

## Mogelzettel: Deskriptive Statistik

Betrachtet wird eine **Grundgesamtheit** quantitativer Daten  $\{x_1, x_2, \dots, x_N\}$ , die *alle* möglichen Daten<sup>1</sup> enthält. Dieser unhandliche Datensatz wird durch "typische" Zahlen charakterisiert.

Arithmetischer **Mittelwert**: 
$$\langle x \rangle := \bar{x} := \frac{1}{N} \sum_{i=1}^N x_i$$

Entsprechend kann der Mittelwert abgeleiteter Größen  $f(x_i)$  (z.B. Kugelvolumina für gegebene Radien) definiert werden:

$$\langle f(x) \rangle := \overline{f(x)} := \frac{1}{N} \sum_{i=1}^N f(x_i)$$

### Klasseneinteilung

Bei einer großen Anzahl von Werten wird häufig mit einer Klasseneinteilung und Histogrammen gearbeitet. Dies erfolgt oft bereits bei der Datenerhebung.

**Klassen**: disjunkte, aufeinanderfolgende Intervalle, dh.  $[x_i^u, x_i^o]$  oder  $(x_i^u, x_i^o]$   $i = 1, \dots, k$

**Klassenbreite**:  $\Delta x_i = x_i^o - x_i^u = x_{i+1}^u - x_i^u$

Als Repräsentant der Klasse wird oft die **Klassenmitte** (Intervallmitte)  $x_i = \frac{x_i^o + x_i^u}{2}$  benutzt. Ist  $n_i$  die Anzahl der Werte in der i-ten Klasse, so definiert man:

$$N = \sum_{i=1}^k n_i, \quad \bar{x} := \frac{1}{N} \sum_{i=1}^k n_i x_i, \quad \overline{f(x)} := \frac{1}{N} \sum_{i=1}^k n_i f(x_i)$$

was eigentlich nur bei symmetrischen Verteilungen korrekte Resultate liefert.

### Alternativen zum arithmetischen Mittelwert:

**Geometrisches Mittel**:  $\langle x \rangle_g = \sqrt[n]{x_1 x_2 \cdots x_n}$  (wichtig für Wachstumsraten)

**quadratisches Mittel**:  $\langle x \rangle_{\text{rms}} = \sqrt{\frac{x_1^2 + \cdots + x_n^2}{n}}$  (z.B. Effektivwerte in der Elektrik)

harmonisches Mittel:  $\langle x \rangle_h = \frac{n}{\frac{1}{x_1} + \frac{1}{x_2} + \cdots + \frac{1}{x_n}}$

Im Alltag sind auch noch der **Median** (Wert in der Mitte) und der **Modus** (häufigster Wert) gebräuchlich.

---

<sup>1</sup>also keine Stichprobe!

Neben dem Mittelwert ist auch die Breite der Verteilung wichtig, was auf die Varianz und Standardabweichung führt:

**Varianz:** 
$$V(x) := \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 = \dots = \frac{1}{n} \sum_{i=1}^n x_i^2 - \left( \frac{1}{n} \sum_{i=1}^n x_i \right)^2 = \overline{x^2} - \bar{x}^2$$

**Standardabweichung:** 
$$\sigma = \sqrt{V(x)}$$
 (hat dieselbe Einheit wie x)

Details sind abhängig von der Verteilung, aber als Faustformel:

$\sigma$  ist eine typische Abweichung vom Mittelwert, Abweichungen um  $2\sigma$  kommen vor, aber Abweichungen von mehr als  $3\sigma$  sind verdächtig!

### Alternative Breitenmaße:

In der Physik ist die **volle Halbwertsbreite (FWHM)**<sup>2</sup> gebräuchlich, da sie einfach zu bestimmen ist und für viele Verteilungen eine unmittelbare Berechnung von  $\sigma$  erlaubt (für Gaußverteilungen gilt:  $\text{FWHM} \approx 2.35\sigma$ ).

Die zum Median gehörenden Breitenmaße sind die Quartile und Dezile; z.B. liegen beim unteren Quartil 25% der Daten darunter und 75% darüber, beim oberen Dezil liegen 90% darunter und 10% darüber.

Seltener: **Spannweite**  $\max(x_i) - \min(x_i)$  und **mittlere absolute Abweichung**  $\frac{1}{n} \sum_i |x_i - \bar{x}|$

### höhere Momente:

**m-tes (zentrales) Moment:** 
$$\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^m$$
 "Momentenentwicklung"

Höhere Momente werden in der Physik eher selten benutzt; Statistiker benutzen auch renormierte Versionen wie *Schiefte* und *Kurtosis*, mit von Autor zu Autor teilweise leicht abweichenden Definitionen.

### Paare von Daten - Korrelation

Datensätze:  $(x_1, y_1), \dots, (x_n, y_n)$  (z.B. Größe und Gewicht)

Information über Zusammenhang zwischen den Größen liefert die **Kovarianz**:

$$\text{cov}(x, y) := \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) = \dots = \overline{xy} - \bar{x}\bar{y}$$

$> 0$  falls tendenziell große x mit großen y auftreten

$< 0$  falls große x mit kleinen y auftreten (oder umgekehrt)

$= 0$  falls kein Zusammenhang besteht

**Korrelationskoeffizient:** 
$$-1 \leq \rho = \frac{\text{cov}(x, y)}{\sigma_x \sigma_y} \leq 1$$

Eine Korrelation sagt nicht notwendigerweise etwas über kausale Zusammenhänge aus! (Zahl der Störche / Geburtenrate)

Eine Verallgemeinerung auf k Variablen führt auf die  $k \times k$  Kovarianzmatrix  $V_{ij}$ .

<sup>2</sup>Full Width at Half Maximum