



**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

<b>1</b>	<b>Hintergrund</b>	<b>4</b>
1.1	Data Profiling . . . . .	4
1.2	Erkennen von Korrelationen in Datens . . . . .	4
1.3	Language Models . . . . .	4
1.4	Nutzung von NLP für Datenbankanalyse . . . . .	4
1.5	Heterogenität in Datensätzen: am Beispiel des US Gender-Pay-Gap	4
<b>2</b>	<b>Methodik</b>	<b>5</b>
2.1	Setup der Studie . . . . .	5
2.1.1	Datenbasis und Vorbereitung . . . . .	5
2.1.2	Vorgehensweise zur Korrelationsanalyse mit LLM's . . . . .	5
2.2	Benchmark-Metriken und Analyse . . . . .	5
2.3	Benchmark-Analyse . . . . .	5
<b>3</b>	<b>Ergebnisse</b>	<b>6</b>
3.1	Darstellung der Ergebnisse . . . . .	6
3.2	Vergleich der Ergebnisse . . . . .	6
<b>4</b>	<b>Diskussion</b>	<b>7</b>
4.1	Diskussion der Ergebnisse . . . . .	7
4.2	Stärken und Limitationen der Arbeit . . . . .	7
4.3	Praktische Anwendungsmöglichkeit . . . . .	7
<b>5</b>	<b>Eidesstattliche Erklärung</b>	<b>8</b>
	Eidesstattliche Erklärung . . . . .	8
	<b>References</b>	<b>9</b>

# List of Figures

# List of Tables

# Chapter 1

## Hintergrund

### 1.1 Data Profiling

### 1.2 Erkennen von Korrelationen in Datensätzen

(Trummer 2023).

### 1.3 Language Models

(Arora et al. 2023).

### 1.4 Nutzung von NLP für Datenbankanalyse

### 1.5 Heterogenität in Datensätzen: am Beispiel des US Gender-Pay-Gap



## Chapter 2

# Methodik

### 2.1 Setup der Studie

#### 2.1.1 Datenbasis und Vorbereitung

#### 2.1.2 Vorgehensweise zur Korrelationsanalyse mit LLM's

### 2.2 Benchmark-Metriken und Analyse

### 2.3 Benchmark-Analyse

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

## Chapter 3

# Ergebnisse

### 3.1 Darstellung der Ergebnisse

### 3.2 Vergleich der Ergebnisse

## Chapter 4

# Diskussion

4.1 Diskussion der Ergebnisse

4.2 Stärken und Limitationen der Arbeit

4.3 Praktische Anwendungsmöglichkeit

## Chapter 5

# 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.