

# Abgabe 1 für Computergestützte Methoden

Gruppe (88), Anas El Fahsi (4242194), Arda Bostancı (4339434),  
Felix Norbert Figge (4361841)

(2. Dezember 2024)

## Inhaltsverzeichnis

<b>1</b>	<b>Der zentrale Grenzwertsatz</b>	<b>2</b>
1.1	Aussage . . . . .	2
1.2	Erklärung der Standardisierung . . . . .	2
1.3	Anwendungen . . . . .	2
<b>2</b>	<b>Bearbeitung zur Aufgabe 1</b>	<b>3</b>
2.1	Berechnung per Tabellenkalkulation . . . . .	3
2.2	Berechnung mit SQL . . . . .	4
<b>3</b>	<b>Github</b>	<b>5</b>

# 1 Der zentrale Grenzwertsatz

Der zentrale Grenzwertsatz (ZGS) ist ein fundamentales Resultat der Wahrscheinlichkeitstheorie, das die Verteilung von Summen unabhängiger, identisch verteilter (*i.i.d.*) Zufallsvariablen (ZV) beschreibt. Er besagt, dass unter bestimmten Voraussetzungen die Summe einer großen Anzahl solcher ZV annähernd normalverteilt ist, unabhängig von der Verteilung der einzelnen ZV. Dies ist besonders nützlich, da die Normalverteilung gut untersucht und mathematisch handhabbar ist.

## 1.1 Aussage

Sei  $X_1, X_2, \dots, X_n$  eine Folge von i.i.d. ZV mit dem Erwartungswert  $\mu = \mathbb{E}(X_i)$  und der Varianz  $\sigma^2 = \text{Var}(X_i)$ , wobei  $0 < \sigma^2 < \infty$  gelte. Dann konvergiert die standardisierte Summe  $Z_n$  dieser ZV für  $n \rightarrow \infty$  in Verteilung gegen eine Standardnormalverteilung:

$$Z_n = \frac{\sum_{i=1}^n X_i - n\mu}{\sigma\sqrt{n}} \xrightarrow{d} \mathcal{N}(0, 1). \quad (1)$$

Das bedeutet, dass für große  $n$  die Summe der ZV näherungsweise normalverteilt ist mit Erwartungswert  $n\mu$  und Varianz  $n\sigma^2$ :

$$\sum_{i=1}^n X_i \sim \mathcal{N}(n\mu, n\sigma^2). \quad (2)$$

## 1.2 Erklärung der Standardisierung

Um die Summe der ZV in eine Standardnormalverteilung zu transformieren, subtrahiert man den Erwartungswert  $n\mu$  und teilt durch die Standardabweichung  $\sigma\sqrt{n}$ . Dies führt zu der obigen Formel (1). Die Darstellung (2) ist für  $n \rightarrow \infty$  nicht wohldefiniert.

## 1.3 Anwendungen

Der ZGS wird in vielen Bereichen der Statistik und der Wahrscheinlichkeitstheorie angewendet. Typische Beispiele sind:

- Berechnung von Konfidenzintervallen
- Durchführung von Hypothesentests

<sup>1</sup>Der zentrale Grenzwertsatz hat verschiedene Verallgemeinerungen. Eine davon ist der Lindeberg-Feller-Zentrale-Grenzwertsatz [[1], Seite 328], der schwächere Bedingungen an die Unabhängigkeit und die identische Verteilung der ZV stellt.

<sup>2</sup>Bei der Schätzung der oberen und unteren Grenze wird eine Transformation in die Standardnormalverteilung mittels des zentralen Grenzwertsatzes verwendet [[2], Kapitel 3.5 unter „Konfidenzintervalle“, Seite 101-102]

<sup>3</sup>Zur Bestimmung der Irrtumswahrscheinlichkeit werden viele Stichproben gezogen bis nach dem zentralen Grenzwertsatz eine Normalverteilung vorliegt [[2], Kapitel 4 unter „Bestimmung der Irrtumswahrscheinlichkeit“, Seite 112]

## 2 Bearbeitung zur Aufgabe 1

In dieser Aufgabe sollen wir aus einem vorgegebenen Datensatz die höchste mittlere Temperatur berechnen – zum einen mithilfe einer Tabellenkalkulation und zum anderen durch das Entwerfen eines Datenbankschemas, das mit den importierten Daten in SQL umgesetzt wird. Anschließend berechnen wir die durchschnittliche Höchsttemperatur mit einer passenden SQL-Abfrage.

### 2.1 Berechnung per Tabellenkalkulation

Als erstes haben wir den Datensatz in Excel importiert und nach unserer Gruppennummer 88 gefiltert, damit wir nur mit den für uns relevanten Teil arbeiten. Anschließend haben wir auf die Spalte *average\_temperature* einen Filter angewendet, der die leeren Zellen ausschließt, damit wir bei der Abfrage nach der höchsten mittleren Temperatur keine Probleme haben. Dann haben wir mit dem MAX(...) Befehl in Excel den gewünschten Wert bestimmt 1.

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	
1															
2					station	date	day_of_year	day_of_week	month_of_year	precipitation	windspeed	min_temp	average_temp	max_temp	count
3					W45 St & Ave	01.01.23	1	1	0	100	42	56	133		
4					W45 St & Ave	02.01.23	2	2	1	648	39	46	55	119	
5					W45 St & Ave	03.01.23	3	3	1	38	582	46	48	50	98
6					W45 St & Ave	04.01.23	4	4	1	649	45	51	61	129	
7					W45 St & Ave	05.01.23	5	5	1	626	45	49	51	107	
					W45 St & Ave	06.01.23	6	6	1	805	39	44	47	121	

Abbildung 1: Zur Berechnung der höchsten mittleren Temperatur wurde über die Spalte L von den Zeilen 2 bis 363 der MAX(...) Befehl verwendet, wodurch wir auf den gesuchten Wert kamen.

Das Ergebnis sind  $83^{\circ}$  F. Die Einheit ist Fahrenheit, da dem Datensatz zu entnehmen ist, dass die Temperaturen in einer Straße in New York gemessen wurden. Um die Temperatur in Grad Celsius zu berechnen nehmen wir uns folgende Formel zur Hilfe:

$$^{\circ}\text{C} = \frac{5}{9}(\text{ }^{\circ}\text{F} - 32). \quad (3)$$

Setzt man  $83^{\circ}$  F in die Formel ein erhält man als Ergebnis  $28,3^{\circ}$  C.

---

<sup>4</sup>Zur Umrechnung von Grad Fahrenheit zu Grad Celsius wurde die oben stehende Formel aus dem Buch von Tipler herangezogen [[3], Kapitel 13.2, Seite 491, Formel 13.2.]

## 2.2 Berechnung mit SQL

Zuallererst haben wir die Zeilen entfernt, bei denen in der Spalte *average\_temperature* kein Wert vorhanden war, um die spätere Abfrage in SQL zu erleichtern. Basierend auf den Prinzipien der Datenbanknormalisierung haben wir folgendes Schema erstellt:

1. Eine dates-Tabelle mit einer id-Spalte als Primärschlüssel und den Spalten für Datum und Station.
2. Eine measurements-Tabelle, die alle Messungen (z. B. Temperaturen, Niederschläge) speichert und über einen Fremdschlüssel (die id) mit der dates-Tabelle verknüpft ist.

So schaffen wir es, Redundanzen zu minimieren und sicherzustellen, dass beide Tabellen die 1. und 2. Normalform erfüllen. Wir haben also zunächst die beiden CSV-Dateien mit den richtigen Spalten erstellt, die Tabellen in SQL definiert und die Daten mithilfe der SQL-Befehle aus dem Skript in die jeweiligen Tabellen importiert (siehe Abbildung 2).

```
sqlite> PRAGMA foreign_keys = ON;
sqlite> CREATE TABLE dates (
    id INTEGER PRIMARY KEY,
    date DATE NOT NULL,
    station TEXT NOT NULL
);
sqlite> CREATE TABLE measurements (
    id INTEGER PRIMARY KEY,
    date_id INTEGER,
    day_of_week INTEGER,
    temperature REAL,
    precipitation REAL,
    windspeed REAL,
    average_min_temp REAL,
    average_max_temp REAL,
    count REAL,
    FOREIGN KEY (date_id) REFERENCES dates(id) -- Verweis auf die id der Tabelle dates
);
sqlite> -- Importieren der CSV Dateien in die jeweiligen Tabellen
sqlite> .separator ;
sqlite> .import /Users/anaeelfahsi/Desktop/Abschae_1/dates_id.csv dates
sqlite> .import /Users/anaeelfahsi/Desktop/Abschae_1/measurements1.csv measurements
```

Abbildung 2: Hier haben wir in Sqlite3 auf Mac die beiden Tabellen erstellt und die jeweiligen CSV-Dateien in die Tabellen importiert

Dann haben wir mithilfe der Abfrage nach der höchsten mittleren Temperatur den Wert 83 erhalten (Abbildung 3).

```
sqlite> -- Abfrage nach der höchsten mittleren Temperatur
sqlite> SELECT MAX(average_temperature) AS max_average_temperature
sqlite> FROM measurements;
83.0
```

Abbildung 3: Abfrage nach der höchsten mittleren Temperatur aus der measurements Tabelle

Wie bei der Berechnung mit Excel bereits erwähnt haben wir Temperaturen in Fahrenheit. Der Wert 83 ist derselbe, wie der den wir bei der Abfrage mit Excel erhalten haben und damit erhalten wir wieder eine Temperatur von 28,3° C.

### 3 Github

Wir haben unser Overleaf Dokument auf Github als ZIP-Datei und nochmal die einzelnen Dateien hochgeladen. Klicken Sie auf diesen Text um den Link zu öffnen.

---

<sup>6</sup>Das Datenbankschema wurde nach der 1. und 2. Normalform erstellt in der es heißt, dass alle Attribute atomar sein müssen und die Tabellen in mehrere Tabellen aufgetrennt werden müssen und mit Fremdschlüsseln in Beziehung gesetzt werden. [[4], unter dem Punkt „Vorüberlegungen: Normalformen“, S. 8.]

<sup>5</sup>Die DDL Befehle haben wir aus dem Skript [[4], unter dem Punkt „Schritt 1: Übersetzung in SQL“, S. 11.]

## Literatur

- [1] Achim Klenke. *Wahrscheinlichkeitstheorie*. Springer, 3. edition, 2013.
- [2] Jürgen Bortz. *Statistik für Human- und Sozialwissenschaften*. Springer, 6. edition, 2005.
- [3] Paul A. Tipler and Gena Mosca. *Tipler Physik für Studierende der Naturwissenschaften und Technik*. Springer, 9. edition, 2024.
- [4] Mike Teßmer. Praktische datenhaltung und -verarbeitung. Online, n.d.  
Zugriffsdatum: 28.11.2024.