

statistische Verfahren WS 2017/2018

## **Projekt 7 - Kriminalität**

Reda Ihtassine (Matrikelnummer)      Ingo Schäfer (165 220)

Jena, am 26. 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 . . . . .	2
<b>3</b>	<b>Resultate</b>	<b>6</b>
<b>4</b>	<b>Diskussion</b>	<b>7</b>

# **Abbildungsverzeichnis**

# **Tabellenverzeichnis**

# 1 Einleitung

Dem Projekt liegen Kriminalitätsdaten des US-amerikanischen Bundesstaats North Carolina zugrunde, welche in dem Zeitraum von 1981 bis 1987 erhoben wurden. Diese Daten wurden schon mehrfach mit Hilfe von verschiedenen Methoden (Hausmann, 2SLS, ...) von Anderen untersucht <sup>1</sup>.

Diese Arbeit versucht ein geeignetes allgemeines lineares Modell zu erarbeiten, mit dem gute Abschätzungen erzielt werden können.

---

<sup>1</sup>1

## 2 Material und Methoden

### 2.1 Material

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 Baltagli <sup>2</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 *fraczahl* 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 (*wser*), Handel (*wtrd*) und Bankgeschäften (*textitwfr*) wieder.

### 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.

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

---

<sup>2</sup><sub>1</sub>

- 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 Parameteranzahl 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 $\frac{1}{4}$ 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

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

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

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

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

1:

Vorgehensweise:

- negative binomialverteilung statt gauß-verteilung, begründen - siehe quelle! - 5 unterschiedliche Herangehensweisen um ein geeignetes Modell zu finden, alle kurz erklären
- besondere Rolle von *region*
- vergleich der modelle funktionsweise knapp erläutern(aic, cross\_validation, cor())
- vorstellen 5 gewinnermodelle, den gewinner

2 :

- beschreibungstest() - funktion(ggf funktionen eingriffigererennamengeben)
- welche einstellungen erzielt gute ergebnisse?
- einmal mit einfachem modell zeigen : mDensity...(< -warum mDensity?) - einmal mit gewinnermodell - - > wie gut ist ergebnis? 1/4 berleitung zu diskussion...

## 4 Diskussion