktion integriert werden. Die Methoden werden für das lineare Modell und das generalisierte lineare Modell mit Exponentialverteilung und log-Link umgesetzt. Bei der Untersuchung von Simulationsbeispielen zeigt sich, dass die allgemeinen Verfahren teilweise nur lokale Extrema finden. Da die Zielfunktionen aber mitunter sehr hochdimensional sind, kann die numerische Optimierung nicht an verschiedenen Stellen systematisch neugestartet werden. Auf Grund dieser Problematik werden heuristische Methoden vorgeschlagen, bei denen iterativ in den Ecken der Datenintervalle neugestartet wird. Diese liefern durchwegs die größten Intervalle der zulässigen Parameter und damit die beste Lösung. Die Methode wird auf ein Anwendungsbeispiel aus der Volkswirtschaft angewandt: Vorerst wird ein Regressionsmodell für den Zusammenhang zwischen dem Bruttoinlandsprodukt und den Ausgaben für Forschung und Entwicklung verschiedener Länder entwickelt. Anschließend werden Intervalle für die Daten konstruiert und die entsprechenden Parameterintervalle berechnet. Mit den vorgestellten und praktisch umgesetzten Verfahren lassen sich oft die exakten Lösungen und in allen untersuchten Fällen sehr gute Approximationen finden. Des Weiteren wird demonstriert, dass

die Ansätze grundsätzlich zur Lösung des Problems führen, wenn die Extrema der Optimierungsprobleme richtig bestimmt werden.

Michael Seitz

Inhaltsverzeichnis

Abbildungsverzeichnis	xi
Tabellenverzeichnis	xv
Algorithmenverzeichnis	xvii
1. Einleitung	1
2. Intervalldaten und generalisierte lineare Modelle	6
2.1. Generalisierte lineare Regression	6
2.2. Regression mit Intervalldaten	8
2.3. Notation	10
2.4. Formulierung als Optimierungsproblem	11
2.5. Optimierung mit Strafterm	12
2.6. Analytische Lösung für die lineare Regression	13
3. Eindimensionaler Parameterraum	16
3.1. Unabhängige Optimierung der Score-Anteile	16
3.2. Lineare Regression	19
3.3. Exponentialverteilung mit log-Link	21
3.4. Direkte Optimierung des Parameters	25
3.5. Optimierung des Parameters mit Strafterm	29
3.6. Heuristischer Algorithmus zur Suche des globalen Extremums	31
3.7. Simulationsstudien	34
3.7.1. Simulationsmodell SLA und SLB	35
3.7.2. Simulationsmodell SLC	37
3.7.3. Simulationsmodell SLD	41
3.7.4. Simulationsmodell SEA und SEB	42
3.7.5. Simulationsmodell SEC	44
3.7.6. Schlussfolgerung aus den Simulationsstudien	45
4. Mehrdimensionaler Parameterraum	54
4.1. Lineare Regression mit Intercept	55

4.2.	Generalisierte lineare Regression mit Exponentialverteilung	58
4.3.	Simulationsstudien	59
4.3.1.	Simulationsmodell MLA und MLB	60
4.3.2.	Simulationsmodell MLD	63
4.3.3.	Simulationsmodell MEA und MEB	64
4.3.4.	Schlussfolgerung aus den Simulationsstudien	65
4.4.	Einhüllende des zweidimensionalen Parameterraums	66
4.5.	Multiple Regression	70
5.	Anwendungsbeispiel	74
5.1.	Zusammenhang der Ausgaben für Forschung und Entwicklung mit dem Bruttoinlandsprodukt	74
5.2.	Regressionsmodell mit Intervalldaten	75
6.	Schluss	83
6.1.	Zusammenfassung	83
6.2.	Ausblick	86
	Literaturverzeichnis	89
A.	Übersicht verwendeter Verfahren	97
A.1.	Numerische Optimierungsverfahren	97
A.2.	Verfahren zur Schätzung der Parameterintervalle für die generalisierte lineare Regression mit Intervalldaten	97
B.	Weiteres Material für die Simulationsbeispiele	99
C.	Weiteres Material für das Anwendungsbeispiel	103
C.1.	Quelle und Generierung der Daten	103
C.2.	Tabelle der Daten	104
C.3.	Weitere Ergebnisse	107

Abbildungsverzeichnis

3.1. Anteile der Score-Funktion für $\beta = 0.5$ bei der linearen Regression ohne Intercept nach x_i und y_i	22
3.2. Score-Anteile bei der linearen Regression ohne Intercept für $\beta = 0.5$ und $\beta = -1$ mit den y -Werten $y = -1, 0, 1$	23
3.3. Drei Beispiele für Intervalldaten mit jeweils zwei Beobachtungen und den entsprechenden Regressionsgeraden ohne Intercept	24
3.4. Anteile der Score-Funktion für $\beta = 0.5$ bei der generalisierten linearen Regression mit Exponentialverteilung und log-Link ohne Intercept nach x_i und y_i	25
3.5. Score-Anteile für $\beta = 0.5$ und $\beta = -1$ bei der generalisierten linearen Regression mit Exponentialverteilung und log-Link	26
3.6. Score-Anteile für $\beta = 0.5$ und $\beta = -1$ bei der generalisierten linearen Regression mit Exponentialverteilung und log-Link	27
3.7. Drei Beispiele für Intervalldaten mit jeweils zwei Beobachtungen und den entsprechenden Regressionskurven ohne Intercept	28
3.8. Zwei Simulationsbeispiele (SLA, $n=20$) und (SLB, $n=20$) für Ober- und Untergrenzen einer linearen Regression ohne Intercept auf Intervalldaten	36
3.9. Simulationsbeispiel (SLC, $n = 20$) für Ober- und Untergrenzen einer linearen Regression ohne Intercept auf Intervalldaten	40
3.10. Beobachtungen mit Abweichungen größer als 0.01 für die Punkte in den Intervallen der einzelnen Beobachtungen für das Simulationsbeispiel (SLC, $n = 1000$)	48
3.11. Simulationsbeispiel (SLC, $n=10$) für Ober- und Untergrenzen einer linearen Regression ohne Intercept auf Intervalldaten	49
3.12. Simulationsbeispiel (SLD, $n=20$) für Ober- und Untergrenzen einer linearen Regression ohne Intercept auf Intervalldaten	50
3.13. Zwei Simulationsbeispiele (SEA, $n = 20$) und (SEB, $n = 20$) mit den Ober- und Untergrenzen einer generalisierten linearen Regression mit Exponentialverteilung ohne Intercept auf Intervalldaten	51

3.14. Unterschiede der Schätzungen für die Daten (SEA, $n = 20$) in der generalisierten linearen Regression mit Exponentialverteilung und log-Link ohne Intercept	52
3.15. Unterschiede der Schätzungen für die Daten (SEC, $n = 20$) in der generalisierten linearen Regression mit Exponentialverteilung und log-Link ohne Intercept	53
4.1. Schätzungen der Parameterintervalle für die Daten (MLB, $n = 20$) in der linearen Regression mit Intercept	64
4.2. Unterschiede für die Intervalldaten (MLA, $n = 20$) in der linearen Regression mit Intercept für β_0	65
4.3. Unterschiede für die Intervalldaten (MLA, $n = 20$) in der linearen Regression mit Intercept für β_1	66
4.4. Schätzungen der Parameterintervalle für die Daten (MLD, $n = 20$) in der linearen Regression mit Intercept	67
4.5. Schätzungen der Parameterintervalle für die Daten (MEA, $n = 20$) in der linearen Regression mit Intercept	72
4.6. Schätzungen der Parameterintervalle für die Daten (MEB, $n = 20$) in der linearen Regression mit Intercept	73
5.1. Streudiagramm des Bruttoinlandsprodukts und der Ausgaben für Forschung und Entwicklung für 82 Länder in PPP US\$ pro Einwohner im Jahr 2008	77
5.2. Die gleiche Situation wie in Abbildung 5.1, wobei die Ausgaben für Forschung und Entwicklung auf einer logarithmischen Skala angegeben sind	78
5.3. Grenzen des Identifizierungsbereichs von β_1 in einem linearen Modell für Intervalldaten des Bruttoinlandsprodukts und der Ausgaben für Forschung und Entwicklung	79
5.4. Die gleiche Situation wie in Abbildung 5.3 mit logarithmischer Skala der unabhängigen Variablen	80
5.5. Grenzen des Identifizierungsbereichs von β_1 in einem generalisierten linearen Modell für Intervalldaten des Bruttoinlandsprodukts und der Ausgaben für Forschung und Entwicklung	81
5.6. Die gleiche Situation wie in Abbildung 5.5 mit logarithmischer Skala der unabhängigen Variablen	82