

ALBERT-LUDWIGS-UNIVERSITÄT FREIBURG

BACHELORARBEIT

Verzerrung der Inferenz bei Verwendung gemischter Modelle in latenten Repräsentationen

Autor:

Yannick Bantel

Professor:

Prof. Dr. Harald Binder

Betreuer:

Clemens Schächter

Abgabedatum:

18. Juni 2024



universität freiburg

Inhaltsverzeichnis

1 Einleitung

In der modernen Datenanalyse spielen gemischte Modelle eine zentrale Rolle, da sie es ermöglichen, sowohl feste als auch zufällige Effekte zu berücksichtigen. Dies macht sie besonders in den Bereichen der Biostatistik, der Sozialwissenschaften und der ökonomischen Modellierung populär.

Im Rahmen dieser Arbeit werden gemischte Modelle auf die Analyse medizinischer Daten angewendet. Mit dem Aufkommen von Big Data und komplexen Datenstrukturen hat sich der Fokus zunehmend auf die effiziente und genaue Extraktion von Informationen aus großen und oft unübersichtlichen Datensätzen verlagert.

In diesem Zusammenhang gewinnen latente Repräsentationen an Bedeutung, da sie es ermöglichen, inhärente Strukturen innerhalb der Daten zu identifizieren und zu nutzen, um tiefere Einblicke zu gewinnen.

Die Integration von gemischten Modellen in latente Repräsentationen birgt jedoch das Risiko einer Verzerrung der Inferenzergebnisse, was die Genauigkeit und Zuverlässigkeit der aus den Daten gezogenen Schlussfolgerungen erheblich beeinträchtigen kann.

Die vorliegende Arbeit widmet sich der Untersuchung von Verzerrungen, die bei der Anwendung gemischter Modelle auf latente Repräsentationen auftreten können. Das Ziel dieser Arbeit ist es, die Mechanismen zu verstehen, die zu diesen Verzerrungen führen, sowie Methoden zu entwickeln, um ihre Auswirkungen zu minimieren.

Das Problem der Verzerrung ist von besonderer Relevanz, da eine fehlerhafte Inferenz zu Fehlentscheidungen führen kann, die in praktischen Anwendungen schwerwiegende Konsequenzen haben können.

Die Arbeit zielt darauf ab, durch eine sorgfältige Analyse und Bewertung von gemischten Modellansätzen in Verbindung mit latenten Repräsentationen einen Beitrag zur Verbesserung der Modellgenauigkeit und der Zuverlässigkeit von Inferenzschlüssen zu leisten.

Die Arbeit ist in mehrere Teile gegliedert, die zunächst die theoretischen Grundlagen von gemischten Modellen und latenten Repräsentationen behandeln. Im Anschluss erfolgt eine Diskussion der Methoden zur Messung und Korrektur von Verzerrungen.

Im Anschluss werden die zuvor theoretisch erörterten Konzepte anhand von empirischen Studien praktisch angewendet und evaluiert. Auf Basis der gewonnenen Erkenntnisse werden abschließend Empfehlungen für die Anwendung dieser Techniken in Forschung und Praxis gegeben.

1.0.1 Motivation

2 Theoretische Grundlagen

Im Vorfeld der Erörterung der Methodik dieser Arbeit ist eine theoretische Aufarbeitung der behandelten Themen unabdingbar. In diesem Kapitel erfolgt eine ausführliche Beschreibung und Behandlung der theoretischen Aspekte dieser Arbeit. Es werden sowohl lineare gemischte Modelle als auch die Theorie hinter Variational Autoencodern eingeführt und beschrieben. Im Folgenden wird insbesondere auf die für die Analyse der Modelle notwendige Theorie eingegangen, wie beispielsweise die Likelihood-Berechnung und der Likelihood-Ratio-Test.

2.1 Variational Autoencoder (VAE)

Die Anwendung der gemischten Modelle auf einer latenten Repräsentation erfolgt mittels Variational Auto-Encoder (VAE). Variational Auto-Encoder sind generative Modelle, welche versuchen die zugrunde liegende Struktur der Inputdaten x im latenten Raum zu modellieren. Im Vergleich zu normalen Autoencodern, welche den latenten Raum durch feste Punkte modellieren, wie es in Abbildung ?? dargestellt ist, wird der latente Raum in der Erweiterung VAE durch eine Wahrscheinlichkeitsverteilung (Normalverteilung) modelliert.

2.1.1 Struktur des VAEs

Die Architektur eines VAE basiert auf zwei neuronalen Netzwerken: einem Encoder, der die Inputdaten x im latenten Raum als Verteilung kodiert, und einem Decoder, der aus den latenten Daten versucht die Originaldaten zu rekonstruieren. Diese Module lernen die wesentlichen Merkmale der Eingabedaten zu extrahieren und eine komprimierte Version dieser Daten zu erzeugen.

VAEs sind für die Modellierung latenter Repräsentationen von großem Interesse, da sie hochdimensionale Datensätze mit Hilfe des Encoders im latenten Raum niedrigdimensional darstellen können. Dies reduziert die Komplexität der Modellierung und ermöglicht es, gemischte Modelle effizienter und genauer zu betreiben. Im Gegensatz zu herkömmlichen Autoencodern ist der VAE in der Lage, nicht nur den Eingabedatensatz zu rekonstruieren, sondern auch neue Inhalte zu generieren. Dies wird durch die verbesserte Repräsentation ermöglicht(vgl. [bigdata-insider-vae]).

Latenter Raum

Variablen, die man nicht direkt messen kann, demnach nicht Teil des erhaltenen Datensatzes sind, bezeichnet man als latente Variablen. Sie werden erst mithilfe der gegebenen Daten erschlossen und ergeben im Verbund den latenten Raum.

Im VAE werden die latenten Variablen z aus der priori-Verteilung $p_\theta(x)$ gezogen, welche eine multivariate Normalverteilung $p_\theta(x) = \mathcal{N}(z; 0, I)$ ist. Die latenten Daten werden aus den Inputdaten durch den Encoder gezogen, welcher die posteriori Verteilung durch eine variable Verteilung $q_\phi(z, x)$ approximiert. Der Encoder erlernt somit zwei Vektoren, nämlich den Mittelwert μ und die Standardabweichung σ^2 der Normalverteilung $q_\phi(z, x) = \mathcal{N}(z; \mu(x), \sigma^2(x))$.

Der Decoder versucht aus den latenten Variablen die Inputdaten x durch die likelihood-Verteilung $p_\theta(x|z)$ zu rekonstruieren. Die Wahrscheinlichkeit, dass die beobachteten Daten aus den latenten Repräsentationen generiert wurden, wird durch dieses Modell modelliert. Auch hier wird typischerweise eine Normalverteilung angenommen, sofern die Daten reellwertig sind. Im Falle binärer Daten wird die Verteilung als Bernoulli-Verteilung modelliert.

Für weiterführende Details wird auf die Publikation [Auto-Encoding Variational Bayes] verwiesen.

2.1.2 Training VAE

Ein VAE stellt ein spezifisches Beispiel einer Variational Inference (VI) dar. Die Zielsetzung einer VI besteht in der Berechnung der Posteriori-Verteilung $p(z|x)$ der latenten Variablen in Abhängigkeit von den Input-Daten. Die Berechnung dieser Verteilung ist besonders bei komplexen Modellen mit Schwierigkeiten verbunden. Infolgedessen approximiert der VAE die Posteriori-Verteilung durch eine einfachere Verteilung $q_\phi(x)$. Das Ziel von VIs besteht folglich in der Minimierung des Rekonstruktionsfehlers sowie der Kullback-Leibler-Divergenz zwischen der tatsächlichen Posteriori-Verteilung $p(z|x)$ und der approximierten Verteilung $q_\phi(x)$.

Definition 2.1.1 (Kullback-Leibler-Divergenz (KL-Divergenz)).

Sei $q_\phi(z|x)$ die approximierte Posteriori-Gauß-Verteilung und $p(z)$ die Priori-Gauß-Verteilung. Dann ist die KL-Divergenz definiert als

$$D_{KL}(q_\phi(z|x)||p(z)) = \int q_\phi(z|x) \log \left(\frac{q_\phi(z|x)}{p(z)} \right) dz$$

Da die Berechnung von $p(z)$ mit Schwierigkeiten verbunden ist, wird stattdessen die Evidence Lower Bound (ELBO) durch Stochastic Gradient Descent (SGD) oder andere ähnliche Verfahren optimiert.

Das Objekt der Optimierung im VAE ist somit der ELBO.

Im Folgenden wird die log-Likelihood der Daten betrachtet, um den ELBO herzuleiten (vgl. [Introduction to VAEs]).

$$\log p_\theta(x) = \log p_\theta(x) * \overbrace{\int q_\phi(z|x) dz}^{=1} \quad (2.1)$$

$$= \int \log p_\theta(x) q_\phi(z|x) dz \quad (2.2)$$

$$= \mathbb{E}_{q_\phi(z|x)} [\log p_\theta(x)] \quad (2.3)$$

$$= \mathbb{E}_{q_\phi(z|x)} \left[\log \left[\frac{p_\theta(x, z)}{p_\theta(z|x)} \right] \right] \quad (2.4)$$

$$= \mathbb{E}_{q_\phi(z|x)} \left[\log \left[\frac{p_\theta(x, z) q_\phi(z|x)}{q_\phi(z|x) p_\theta(z|x)} \right] \right] \quad (2.5)$$

$$= \underbrace{\mathbb{E}_{q_\phi(z|x)} \left[\log \left[\frac{p_\theta(x, z)}{q_\phi(z|x)} \right] \right]}_{= \mathcal{L}_{\theta, \phi}(x) \text{ (ELBO)}} + \underbrace{\mathbb{E}_{q_\phi(z|x)} \left[\log \left[\frac{q_\phi(z|x)}{p_\theta(z|x)} \right] \right]}_{= D_{KL}(q_\phi(z|x)||p_\theta(z|x))} \quad (2.6)$$

Der zweite Term in der Gleichung ist die nicht negative Kullback-Leibler-Divergenz (KL-Divergenz) zwischen $q_\phi(z|x)$ und $p_\theta(z|x)$. Der erste Term in Gleichung ?? stellt den Evidence-Lower-Bound, kurz ELBO, dar:

Definition 2.1.2 (Evidence Lower Bound (ELBO) für VAEs).

Sei $q_\phi(z|x)$ das Encoder Modell und $p_\theta(x, z)$ das Decoder Modell. Der ELBO ist definiert durch

$$\mathcal{L}_{\theta, \phi}(x) = \mathbb{E}_{q_\phi(z|x)} [\log p_\theta(x, z) - \log q_\phi(z|x)]$$

Die Umstellung der Gleichung ?? zeigt, dass der ELBO eine untere Schranke für die log-likelihood der Daten darstellt, da die KL-Divergenz nicht negativ ist:

$$\mathcal{L}_{\theta, \phi}(x) = \log p_\theta(x) - D_{KL}(q_\phi(z|x)||p_\theta(z|x)) \quad (2.7)$$

$$\leq \log p_\theta(x) \quad (2.8)$$

Die Gleichung kann alternativ mit der Jensenschen Ungleichung (vgl. **[JensenscheUngleichung]**) wie folgt umgestellt werden:

$$\log p_\theta(x) = \log \int p_\theta(x, z) dz \quad (2.9)$$

$$= \log \int p_\theta(x, z) \frac{q_\phi(z|x)}{q_\phi(z|x)} dz \quad (2.10)$$

$$= \log \mathbb{E}_{q_\phi(z|x)} \left[\frac{p_\theta(x, z)}{q_\phi(z|x)} \right] \quad (2.11)$$

$$\stackrel{\text{Jensen}}{\geq} \mathbb{E}_{q_\phi(z|x)} \left[\log \frac{p_\theta(x, z)}{q_\phi(z|x)} \right] = \mathcal{L}_{\theta, \phi}(x) \quad (2.12)$$

Es ist sofort ersichtlich, dass die log-likelihood $p_\theta(x)$ durch Maximieren der ELBO bzgl. θ und ϕ selbst maximiert wird. Folglich verbessert sich die Qualität unseres generatives Modells. Gleichzeitig wird dadurch die KL-Divergenz der Approximation $q_\phi(z|x)$ an den wahren Posteriori $p_\theta(z|x)$ minimiert. Daher wird die Approximation $q_\phi(z|x)$ optimiert. Die ELBO kann durch stochastische Gradientenverfahren optimiert werden.

Die Berechnung des Gradienten von $\mathcal{L}_{\theta, \phi}(x) = \mathbb{E}_{q_\phi(z|x)} [\log p_\theta(x, z) - \log q_\phi(z|x)]$ bezüglich θ kann problematisch sein. Der Monte-Carlo-Schätzer für Gradienten ist eine gängige Methode und kann wie folgt definiert werden (vgl. **[MonteCarloEstimation]**):

$$\nabla_\phi \mathbb{E}_{q_\phi(z)} [f(z)] = \mathbb{E}_{q_\phi(z)} [f(z) \nabla_{q_\phi(z)} \log q_\phi(z)] \approx \frac{1}{L} \sum_{l=1}^L f(z) \nabla_{q_\phi(z^{(l)})} \log q_\phi(z^{(l)}) \quad (2.13)$$

wobei $z^{(l)} \sim q_\phi(z|x^{(i)})$ ist. In der Regel stellt dies keine Herausforderung dar. Allerdings ist $q_\phi(z|x)$ nun auch abhängig von ϕ , sodass der Monte Carlo Gradienten-Schätzer nicht direkt berechnet werden kann. Zur Lösung dieses Problems wird der sogenannte Reparameterization-Trick eingesetzt, welcher die Zufallsvariable transformiert um die Gradienten-Berechnung zu vereinfachen.

Reparametrisierungs Trick

Der Reparameterisierungs-Trick ist eine Methode zur Vereinfachung der Gradientenberechnung in Variational-Autoencodern. Er ermöglicht eine effizientere Berechnung der Gradienten der Evidence Lower Bound und somit eine effizientere Optimierung der ELBO. Der Reparameterization Trick transformiert die Zufallsvariable z in eine deterministische Funktion von einer Hilfsvariablen ϵ , was es erleichtert Backpropagation durchzuführen. Sei also die latente Variable z , die aus $q_\phi(z|x)$ gezogen wurde, gegeben. Die latente Variable z wird nun als deterministische Funktion einer Hilfsvariablen ϵ , unabhängig von x und ϕ , ausgedrückt. Die Transformation sieht dann wie folgt aus

$$z = g_\phi(\epsilon, x)$$

$g_\phi(\epsilon, x)$ ist dabei eine differenzierbare Funktion und ϵ eine Zufallsvariable mit einer bekannten Verteilung (z.B. $\epsilon \sim \mathcal{N}(0, I)$).

Im Falle einer Gaußverteilung $z \sim \mathcal{N}(\mu, \sigma^2)$ könnte die Umparametrisierung wie folgt aussehen

$$z = \mu + \sigma \odot \epsilon \quad \text{mit } \epsilon \sim \mathcal{N}(0, I).$$

Dies ist in Abbildung ?? veranschaulicht. Dadurch können die Gradienten bezüglich μ und σ effizient berechnet werden, da der Erwartungswert über $q_\phi(z|x)$ sich nun als Erwartungswert über ϵ schreiben

lässt.

Wir können den Trick direkt auf Gleichung ?? anwenden (vgl. [MonteCarloEstimation]). Es gilt

$$\nabla_{\phi} \mathbb{E}_{q_{\phi}(z)} [f(z)] = \nabla_{\phi} \int q_{\phi}(z) f(z) dz \quad (2.14)$$

$$= \nabla_{\phi} \int q_{\epsilon}(z) f(g_{\phi}(\epsilon; \phi)) d\epsilon \quad (2.15)$$

$$= \nabla_{\phi} \mathbb{E}_{q(\epsilon)} [g_{\phi}(\epsilon, z)] \quad (2.16)$$

2.2 Gemischte Modelle

Ein gemischtes Modell stellt ein statistisches Verfahren zur Datenanalyse dar, welches sowohl feste als auch zufällige Effekte (fixed and random effects) modelliert. Gemischte Modelle finden insbesondere bei der Analyse von longitudinalen und cluster-spezifischen Daten Anwendung, welche aus zeitlich wiederholten Beobachtungen (y_{it}, x_{it}) , $t = 1, \dots, T_i$ für jedes Individuum $i = 1, \dots, n$ bestehen. Die Variable y kennzeichnet dabei eine Antwortvariable, während x ein Vektor von Kovariablen darstellt. Ein Beispiel für einen solchen Datensatz ist ein medizinischer Datensatz,

$$(y_i, x_i) = (y_{i1}, \dots, y_{iT_i}, x_{i1}, \dots, x_{iT_i})$$

bei dem y_{ij} eine Beobachtung an Individuum i zum Zeitpunkt t_{ij} bezeichnet und T_i ist die Anzahl an Beobachtungen.

Zur Einführung der gemischten Modelle folgen wir den Notationen in [fahrmeir-2001-multivariate] und [fahrmeir-2011-regression]. Longitudinal und cluster-spezifische Daten weisen zwei Ebenen auf. Im Folgenden betrachten wir das Beispiel des medizinischen Datensatzes. Die erste Ebene bezieht sich dabei auf die Daten innerhalb einer Gruppe oder eines Individuums. In diesem Fall umfasst die erste Ebene den Patienten als Individuum mit seinen unterschiedlichen Werten für die Tests entlang der Zeitreihe T_i . Auf der allgemeineren zweiten Ebene erfolgt eine Betrachtung aller Patienten.

Im Rahmen eines gemischten Modells wird auf der ersten Ebene angenommen, dass die Antwortvariablen linear von den unbekannten bevölkerungsspezifischen festen Effekten β und den unbekannten cluster-spezifischen zufälligen Effekten b_i abhängen.

Die folgende Gleichung beschreibt das Modell:

$$y_{it} = x_{it}^t \beta + z_{it}^t b_i + \epsilon_{it} \quad (2.17)$$

Innerhalb des Modells werden die Designvektoren z_{it} und w_{it} als unabhängige Variablen definiert, wobei z_{it} beispielsweise die Testwerte in einem medizinischen Datensatz repräsentiert. Die Zufallsvariable ϵ_{it} hingegen ist unkorreliert und folgt einer normalverteilten Wahrscheinlichkeitsdichte mit Erwartungswert $\mathbb{E}(\epsilon_{it}) = 0$ und Varianz $\text{Var}(\epsilon_{it}) = \sigma^2$. Der Ausdruck a^t bezeichnet den transponierten Vektor, bzw. die transponierte Matrix a .

Betrachtet man nun die zweite Ebene, so werden die zufälligen Effekte b_i zwischen den verschiedenen Individuen gemäß einer Mischverteilung mit $\mathbb{E}(b_i) = 0$ unabhängig variieren. Es wird angenommen, dass die zufälligen Effekte b_i unabhängig und identisch normalverteilt sind,

$$b_i \sim \mathcal{N}(0, Q) \quad (2.18)$$

mit der $(q \times q)$ Kovarianzmatrix $\text{Cov}(b_i) = Q > 0$, welche symmetrisch und positiv semi-definit ist. Eine ausführliche Beschreibung findet sich in [pinheiro2000] (Kapitel 2.2.1).

Aufgrund dieser Überlegungen lässt sich nun das Model ?? in einer allgemeineren Form beschreiben:

Definition 2.2.1 (Lineares gemischtes Modell für Longitudinal- oder Clusterdaten).

Seien $X_i = (x_{i1}, \dots, x_{iT_i})$ und $Z_i = (z_{i1}, \dots, z_{iT_i})$ bekannte Designmatrizen für die festen und zufälligen Effekte. Seien β ein p -dimensionaler Vektor von festen Effekten und b_i ein q -dimensionaler Vektor von zufälligen Effekten und sei $\epsilon_i = (\epsilon_{i1}, \dots, \epsilon_{iT_i})$ der normalverteilte Fehlervektor.

Ein lineares gemischtes Modell für den T_i -dimensionalen Antwortvektor der i -ten Gruppe wird durch

$$y_i = X_i * \beta + Z_i * b_i + \epsilon_i$$

$$b_i \sim \mathcal{N}(0, Q), \epsilon_i \sim \mathcal{N}(0, R = \sigma_\epsilon^2 I)$$

definiert.

Die Daten der zufälligen und festen Effekte werden in einer Designmatrix (Datenmatrix) gespeichert. Die Parametervektoren β (für die festen Effekte) und b_i (für die zufälligen Effekte) initialisieren den Einfluss der Daten auf den Antwortvektor. Um auch für immer auftretende Messfehler oder unerwartete Einflüsse gewappnet zu sein, wird ein zufälliges Rauschen ϵ hinzugefügt.

Aufgrund des normalverteilten Fehlervektors können nun auch ein marginales Modell als multivariates heteroskedastisches lineares Regressionsmodell definiert werden. Dieses Modell ist für die Berechnung der Likelihood-Inferenz von entscheidender Bedeutung.

Definition 2.2.2 (Marginale gemischtes Modell).

Seien die Annahmen von ?? gegeben. Das marginale gemischte Modell ist definiert als

$$y_i = X_i \beta + \epsilon_i^*,$$

mit dem multivariaten Fehlervektor $\epsilon_i^* = (\epsilon_{i1}^*, \dots, \epsilon_{iT_i}^*)$ mit $\epsilon_{it}^* = z_{it}^T b_i + \epsilon_i$. Die ϵ_{it}^* sind dabei unabhängig und identisch verteilt (i.i.d.),

$$\epsilon_i^* \sim \mathcal{N}(0, V_i), \quad \text{mit } V_i = \sigma_\epsilon^2 I + Z_i Q Z_i^T \quad (2.19)$$

Die einzelnen Cluster/Gruppen können zu einem einzigen allgemeinen linearen gemischten Modell zusammengefasst werden.

Definition 2.2.3 (Allgemeines lineares gemischtes Modell).

Ein lineares gemischtes Modell ist definiert durch

$$y = X\beta + Zb + \epsilon$$

mit

$$\begin{pmatrix} b \\ \epsilon \end{pmatrix} \sim \mathcal{N} \left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} Q & 0 \\ 0 & R = \sigma_\epsilon^2 I \end{pmatrix} \right)$$

gegeben. Dabei sind X , bzw Z die Designmatrizen der festen, bzw zufälligen Effekte, β und b die Parametervektoren der festen und der zufälligen Effekten und ϵ der Fehlervektor.

In Konsequenz dessen lässt sich das marginale Modell verallgemeinern zu:

$$y = X\beta + \epsilon^* \quad (2.20)$$

wobei $\epsilon^* = Zb + \epsilon$ ist mit $\epsilon^* \sim \mathcal{N}(0, V)$ und $V = R + ZQZ^T$.

2.3 Likelihood Inferenz

Um die Verzerrung der Inferenz messen zu können, ist es zunächst erforderlich, die Theorie zur Likelihood-Inferenz von gemischten Modellen einzuführen. Dies umfasst sowohl die Schätzung der Parameter der zufälligen Effekte b_i als auch die Schätzung der Parameter β , σ_ϵ und Q . Um die Verzerrung zu quantifizieren, werden wir ein vollständiges gemischtes Modell mit einem reduzierten Modell ohne einen festen Effekt vergleichen. Dazu wird üblicherweise der sogenannte Likelihood-Ratio-Test (LRT) verwendet. Wie dieser Test genau funktioniert und wie der LRT durchgeführt wird, werden wir später erläutern. Zuvor benötigen wir noch etwas Theorie zur Likelihood-Berechnung.

2.3.1 Likelihood Berechnung gemischter Modelle

Im Folgenden wird die Schätzung der unbekannten Parameter erörtert. Der Vorliegende Ansatz basiert auf den Ausführungen von [fahrmeir-2011-regression].

Die Berechnung der Schätzer erfolgt mittels Maximum-Likelihood-Methode. Als Alternative kann die restringierte ML-Methode heran gezogen werden, die jedoch nicht für den Likelihood-Ratio-Test geeignet ist. Daher erfolgt die Berechnung der Parameter mittels der ML-Methode.

Die Schätzung der Parameter in einem gemischten Modell ist jedoch mit gewissen Schwierigkeiten verbunden. Neben dem β sind auch b_i , Q und σ_ϵ unbekannt. Daher ist es erforderlich, sowohl die festen und zufälligen Effekte als auch die unbekannten Parameter in Q und σ_ϵ , die wir als θ bezeichnen, zu schätzen. Dies bedingt eine geschachtelte Schätzung.

Im Folgenden wird zunächst angenommen, dass die Kovarianzen R , bzw. σ_ϵ , und Q bekannt sind. In diesem Zusammenhang ist auch V gemäß ?? bekannt. Für die Schätzung von β , ausgehend vom marginalen Modell, bietet sich

$$\hat{\beta} = (X^t V^{-1} X)^{-1} X^t V^{-1} y \quad (2.21)$$

an. Dieser Kleinste-Quadrate-Schätzer für β ergibt sich aus dem verallgemeinertem Kleinste-Quadrate-Kriterium (vgl. [KQ-Schätzer]), welches folgenden Term

$$(y - X\beta)^t V^{-1} (y - X\beta)$$

bezüglich β minimiert. Siehe hierzu auch [fahrmeir-2011-regression] (Kap. 3).

Der KQ-Schätzer ist gleichzeitig der log-Likelihood Schätzer unter der Normalverteilungsannahme. Die log-Likelihood für β aus dem marginalen Modell sieht folgendermaßen aus:

$$l(\beta) = -0.5 * (\log(|V|) + (y - X\beta)^t V^{-1} (y - X\beta) + N * \log(2\pi)).$$

Ableiten nach β ergibt den KQ-Schätzer aus ??.

$$\frac{d}{d\beta} l(\beta) = X^t V^{-1} (y - X\beta) \stackrel{!}{=} 0 \Rightarrow \hat{\beta} = (X^t V^{-1} X)^{-1} X^t V^{-1} y$$

Siehe hierzu auch [fahrmeir-2011-regression] (Kap. 3).

Gemäß dem Gauß-Markov-Theorem stellt $\hat{\beta}$ den besten linearen erwartungstreuen Schätzer (BLUE, best linear unbiased estimator) für die fixen Effekte dar. Zur Ermittlung des Schätzers ist lediglich eine Schätzung der Parameter in V sowie der Einsatz des Schätzers \hat{V} von V in $\hat{\beta}$ erforderlich.

Für den Schätzer von b verwenden wir den bedingten Erwartungswert $E(b|y)$ von b , gegeben die Daten y , welcher unter der Normalverteilungsannahme der beste Schätzer ist (vgl. [fahrmeir-2011-regression] Kap. 6.3.1).

Betrachtet man nun die gemeinsame Verteilung von b und y , welche folgendermaßen dargestellt wird:

$$\begin{pmatrix} y \\ b \end{pmatrix} \sim N \left(\begin{pmatrix} X\beta \\ 0 \end{pmatrix}, \begin{pmatrix} V & ZQ \\ QZ^t & Q \end{pmatrix} \right)$$

In Anbetracht dessen erhalten wir $E(b|y) = QZ^t V^{-1} (y - X\beta)$.

Ersetzt man nun β durch den Schätzer $\hat{\beta}$ erhält man den Schätzer für die zufälligen Effekte

$$\hat{b} = \hat{Q}Z^t \hat{V}^{-1} (y - X\hat{\beta}).$$

Der Schätzer \hat{b} ist der beste lineare unverzernte Schätzer (BLUP, best linear unbiased prediction)

Definition 2.3.1 (Schätzer für feste und zufällige Effekte).

Sei $y = X\beta + Zb + \epsilon$ ein lineares gemischtes Modell und $y = X\beta + \epsilon^*$ das zugehörige Marginale nach ?? . Dann ist

$$\hat{\beta} = (X^t \hat{V}^{-1} X)^{-1} X^t \hat{V}^{-1} y$$

ein Schätzer für die festen Effekte und

$$\hat{b} = \hat{Q} Z^t \hat{V}^{-1} (y - X\hat{\beta})$$

ein Schätzer für die zufälligen Effekte.

Wie bereits erwähnt, soll der Parametervektor θ alle unbekannten Parameter in $V = V(\theta)$, $Q = Q(\theta)$ und $\sigma_\epsilon = \sigma_\epsilon(\theta)$ enthalten. Anhand des Schätzers $\hat{\theta}$ lassen sich der Kovarianzschätzer sowie die Schätzer der festen und zufälligen Effekte berechnen. Die ML-Methode für θ basiert auf dem marginalen Modell

$$y \sim \mathcal{N}(X\beta, V(\theta)).$$

Die Log-Likelihood von β und θ ist gegeben durch

$$l(\beta, \theta) = -\frac{1}{2} (\log(|V|) + (y - X\beta)^t V^{-1} (y - X\beta)).$$

Maximieren von $l(\beta, \theta)$ bezüglich β für festes θ ergibt

$$\hat{\beta} = (X^t V^{-1} X)^{-1} X^t V^{-1} y.$$

Setzt man nun $\hat{\theta}$ in $l(\beta, \theta)$ ein, so erhält man die Profil-Log-Wahrscheinlichkeit

$$l(\theta)_p = -\frac{1}{2} (\log(|V|) + (y - X\hat{\beta})^t V^{-1} (y - X\hat{\beta})).$$

Folglich erhält man den ML-Schätzer $\hat{\theta}_{ML}$ durch Maximierung von $l(\theta)_p$.

Definition 2.3.2 (Kovarianz-Schätzer).

Sei $y = X\beta + Zb + \epsilon$ ein lineares gemischtes Modell mit

$$\begin{pmatrix} b \\ \epsilon \end{pmatrix} \sim \mathcal{N} \left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} Q & 0 \\ 0 & R \end{pmatrix} \right)$$

und sei θ der unbekannte Parametervektor von Q, R und $V = \text{Var}(y)$.

Dann ist $\hat{\theta}_{ML}$ der ML-Schätzer für θ , den man durch maximieren von

$$l(\theta)_p = -\frac{1}{2} (\log(|V|) + (y - X\hat{\beta})^t V^{-1} (y - X\hat{\beta}))$$

erhält. Dabei ist

$$\hat{\beta} = (X^t V^{-1} X)^{-1} X^t V^{-1} y$$

Mit dem Schätzer \hat{V} lassen sich die Schätzer der festen und zufälligen Effekte nun berechnen.

Um die Verzerrung der Inferenz messen zu können, müssen wir die log-Likelihood Werte berechnen können, um diese in den Likelihood-Ratio-Test einzusetzen. Der log-Likelihood Wert eines gemischten Modell ergibt sich aus der Maximum-Likelihood (ML)-Methode und ist folgendermaßen definiert:

Definition 2.3.3 (log-Likelihood Wert für ein gemischtes Modell).

Sei $r = y - X(X^t V^{-1} X)^{-1} X^t V^{-1} y$ und p der Rang von X

$$l_{ML}(Q, R) = -0.5 * (\log(|V|) + r^t V^{-1} r + N * \log(2\pi))$$

$$l_{REML}(Q, R) = -0.5 * (\log(|V|) + X^t V^{-1} X + r^t V^{-1} r + (N - p) * \log(2\pi))$$

$l_{REML}(Q, R)$ ist die eingeschränkte log-Likelihood, der sich aus der Methode "Restricted Maximim Likelihood" ergibt und entspricht im Wesentlichen der normalen log-Likelihood mit Ausnahme einer Differenz. Bei der "Restricted Maximim Likelihood" werden im Gegensatz zu der Methode "Maximum Likelihood" die Freiheitsgrade, die für die Schätzung fester Effekte bei der Schätzung von Varianzkomponenten verwendet werden, berücksichtigt. Im Gegensatz zum ursprünglichen Datenvektor basiert die eingeschränkte Maximum-Likelihood-Methode auf linearen Kombinationen der Beobachtungen, die so gewählt sind, dass diese Kombinationen invariant zu den Werten der festen Effektparametern sind.

Diese linearen Kombinationen sind äquivalent zu den Residuen, die nach der Anpassung durch normale kleinste Quadrate (gewichtet bei Angabe einer Regressionsgewichtung) lediglich den festen Effektanteil des Modells berechnen. Das Verfahren führt somit eine Maximierung in einem eingeschränkten Vektorraum durch.

2.3.2 Likelihood-Ratio-Test

Die Berechnung der Likelihood-Ratio-Test-Statistik (LRT-Statistik) ist relativ einfach, sofern die Theorie der ML-Methode vergegenwärtigt wird. Zur Erinnerung: Der Vergleich eines reduziertes Modells mit dem vollständigen Modell dient der Evaluierung des Einflusses einer Störgröße. Zur Durchführung dieser Analyse dient der Likelihood-Ratio-Test. Er ermöglicht den Vergleich eines einfacheren Modells (Nullmodell) mit einem komplexeren Modell (alternatives Modell), indem er die Likelihoods, bzw. die log-Likelihoods, der beiden Modelle vergleicht. Dies ist zum Beispiel nützlich um den Einfluss eines zusätzlichen Parameters zu beurteilen.

Der Likelihood-Ratio-Test wird wie folgt definiert:

Definition 2.3.4 (Likelihood-Ratio-Test (LRT)).

Sei L_{full} der Likelihood-Wert des vollständigen Modells sowie L_{red} der Likelihood-Wert des reduzierten Modells. Es sei i die Anzahl der Freiheitsgrade.

Dann ist die LRT Statistik gegeben durch

$$LRT = 2(\log L_{full} - \log L_{red})$$

Sofern die Größen L_{full} und L_{red} gemäß der Definition initialisiert sind, gilt $L_{full} > L_{red}$. Insbesondere gilt $\log(L_{full}) > \log(L_{red})$. Sofern die Log-Likelihood-Werte der Modelle bereits als L_{full} und L_{red} gegeben sind, lässt sich die LRT-Statistik durch $2(L_{full} - L_{red})$ berechnen.

Die Teststatistik des Likelihood-Ratio-Tests ergibt sich letztendlich, indem wir die LRT-Werte als Histogramm darstellen, und folgt einer χ^2 -Verteilung. Eine Chi-Quadrat-Verteilung mit k Freiheitsgraden ist folgendermaßen definiert:

Definition 2.3.5 (χ^2 -Verteilung).

Sei X_1, X_2, \dots, X_k eine Folge von unabhängigen standardnormalverteilten Zufallsvariablen, also $X_i \sim N(0, 1)$ für $i = 1, \dots, k$. Dann ist die Zufallsvariable

$$Y = \sum_{i=1}^k X_i^2$$

Chi-Quadrat-verteilt mit k Freiheitsgraden. Wir schreiben:

$$Y \sim \chi^2(k)$$

Die Wahrscheinlichkeitsdichtefunktion der χ^2 -Verteilung mit k Freiheitsgraden ist gegeben durch:

$$f(x; k) = \begin{cases} \frac{1}{2^{k/2}\Gamma(k/2)} x^{k/2-1} e^{-x/2} & x > 0, \\ 0 & x \leq 0, \end{cases}$$

wobei $\Gamma(\cdot)$ die Gamma-Funktion ist.

Die χ^2 -Verteilung bietet einen Vergleichswert für die Interpretation der Ergebnisse. Somit können wir feststellen, wie signifikant der Einfluss des zusätzlichen Parameters ist und ob die Anwendung gemischter Modelle im latenten Raum die Inferenz verzerrt. Um dies anschaulich darzustellen legt man das Histogramm der Teststatistik unter die χ^2 Verteilung. Dies erleichtert die Analyse, ob die Inferenz verzerrt ist. Die χ^2 -Verteilungen sind in Abbildung ?? veranschaulicht.

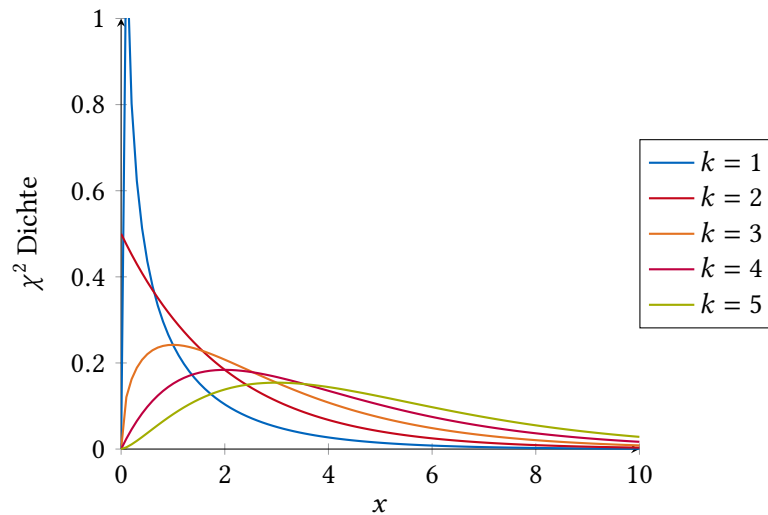


Abbildung 2.1 Chi-Quadrat-Verteilung für verschiedene Freiheitsgrade k .

3 Minimierung der Verzerrung

3.1 Datenbeschaffung

Quellen und Typen der verwendeten Daten

3.2 Modellierungstechniken

Beschreibung der spezifischen gemischten Modelle und der Techniken zur Gewinnung latenter Repräsentationen.

3.3 Analysemethoden

Verfahren zur Untersuchung der Verzerrung in den Inferenzergebnissen.

4 Experimente und Ergebnisse

Wir haben nun die nötigen theoretischen Kenntnisse, um die Verzerrung der Inferenz zu messen. Um die Verzerrung zu messen haben wir uns für eine Likelihood-Ratio-Test Statistik entschieden. Dabei fügen wir dem Datensatz einen neuen Parameter hinzu, welcher keinen Einfluss auf die Daten haben sollte. Bei unserem Datensatz haben wir jedem Patienten zufällig ein Geschlecht hinzugefügt. Das Geschlecht hat keinen Einfluss auf die Testergebnisse und demnach auch keinen Einfluss auf die Response-Variable. Nun können wir das vollständige, mit dem Geschlecht ergänzte, Modell mit dem reduzierten, ohne dem Geschlecht, Modell vergleichen, in dem wir den Likelihood-Ratio-Test anwenden. Um einen Richtwert zu haben, haben wir zuerst ein selbst überlegtes Simulationsdesign erstellt, welches einen niedrigdimensionalen medizinischen Datensatz simuliert. Dabei handelt es sich um einen Datensatz der die Herzgesundheit berechnet. Für jeden der $n=500$ Patienten wird Health-Score basierend auf verschiedensten Einflussfaktoren berechnet. Jeder der Patienten ist erkrankt und erhält nach frühestens drei Jahren eine Behandlung, die den Einfluss der einzelnen Parametern leicht verbessert. Der Start der Behandlung wird zufällig nach drei Jahren ausgemacht. Die Daten werden über 10 Jahre erhoben. Der Health-Score setzt sich aus dem diastolischen- und systolischen Blutdruck, dem Cholesterinspiegel, dem Triglyceride-Wert, dem Creatininspiegel und dem BMI zusammen. Die Werte werden für jeden Patienten aus einer Normalverteilung gezogen.

$$bp_{sys} \sim \mathcal{N}(120, 10)$$

$$bp_{dia} \sim \mathcal{N}(80, 10)$$

$$cholesterol \sim \mathcal{N}(200, 30)$$

$$triglyceride \sim \mathcal{N}(150, 20)$$

$$creatinin \sim \mathcal{N}(1, 0.2)$$

$$bmi \sim \mathcal{N}(25, 4)$$

Insbesondere wird jedem Patienten ein zufälliges Alter zwischen 30 und 60 Jahren zugeteilt, welches ebenfalls einen minimal negativen Effekt auf den Health-Score hat. Als zusätzlichen Effekt, welcher keinen Einfluss auf die Response-Variable (in diesem Fall der Health-Score) hat, wird jedem Patienten ein Geschlecht zugeteilt. Für immer zufällig auftretende Effekte wird ein Random Slope und ein Random Intercept in den Health-Score hinzugefügt.

Basierend auf diesem Simulationsdesign, welches einem Gemischten Modell folgt, können wir nun einen Likelihood-Ratio-Test durchführen. Wir trainieren dazu ein vollständiges gemischtes Modell und ein reduziertes gemischtes Modell ohne den Effekt 'Geschlecht'. Mit den log-Likelihood Werten für die trainierten Modelle führen wir den Likelihood-Ratio-Test durch. Nach 500 Wiederholungen ergibt sich ein Histogramm der Likelihood-Ratio-Test Statistik, welches, wie zu erwarten, einer χ^2 Verteilung folgt. Wie wir in Abb sieht folgt das Histogramm der Test Statistik der Roten Kurve, welche die χ^2 Verteilung beschreibt, ohne Verzerrung. Bis auf einzelne Ausnahmen, welche durch Instabilitäten der Berechnung immer verursacht werden können, sind die Ergebnisse immer unter der χ^2 Verteilung. Dies war allerdings auch so zu erwarten, da wir ein ganz normales gemischtes Modell betrachtet hatten.

4.0.1 Gemischtes Modell auf latenter Datenwolke mit separatem Training

Nun wollen wir das gemischte Modell in einer latenten Repräsentation betrachten. Dazu wählen wir, wie schon angeführt, einen Variational-Autoencoder. Außerdem haben wir nun einen hochdimensionalen medizinischen Datensatz gegeben. Dieser ist einem echten Datensatz bestmöglich nachgebaut, allerdings kann hier aus Datenschutzgründen kein wirklich echter Datensatz benutzt werden. Wir wählen

zunächst ein recht simples Encoder-Modell mit zwei Schichten und einer zweidimensionalen latenten Dimension. Wir trainieren zuerst den VAE separat von den gemischten Modellen. Dazu optimieren wir in der Loss-Funktion den Reconstruction Loss und die KL-Divergenz. Das vollständige und reduzierte gemischte Modell trainieren wir dann auf der latenten Datenwolke jeweils nach dem abgeschlossenen Training des VAEs. So erhalten wir wieder zwei log-likelihood Werte welche wir mit dem Likelihood-Ratio-Test auswerten können. Fassen wir alle LRT-Werte gleichermaßen wie zuvor in einem Histogramm zusammen und vergleichen mit der χ^2 Verteilung, so sehen wir, dass eine bedeutende Masse dieses Mal über der χ^2 Verteilung liegt. Wir sehen also, dass es zu einer Verzerrung der Inferenz kommt.

4.0.2 Gemischtes Modell auf latenter Datenwolke mit gleichzeitigem Training

Wenn wir nun versuchen das gemischte Modell zusammen mit dem VAE in einer einzigen Loss-Funktion zu trainieren, sehen wir recht schnell, dass wir so nicht zu einem gewünschten Ergebnis kommen. Bei einem gemeinsamen Training unter der Voraussetzung, dass in der Loss-Funktion alle Parameter gleich gewichtet sind, geht die χ^2 Verteilung komplett verloren. Fügen wir der Lossfunktion auch nur den Mean-Squared-Error zwischen dem Encoder Output und dem Output des vollständig trainierten gemischten Modells hinzu und gewichten diesen nur minimal, so verlieren wir schon die χ^2 -Verteilung. Das bedeutet sobald der Encoder Einfluss auf das gemischte Modell hat, geht die gewünschte Verteilung verloren. Dementsprechend geht die Verteilung auch verloren, wenn wir

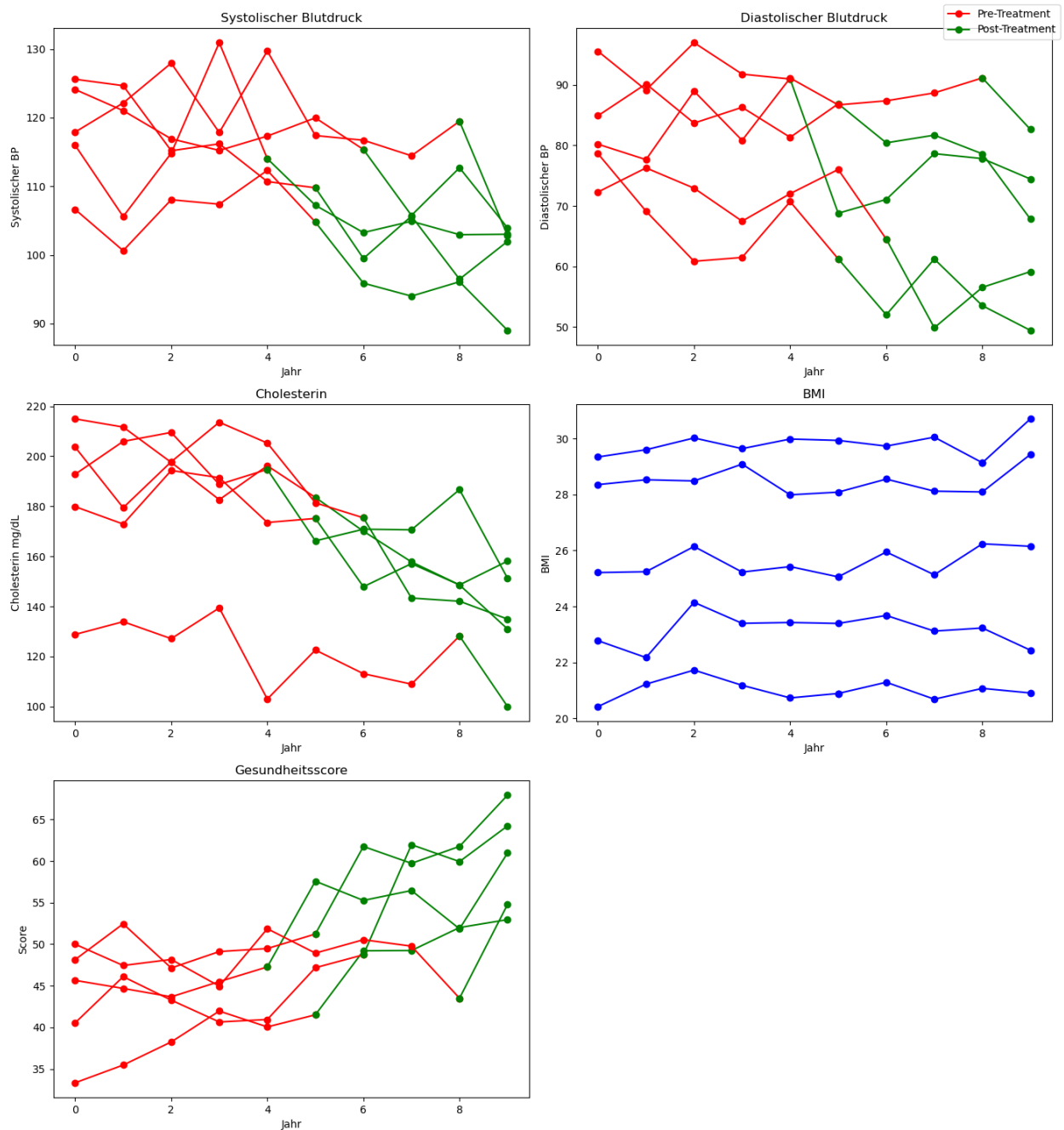


Abbildung 4.1 Simulierte Datensätze für 20 zufällig ausgewählte Patienten

5 Methodik

5.1 Vorgehen

In den ersten Wochen habe ich mir selbst ein Simulationsdesign für einen longitudinalen medizinischen Datensatz ausgedacht und basierend darauf ein gemischtes Modell gefittet. Mit diesen simulierten Daten habe ich ein reduziertes Modell mit dem vollständigen Modell verglichen. Die LRT Statistik habe ich dann in einem Histogramm dargestellt.

Wir fügen dem gemischten Modell einen festen Effekt hinzu, welcher keinen Einfluss auf die Trajektorie haben soll. In unserem Fall ist dieser feste Effekt das Geschlecht, welches keinen Einfluss auf den Verlauf einer Krankheit haben sollte.

Mein zweites Projekt ist nun einen hoch dimensionalen medizinischen Datensatz durch den Encoder eines Variational Autoencoders im latenten Raum zu repräsentieren und dort mit einem gemischten Modell darzustellen. Ähnlich wie zuvor will ich wieder eine LRT Statistik erhalten, in dem ich ein reduziertes Modell mit dem vollständigen Modell vergleiche. Dazu trainiere ich in einer Schleife den Encoder und das gemischte Modell für jeden Iterationsschritt neu und vergleiche die negativen Maximum Likelihood-Werte (ML-Werte) durch den Likelihood Ratio Test. Am Ende der Schleife erhalte ich wieder eine LRT Statistik, welche durch ein Histogramm dargestellt wird. Im Optimalfall ähnelt das Histogramm einer Chi-Quadrat-Verteilung mit einem Freiheitsgrad (Da das reduzierte Modell nur einen festen Effekt, das Geschlecht, weniger hat).

5.2 Experimentelles Design

Aufbau der experimentellen Tests und Simulationen.

5.3 Durchführung

Beschreibung der durchgeführten Experimente und verwendeten Parameter.

5.4 Analyse der Ergebnisse

Diskussion der Ergebnisse im Hinblick auf die Verzerrung der Inferenz.

6 Diskussion

6.1 Interpretation der Ergebnisse

Tiefere Analyse der Ergebnisse und ihrer Implikationen.

6.2 Vergleich mit bestehenden Arbeiten

Wie sich die Ergebnisse zu bereits veröffentlichten Forschungen verhalten.

6.3 Limitationen und Herausforderungen

Kritische Betrachtung der Grenzen der Studie und mögliche Probleme.

7 Fazit

Zusammenfassung der wichtigsten Erkenntnisse Praktische Implikationen: Wie die Ergebnisse in der Praxis angewendet werden können. Empfehlungen für zukünftige Forschungen: Vorschläge für weiterführende oder ergänzende Studien.

8 Anhang

A Appendix

A.1 Supporting Data

A.2 Some Code Listings