

Seminararbeit
M.Sc.-Studiengang ”Betriebswirtschaftslehre”

UNIVERSITÄT HAMBURG
FAKULTÄT FÜR BETRIEBSWIRTSCHAFT
LEHRSTUHL FÜR STATISTIK MIT ANWENDUNG IN DER
BETRIEBSWIRTSCHAFTSLEHRE
PROF. DR. MARTIN SPINDLER

**Data Profiling mit Large Language Models -
Eine Studie am Beispiel des US
Gender-Pay-Gap**

Gutachter: Prof. Dr. Martin Spindler

eingereicht von:

Vincent-Konstantin Kapp

7778027

Methfesselstraße 16

20257 Hamburg

0152-56761933

betreut von:

Prof. Dr. Martin Spindler

Abgabedatum, Ort:

07.01.2025, Hamburg

Einleitung

Abstract

Datengetriebene Prozesse können Unternehmen nicht nur dabei helfen, ihre Prozesse zu optimieren, sondern auch bei der Entwicklung neuer Strategien. Um festzustellen, welche Daten in den Datenbanken vorliegen und wie diese zueinander stehen, braucht es Mitarbeiter*innen, die einzeln überprüfen, welche Daten vorhanden sind und in welchen Zusammenhängen sie basierend auf Common Sense existieren. Data Profiling stellt Unternehmen mit historisch gewachsenen Datenbanksystemen vor Herausforderungen für die datengestützte Transformation. Können die aktuellen Entwicklungen von Large Language Modellen (LLMs) dazu beitragen, automatisiert die Zusammenhänge von Spalten zu erkennen? Um diese Frage zu beantworten, widmet sich diese Seminararbeit, aufbauend auf der Arbeit von Trummer (2024), der Untersuchung, wie gut LLMs anhand einer Case Study die Korrelation von Daten erkennen können. Dazu wird ein LLM auf Basis von GPT-3 trainiert und mit einem Benchmark-Verfahren verglichen. Die Ergebnisse zeigen, dass LLMs eine hohe Genauigkeit bei der Korrelationserkennung aufweisen und damit Unternehmen bei der Datenanalyse unterstützen können.

Table of contents

1	Hintergrund	4
1.1	Data Profiling	4
1.2	Korrelationserkennung in Datenbanken	4
1.3	Sprachmodelle und deren Anwendung in der Datenanalyse	4
2	Methodik der Fallstudie	5
2.1	Datenbasis und Vorbereitung	5
2.2	Vorgehensweise zur Korrelationsanalyse mit LLM	5
2.3	Benchmark und Bewertungsmetriken	5
3	Durchführung der Fallstudie: Gender-Pay-Gap-Analyse	6
3.1	Fragestellung	6
3.2	Implementierung des LLM für Korrelationserkennung	6
3.3	Benchmarking-Ergebnisse	6
4	Vergleich und Auswertung der Ergebnisse	7
4.1	Vergleich der Ergebnisse mit Prompoting	7
4.2	Szenariobasierte Auswertung	7
4.3	Bewertung der Korrelationserkennung	7
5	Diskussion der Ergebnisse und Implikationen	8
5.1	Interpretation der Ergebnisse	8
5.2	Limitationen	8
5.3	Praktische Anwendungsmöglichkeit	8
6	Fazit und Ausblick	9
7	Eidesstattliche Erklärung	10
	Eidesstattliche Erklärung	10
	References	11

List of Figures

List of Tables

Chapter 1

Hintergrund

1.1 Data Profiling

1.2 Korrelationserkennung in Datenbanken

(Trummer 2023).

1.3 Sprachmodelle und deren Anwendung in der Datenanalyse

(Arora et al. 2023).

Chapter 2

Methodik der Fallstudie

2.1 Datenbasis und Vorbereitung

2.2 Vorgehensweise zur Korrelationsanalyse mit LLM

2.3 Benchmark und Bewertungsmetriken

Chapter 3

Durchführung der Fallstudie: Gender-Pay-Gap-Analyse

Meine Überlegung ist es das Ergebnisse von (Bach, Chernozhukov, and Spindler 2024).

3.1 Fragestellung

3.2 Implementierung des LLM für Korrelationserkennung

3.3 Benchmarking-Ergebnisse

Chapter 4

Vergleich und Auswertung der Ergebnisse

4.1 Vergleich der Ergebnisse mit Prompoting

4.2 Szenariobasierte Auswertung

4.3 Bewertung der Korrelationserkennung

Chapter 5

Diskussion der Ergebnisse und Implikationen

5.1 Interpretation der Ergebnisse

5.2 Limitationen

5.3 Praktische Anwendungsmöglichkeit

Chapter 6

Fazit und Ausblick

Chapter 7

Eidesstattliche Erklärung

Eidesstattliche Erklärung

“Ich versichere, dass ich die vorstehende Arbeit selbstständig und ohne fremde Hilfe angefertigt und mich anderer als der im beigefügten Verzeichnis angegebenen Hilfsmittel nicht bedient habe. Alle Stellen, die wörtlich oder sinngemäß aus Veröffentlichungen übernommen wurden, sind als solche kenntlich gemacht. Alle Internetquellen sind der Arbeit beigefügt. Des Weiteren versichere ich, dass ich die Arbeit vorher nicht in einem anderen Prüfungsverfahren eingereicht habe und dass die eingereichte schriftliche Fassung der auf dem elektronischen Speichermedium entspricht.”

~\[20mm] Ort, Datum

Unterschrift

References

- Arora, Simran, Brandon Yang, Sabri Eyuboglu, Avanika Narayan, Andrew Hojel, Immanuel Trummer, and Christopher Ré. 2023. “Language Models Enable Simple Systems for Generating Structured Views of Heterogeneous Data Lakes.” Journal Article. *arXiv Preprint arXiv:2304.09433*.
- Bach, Philipp, Victor Chernozhukov, and Martin Spindler. 2024. “Heterogeneity in the US Gender Wage Gap.” *Journal of the Royal Statistical Society Series A: Statistics in Society* 187 (1): 209–30.
- Trummer, Immanuel. 2023. “Can Large Language Models Predict Data Correlations from Column Names?” Journal Article. *Proceedings of the VLDB Endowment* 16 (13): 4310–23.