

statistische Verfahren WS 2017/2018

## **Projekt 7 - Kriminalität**

Reda Ihtassine (155 685)      Ingo Schäfer (165 220)

Jena, am 27. März 2018

# Inhaltsverzeichnis

<b>1</b>	<b>Einleitung</b>	<b>1</b>
<b>2</b>	<b>Material und Methoden</b>	<b>2</b>
2.1	Material . . . . .	2
2.2	Methoden . . . . .	3
<b>3</b>	<b>Resultate</b>	<b>6</b>
3.1	Modellwahl . . . . .	6
3.2	Simulationsaufgabe . . . . .	8
<b>4</b>	<b>Diskussion</b>	<b>9</b>

# **Abbildungsverzeichnis**

# **Tabellenverzeichnis**

# 1 Einleitung

Statistiken sind ein wichtiges Mittel, um die Werte und Trends der Kriminalität zu schätzen, die Kosten für Auswirkungen auf die Gesellschaft zu bewerten und darüber die Strafverfolgungsansätze zu optimieren, um die Kriminalität im folgenden zu verhindern. Um ein ökonomisches Kriminalitätsmodell zu schätzen, können die Eigenschaften der Counties nicht ignoriert werden. Diesem Projekt liegen solche Daten des US-amerikanischen Bundesstaats North Carolina zugrunde, welche in dem Zeitraum von 1981 bis 1987 erhoben wurden. Sie wurden u.a. in der Arbeit von Baltagli<sup>1</sup> sowie von Cornwell und Trumbull<sup>2</sup> veröffentlicht.

Der übliche Hausman-Test, der auf dem Unterschied zwischen fixierten und zufälligen Effekten basiert, kann zu einer irreführenden Inferenz führen, wenn es endogene Regressoren des konventionellen simultanen Gleichungstyps gibt<sup>1</sup>. Daher ist es das Ziel dieser Projektarbeit ein geeignetes statistisches Modell für die Zahl der Verbrechen mithilfe von anderen Kriterien zu entwickeln. Dabei betrachten wir insbesondere die qualitative Einflussgröße *region* und deren mögliche Wechselwirkungen mit anderen Prädiktoren. Der zweite Teil dieser Arbeit beschäftigt sich mit der Untersuchung des Einflusses des Stichprobenumfangs auf die Genauigkeit der Approximation der tatsächlichen Kovarianzmatrix des Maximum-Likelihood-Schätzers durch die asymptotische Kovarianzmatrix.

Diese Arbeit gliedert sich in drei Kapitel: Im Kapitel Material und Methoden wird zunächst das Material aus der Datei *crimes.csv* und die verwendeten Methoden beschrieben. Im Kapitel Resultate werden die numerischen Ergebnisse vorgestellt. Im letzten Kapitel erfolgt die Diskussion und Interpretation der Ergebnisse hinsichtlich der Aufgabenstellung und der praktischen Anwendbarkeit der ausgewählten Modelle.

---

<sup>1</sup>Vgl.: Badi H. Baltagli, estimating an economic model of crime using panel data from North Carolina, journal of applied econometrics, S.: 543

<sup>2</sup>Vgl.: Cornwell C, Trumbull WN. 1994. Estimating the economic model of crime with panel data. Review of Economics and Statistics 76: 360 - 366.

## 2 Material und Methoden

### 2.1 Material

Cornwell und Trumbull (1994), nachstehend (CT), schätzten ein Wirtschaftsmodell der Kriminalität unter Verwendung von Paneldaten über 90 Countys in North Carolina zwischen 1981 und 1987. Becker (1963) und Ehrlich (1973) unter anderem folgen das empirische Modell, der die Kriminalitätsrate (FBI-Index, der die Anzahl der Verbrechen geteilt durch die Kreisbevölkerung misst) auf eine Reihe von erklärenden Variablen bezieht, diese Variablen einschließlich abschreckende Variablen sowie Variablen, die die Chancen zum Legal zurückzukehren berechnen. Alle Variablen, außer der Regional- und Zeitdummy Variablen, sind in Protokollen gespeichert <sup>3</sup>. Der Datensatz besteht aus einer .csv-Datei. In ihr sind die unterschiedlichen 90 Counties von North Carolina zeilenweise aufgelistet. Die Spalten sind (mögliche) Eigenschaftsvektoren. In der Arbeit von Baltagi <sup>4</sup> werden noch einige Eigenschaften mehr aufgelistet, als in dieser Arbeit betrachtet wurden. Daher hier eine kleine Übersicht über alle möglichen Einflussgrößen: Alle Eigenschaftsvektoren sind logarithmisch mit Ausnahme der Region und der Zeit. Die erste Spalte beinhaltet die Zielgröße *crimes*, also die Anzahl aller Straftaten in dem jeweiligen County über den Zeitraum von 1981-1987.

Weiterhin wurde die Arrestwahrscheinlichkeit  $P_A$  hinzugefügt. Sie berechnet sich aus  $P_A = \frac{\text{Arrestierungen}}{\text{textDelikte}}$ . Sie wird abgekürzt *prbarr* geschrieben. Daneben gibt es auch die Überzeugungswahrscheinlichkeit  $P_C$ . Sie gibt das Verhältnis zwischen tatsächlichen Arrestierungen und den gestandenen Straftaten an und wird daher berechnet mit  $P_P = \frac{\text{Anzahl tatsächlicher Arrestierungen}}{\text{Anzahl gestandener Straftaten}}$ . Sie wird bezeichnet als *prbpris*.

Eine weitere Eigenschaft ist die Fähigkeit des Countys ein Verbrechen auch zu ermitteln. In dem Datensatz spiegelt sich dies in der Variable *polpc* wieder. Sie gibt das Polizei-pro-Kopf-Verhältnis an.

Ein weiteres wichtiges Merkmal ist die Bevölkerungsdichte (*density*). Sie stellt das Verhältnis *fraca*anzahl bevölkerung/Fläche des Countys in square miles dar.

Darüber hinaus wird das Verhältnis von Minderheiten zu der Gesamtanzahl Einwohner in der Variable *pctmin* ausgedrückt.

*pctymale* ist eine Eigenschaft, die den Anteil der jungen männlichen Bevölkerung zur Gesamtbevölkerung anzeigt.

Die letzten fünf Variablen geben den durchschnittlichen Bruttolohn in den Bereichen Baugewerbe (*wcon*, Staatsangestellte (*wsta*), Dienstleistungssektor (*user*), Handel (*wtrd*) und Bankgeschäften (*textitwfr*) wieder.

<sup>3</sup>Vgl.: Badi H. Baltagi, estimating an economic model of crime using panel data from North Carolina, journal of applied econometrics, S.: 543 f.

<sup>4</sup>1

## 2.2 Methoden

Um ein geeignetes Modell aus den oben beschriebenen Merkmalen zu finden, wurden fünf unterschiedliche Herangehensweisen vorgeschlagen, um ein Modell zu finden, das möglichst geringe Fehler aufweist.

- explorative Herangehensweise (ausprobieren)
- Vergleich aller Modelle mit nur einem Merkmal
- Verwendung von `step()` und anschließende Minimierung des Modells
- strukturierte Suche nach einem geeigneten Modell
- Verwendung von `cor()`

Am Ende einer jeden Herangehensweise wurde ein bestes Modell vorgeschlagen. Diese wurden dann anschließend miteinander verglichen, um ein bestmögliches Modell zu bestimmen.

Hauptsächlich wurden zwei Gütekriterien verwendet.

Zum einen *Akaike's Information Criterion* (AIC), welches die logarithmische Fehlerabweichung des Schätzers mit der Anzahl der verwendeten Merkmale bestraft.

$$\text{AIC} := -2 * \ln(\hat{\Theta}_n) + 2p \quad (1)$$

AIC spiegelt den Kompromiss zwischen Verbesserung der Modellanpassung durch erhöhte  $p$  und erhöhte Ungenauigkeit durch Schätzung vieler Parameter wieder.

In einigen Fällen wurde auch die *Devienz* betrachtet, um die Güte mehrerer Modelle miteinander zu vergleichen. Hier geht man von einem saturierten Modell aus. Dies ist das komplexeste Modell für einen Datensatz, dass durch Erhöhung der Parameterzahl erzeugt werden kann. In vielen Fällen hat das saturierte Modell daher so viele Parameter wie Beobachtungen. Falls Einflussvektoren mehrfach vorkommen, besitzt das saturierte Modell weniger Parameter. Das ist typischerweise der Fall für Experimente mit qualitativen Einflussgrößen.

Hier wird die Likelihood-Quotienten-Statistik zum Vergleich eines Modells  $M$  mit dem saturierten Modell

$$T(\underline{Y}) = 2(l_{\text{saturiert}} - l_M) \quad (2)$$

betrachtet.

Die Likelihood-Quotienten-Statistik ist asymptotisch  $\chi^2$  - verteilt. Dabei ist  $r$  die Differenz der Parameterzahlen. Deswegen funktioniert hier der Likelihood-Quotienten-Test nicht, da für  $n \rightarrow \infty$  die Anzahl der Freiheitsgrade auch typischerweise unbeschränkt wächst.

Die Größe

$$D(M) = 2(l_{\text{saturiert}} - l_M) \quad (3)$$

heißt Devienz des Modells  $M$ .

Dabei ist zu beachten, dass ein Modell  $M$  ein geeignetes Modell ist, falls die Devienz

von  $M$  ungefähr so groß ist wie die ungefähre Anzahl Parameter von  $M$ .

$$D(M) \approx n - |M| \quad (4)$$

Als anderes Gütekriterium wurde das Quadrat der erwarteten Fehlerabweichungen (*SPSE*) im Kreuzvalidierungsverfahren berechnet.

Dazu wurde der gesamte ausgewählte Datensatz in einen Trainings- und einen Testdatensatz aufgeteilt. Das "beste Modell" ist dasjenige, dass im Mittel den kleinsten geschätzten erwarteten Prognosefehler liefert. Dabei wird typischerweise eine  $l$ -fache Kreuzvalidierung durchgeführt:

Es gibt einen Testdatensatz  $I = 1 \dots n$ . Dieser wird in  $l$  etwa gleichgroße Indexmengen  $I_1, \dots, I_l$  zerlegt.

In jedem  $j$ -ten Schritt wird ein  $I_j$  als Testdatensatz gewählt. Alle anderen Indexmengen bilden den Trainingsdatensatz.

Nun wird der erwartete Prognosefehler geschätzt:

$$\sum_{i \in I_j} (y_i - \underline{x}_i^{(M)T} \underline{\hat{\beta}}^{(m-j)})^2 = SPSE_j^{(M)} \quad (5)$$

Dabei ist  $\underline{\hat{\beta}}^{(m-j)}$  die auf  $I/I_j$  basierende Schätzung.

Zuletzt werden alle Teilschätzungen zu einer Schätzung für SPSE zusammen kombiniert:

$$SPSE^{(M)} := \sum_{j=1}^l (SPSE_j^{(M)}) \quad (6)$$

Zu bemerken ist, dass jede Beobachtung einmal in einem Testdatensatz verwendet wird. Außerdem ist die Abhängigkeit von der konkreten Zerlegung nur reduziert, aber nicht verschwunden. Es gibt einen Spezialfall, wenn  $l = n$ . Das heißt, dass der gesamte Testdatensatz in  $n$  Teildatensätze zerlegt wird. Jede Beobachtung wird mit der Prognose basierend auf  $(n - 1)$  Beobachtungen verglichen. Dies ist auch bekannt als *leave-one-out-cross-validation*. Als Faustregel empfiehlt es sich  $l \approx 10$  zu wählen.

In der Simulationsaufgabe des Projektes sollte der Einfluss des Stichprobenumfangs auf die Genauigkeit der Approximation der tatsächlichen Kovarianzmatrix des Maximum-Likelihood-Schätzers durch die asymptotische Kovarianzmatrix untersucht werden. Dazu wurde zunächst ein möglichst einfaches wahres Modell angenommen. Anhand dessen wurde aus dem gesamten Datensatz eine beliebig große Teilmenge  $T$  entnommen. Aus den Daten von  $T$  wurde dann eine Designmatrix gebildet, die Grundlage für die darauf folgenden Generierungen fÄ<sub>4</sub>r Pseudozufallszahlen war. Meistens wurde die Größe dieser zu untersuchenden Teilmenge auf 30 gesetzt. (Der gesamte Datensatz umfasst 90 Subjekte.) Jedoch wurden auch andere Größen überprüft. Aus der auf diese Art und Weise gebildeten Designmatrix, wurden nun wiederum unterschiedlich große Stichproben ausgewählt und die daraus berechneten  $\beta_0$  und  $\beta_1$  Werte in einer weiteren Matrix



gespeichert. Aus dieser Matrix wurden dann die Varianz und die Kovarianz für die tatsächliche Kovarianzmatrix berechnet.

Da

$$F^{\frac{T}{2}}(\hat{\underline{\beta}})(\hat{\underline{\beta}}_n - \underline{\beta}) \xrightarrow[n \rightarrow \infty]{d} N(0, I) \quad (7)$$

gilt, gilt für die Approximation von  $\hat{\underline{\beta}}_n$  bei festem  $n$ :

$$\hat{\underline{\beta}}_n \approx N(\underline{\beta}, I^{-1}(\underline{\beta})) \quad (8)$$

Die Kovarianzmatrix  $\mathbb{X}$  ist die inverse Fisher-Matrix  $I$ . Daher hat  $I(\underline{\beta})$  die kanonische Linkfunktion einfachen Gestalts

$$I(\underline{\beta}) = \mathbb{X}^T V \mathbb{X} \quad (9)$$

Dabei ist  $V$  eine Diagonalmatrix, welche in der Spur die Varianzen hält. Anhand dessen wurde die asymptotische Kovarianzmatrix berechnet. Die daraus herausgehenden Resultate wurden dann bei einem kleiner werdenden  $n$  auch immer geringer, sodass die Ergebnisse immer in Relation zueinander verglichen wurden.

## 3 Resultate

### 3.1 Modellwahl

Wie bereits erwähnt, wurden fünf unterschiedliche Herangehensweisen betrachtet, um ein geeignetes Modell zu finden.

**Wahl der Verteilung** negative binomialverteilung statt gauß-verteilung, begründung - siehe quelle!

**Die besondere Rolle von der Einflussgröße region**

**Herangehensweisen**

**explorative Herangehensweise** Um ein gutes Gefühl für die Merkmalsvektoren zu bekommen, wurden zunächst einige Modelle ausprobiert und mittels AIC verglichen. Damit ein Vergleichswert nach dem Akaike-Maß vorhanden war, wurde ein komplettes Modell angenommen, das aus allen vorhandenen Merkmalen besteht. Dieses Modell heißt *mAll*. Die entsprechende Formel sieht so aus:

$$crimes = prbarr + prbpris + polpc + density + area + taxpc + region + pctmin + pctymale + wcon + wsta + wstc \quad (10)$$

Es wurde bewusst darauf verzichtet in diesem 'gesamten' Modell die Intersections (Wechselwirkungen) der einzelnen Merkmale zu betrachten. Grund dafür ist, dass das Akaike-Maß Modelle mit vielen Einflussgrößen mehr bestraft, als solche die weniger besitzen. Da bei dieser Untersuchung das Akaike-Maß das am häufigsten verwendete Kriterium war, sollte also das erste Modell, mit dem die anderen verglichen wurden, nicht einen großen negativen Wert aufweisen, so wie das in diesem Fall der Fall gewesen wäre. (Der Akaike-Wert des Modells, das alle Merkmale und alle Wechselwirkungen zwischen diesen betrachtet, beträgt -3441.465. In diesem Modell gibt es 91 Freiheitsgrade.) Die Daten *crimes.data*, welche der Funktion `glm.nb(formula, data = crimes.data)` während der gesamten Untersuchung gegeben wurden, wurden nicht verändert. Es handelt sich hierbei immer um den gesamten Datensatz aus der Datei *crimes.csv*.

Im Folgenden wurde bemerkt, dass diese Merkmale durchaus gruppiert betrachtet werden können. Daher bestand die erste Idee darin, die unterschiedlichen Gruppierungen je Modell zu betrachten: Die ersten beiden Merkmale (*prbarr* und *prbpris*) geben beide Verhältnisse zum Anteil aller Straftäter in einem County an. Daher wurde ein Modell aus diesen beiden Einflussgrößen betrachtet.

$$crimes = prbarr : prbpris \quad (11)$$

Die Einflussgrößen *density* und *area* sind beides räumliche Merkmale. Auch sie wurden in einem Modell zusammengefasst. Wie in [3.1](#)

**Vergleich aller Modelle mit jeweils nur einem Merkmal**

**Verwendung von `step()` und anschließende Minimierung des Modells**

**strukturierte Suche nach einem geeigneten Modell**

**Verwendung von `cor()`**

**Die Gewinnermodelle**

## 3.2 Simulationsaufgabe

Beschreibung simulation()

Auswertung der Ergebnisse

einfaches Modell: mDensity

**Ergebnisse mit Gewinnermodell aus der ersten Aufgabe** Hier leite ich zur Diskussion über.

## 4 Diskussion

- 'sieger'modelle sind recht gut, aber immer noch sehr große abweichungen zu den tatsächlichen werten. -

## **Literatur**