

Gebäuchliche Bezeichnungen:	
<ul style="list-style-type: none"> • <i>Merkmal</i>: interessierendes Kennzeichen, Charakteristikum bzw. interessierende Größe; ein Merkmal ist dabei <i>qualitativ</i>: nicht durch Messen oder Zählen erfassbar <i>quantitativ</i>: durch Messen oder Zählen erfassbar • <i>Merkmalsträger</i>: Objekte (Elemente), welche die in Frage kommenden Merkmale aufweisen • <i>Merkmalsausprägung (Beobachtungs- bzw. Merkmalswert) x_i</i>: eine konkrete Beschreibung oder Realisierung bei einem qualitativen bzw. ein konkretes Mess- oder Zählergebnis bei einem quantitativen Merkmal Bemerkung: Schulnoten, Güteklassen usw. sind trotz der Zahlenangaben <i>keine quantitativen sondern qualitative</i> Merkmale (die Zahlen dienen der Reihung)! • <i>diskretes</i> quantitatives Merkmal: nur endlich viele, meist ganzzahlige Werte sind möglich (Zählergebnisse) <i>stetiges</i> quantitatives Merkmal: beliebige reelle Zahlen eines Intervalls können Beobachtungswerte sein (Messergebnisse) • <i>Grundgesamtheit</i>: die Menge aller (vor allem: sehr vieler) Merkmalsträger mit dem interessierenden Merkmal • <i>Umfang N der Grundgesamtheit</i>: Anzahl der Elemente der Grundgesamtheit • <i>Stichprobe</i>: eine <i>zufällige</i> Auswahl von Merkmalsträgern aus der Grundgesamtheit <i>zufällig</i>: gleiche Auswahlchancen für alle Objekte (oft problematisch, die Chancen-Ungleichheit zu erkennen) • <i>Stichprobenumfang n</i>: Anzahl der Elemente der Stichprobe • <i>Urliste</i>: bezüglich der Merkmalsausprägungen noch ungeordnete Datenerhebung • <i>Strichliste</i>: bezüglich der Beobachtungsergebnisse geordnete Datenliste; jedes Auftreten einer bestimmten Ausprägung wird durch einen Strich registriert (Blockbildung ist dabei empfehlenswert) • <i>Absolute Häufigkeit (bzw. Besetzungszahl) n_i</i>: die Anzahl der Objekte mit der Merkmalsausprägung x_i • <i>Relative Häufigkeit $h(x_i)$ bzw. h_i</i>: der Anteil der Objekte mit der Ausprägung x_i am Gesamtumfang der Daten (wird oft in Prozent angegeben) Bemerkung: relative Häufigkeiten sind nie negativ und nie größer als eins: $0 \leq h_i \leq 1$ • <i>Absolute Summenhäufigkeit (kumulierte Häufigkeit) N_k</i>: Summe der Besetzungszahlen von n_1 bis n_k • <i>Relative Summenhäufigkeit (relative kumulierte Häufigkeit) $H(x_k)$ bzw. H_k</i>: Summe der relativen Häufigkeiten von h_1 bis h_k Bemerkung: Die relative Summenhäufigkeit ist eine von null bis eins monoton steigende Funktion. 	<ul style="list-style-type: none"> a) Wohnort; Haarfarbe b) Länge; Druckfehleranzahl a) Schüler b) Bolzen; Buchseiten a) Linz; blond b) 56 mm; 2 Fehler <p>Fehleranzahl; defekt – fehlerfrei</p> <p>Längen; Zeitdauer; Zugfestigkeit</p> <ul style="list-style-type: none"> a) alle Schüler Oberösterreichs b) alle Bolzen einer gewissen Art; eines bestimmten Romans a) alle Volksschüler b) eine Bolzen-Tagesproduktion; alle Seiten des 3. Kapitels a) Linz, Zwettl, Linz, Steyr ... b) 59 mm, 56 mm, 56 mm, 57 mm, ... b) 56 mm 57 mm 58 mm b) $x_i = 56 \text{ mm} \rightarrow n_i = 12$ <div style="border: 1px solid black; padding: 5px; width: fit-content; margin: 10px auto;"> $h(x_i) = h_i := \frac{n_i}{n} \cdot 100\%$ </div> <p style="text-align: center;">$0 \leq n_i \leq n \rightarrow 0 \leq \frac{n_i}{n} \leq 1$</p> <div style="border: 1px solid black; padding: 5px; width: fit-content; margin: 10px auto;"> $N_k := \sum_{i=1}^k n_i$ </div> <div style="border: 1px solid black; padding: 5px; width: fit-content; margin: 10px auto;"> $H(x_k) = H_k := \sum_{i=1}^k h_i \cdot 100\%$ </div> <p style="text-align: center;">$H_{k+1} = H_k + h_{k+1} \geq H_k$</p>

Bei sehr vielen verschiedenen, möglichen oder tatsächlichen Ausprägungen eines quantitativen (stetigen) Merkmals wird das Intervall reeller Zahlen, in dem alle Messwerte liegen, zweckmäßigerweise in Teilintervalle (so genannte *Klassen*) unterteilt, die Daten werden *klassiert*. Dazu einige

Bemerkungen:

- Da bei stetigen Merkmalen gemessen wird, ist von vornherein eine Klassierung durch die Messgenauigkeit gegeben. Diese ist hier aber nicht gemeint.
- Im Allgemeinen erfolgt eine Klassierung ab einem Stichprobenumfang von $n \geq 50$, darunter ist eine Klasseneinteilung unüblich. Dabei ist es gleichgültig, ob dieser Umfang tatsächlich vorliegt oder ob dieser Umfang möglich ist (Bestimmung des Geburtsgewichts von Säuglingen; dabei könnten die Ergebnisse beispielsweise zwischen 800g und 4500g liegen, bei einer Genauigkeit von 1g sind daher sehr viele Ausprägungen möglich).
- Klassenbildung schafft Übersichtlichkeit, gleichzeitig gehen aber Informationen verloren. Man wird im Einzelnen die Vor- und Nachteile gegeneinander abwägen.
- Die Anzahl der Klassen richtet sich nach dem (möglichen) Stichprobenumfang. In der Praxis üblich sind 7 bis 20 Klassen, als Faustregel gilt:

$$\text{Anzahl der Klassen} \approx \sqrt{n}$$

- Die Klassen sollten nach Möglichkeit gleich breit sein.

Empfohlene Vorgangsweise bei der Klassenbildung:

- Die Unterteilung muss so erfolgen, dass *jeder Messwert eindeutig einer Klasse zugeordnet* werden kann. Als Klassengrenzen wählt man daher zweckmäßigerweise Rundungsgrenzen der Messwerte, da diese „um eine Stelle genauer“ als die Messergebnisse selbst sind.

- Die *Klassenweite* w errechnet man aus den Rundungsgrenzen ($o(u)x_i \dots$ obere (untere) Rundungsgrenze von x_i) des größten bzw. kleinsten Messwerts nach den Formeln

$$\left. \begin{aligned} w &= \frac{o \cdot x_{\max} - u \cdot x_{\min}}{\sqrt{n}} && \text{für } 50 \leq n \leq 400 \\ w &= \frac{o \cdot x_{\max} - u \cdot x_{\min}}{20} && \text{für } n > 400 \end{aligned} \right\} \begin{array}{l} w \text{ auf die Genauigkeit} \\ \text{der Messwerte} \\ \text{runden} \end{array}$$

- *Anzahl k der Klassen:* Durch Runden von $\frac{o \cdot x_{\max} - u \cdot x_{\min}}{w}$ auf die *nächstgrößere* ganze Zahl erhält man die Anzahl der Klassen.

- *Klassengrenzen:* Der linke Rand der ersten und der rechte Rand der letzten Klasse werden so gewählt,
 - dass die Klassenränder Rundungsgrenzen sind und
 - dass das von ihnen gebildete Intervall der „Länge“ $k \cdot w$ das Messintervall möglichst symmetrisch umfasst.

- Als *Repräsentanten* der Klassen werden die Klassenmitten gewählt. Sie dienen – als Stellvertreter aller Messwerte innerhalb einer Klasse – zur näherungsweisen Berechnung der Stichprobenkennwerte (*praktische* Kennwerte).

Messwerte sind Zahlen begrenzter Genauigkeit: 3,2 repräsentiert das Intervall [3,15; 3,25[

↑ ↑
Rundungsgrenzen

Hinweis: In der Praxis wird auch „nichtmathematisch“ gerundet (z.B. Zeitangabe bei Digitaluhren)!

Beispiel:

$$x_{\max} = 531 \text{ g}, \quad x_{\min} = 422 \text{ g}, \quad n = 60$$

$$\rightarrow w = \frac{531,5 \text{ g} - 421,5 \text{ g}}{\sqrt{60}} \approx 14 \text{ g}$$

$$\rightarrow \frac{531,5 \text{ g} - 421,5 \text{ g}}{14 \text{ g}} \approx 7,9 \approx 8 = k$$

$$k \cdot w = 8 \cdot 14 \text{ g} = 112 \text{ g}$$

$$\text{Rand: } \frac{112 \text{ g} - (531 \text{ g} - 422 \text{ g})}{2} = 1,5 \text{ g}$$

$$\rightarrow \begin{cases} 1. \text{ Klasse: } [420,5 \text{ g}; 434,5 \text{ g}[\\ 8. \text{ Klasse: } [518,5 \text{ g}; 532,5 \text{ g}[\end{cases}$$

427,5g ist der Repräsentant der ersten Klasse (ist aber selbst gar kein Messwert).

Bemerkungen:

- Die Anzahl der Objekte einer Stichprobe, die innerhalb einer Klasse liegen, ergibt die absolute Häufigkeit dieser Klasse und damit das „Gewicht“ ihres Repräsentanten.
- Die Repräsentanten müssen selbst nicht unbedingt Messwerte sein.

Graphische Darstellungen: Die Besetzungszahlen alleine vermitteln ein recht unanschauliches Bild von der *Häufigkeitsverteilung* einer Stichprobe. Besser geeignet dafür sind graphische Darstellungen.

- *Koordinatensystem:* Es kann nur bei einem Merkmal verwendet werden, das entweder quantitativ ist oder dessen Ausprägungen gereiht werden können. Dabei trägt man die absoluten oder relativen Häufigkeiten in Abhängigkeit von den Merkmalsausprägungen bzw. den Klassenmitten auf, und zwar als

- Punkte (unüblich): *Punktendiagramm*; verbindet man die Punkte miteinander, so erhält man ein *Liniendiagramm* (*Häufigkeitspolygon*).

Achtung: Letzteres täuscht bei diskreten Merkmalen Steigtigkeit vor!

- senkrechte Striche oder Balken, deren Längen den Häufigkeiten entsprechen: *Stab-* oder *Balkendiagramm*.
Achtung: Bei Klassenbildung geht der Klasseindruck verloren, stetige Merkmale erscheinen diskret!

- aneinander gereichte Rechtecke, deren Flächen (bei gleichen Breiten der Rechtecke auch die Höhen) den Häufigkeiten entsprechen: *Histogramm* (*Staffelbild*); gut geeignet für klassierte Daten, die Rechtecksbreiten entsprechen dabei den Klassenweiten.

Bemerkung: Die Gesamtfläche entspricht dem Stichprobenumfang n bei absoluten bzw. 1 bei relativen Häufigkeiten.

- *Kreisdiagramm:* Diese Darstellungsform wird vorzugsweise für qualitative Merkmale verwendet. Die Kreisfläche wird entsprechend den Häufigkeiten in Sektoren unterteilt (neben den Flächen entsprechen auch die Winkel den Häufigkeiten).

- *Bildhafte Graphiken, Piktogramme:* Anhand von Bildern und Symbolen werden die Häufigkeiten qualitativer Merkmale meist durch die Größe (Länge, Fläche) und die Anzahl der verwendeten Objekte dargestellt.

Die Strichliste (zum vorigen Beispiel) hat z. B. folgendes Aussehen:

422 g | , 425 g ||| , 426 g | ,

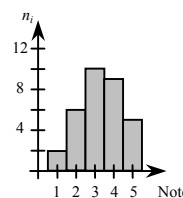
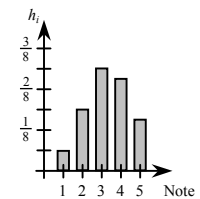
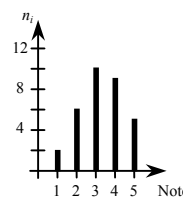
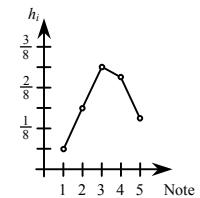
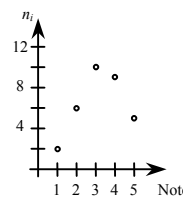
432 g | , 435 g || ,

Dann hat die 1. Klasse (und ihr Repräsentant) die absolute Häufigkeit $n_1 = 6$

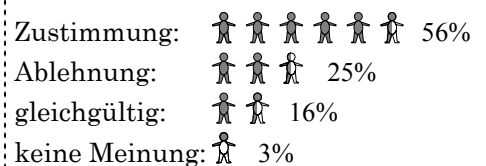
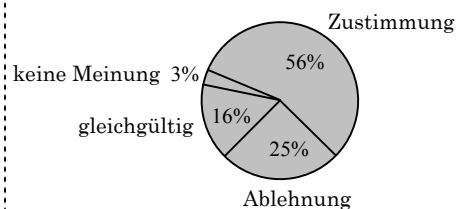
Beispiel:

Notenverteilung einer Schularbeit:

Note	1	2	3	4	5
Anzahl	2	6	10	9	5



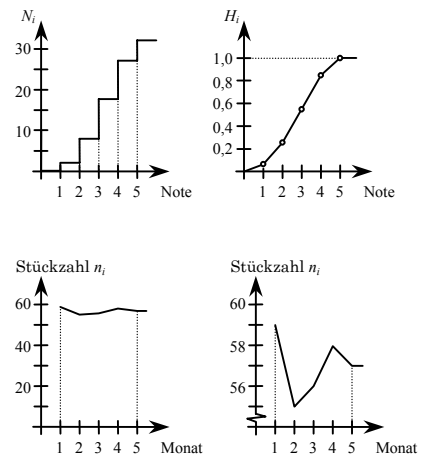
Ergebnis einer Meinungsumfrage:



Bemerkungen:

- Neben der Darstellung von absoluter bzw. relativer Häufigkeit ist auch die Veranschaulichung der Summenhäufigkeit gebräuchlich: *Summenhäufigkeitsverteilung* (üblicherweise als Liniendiagramm oder Histogramm ausgeführt); es dient für Fragestellungen der Art „wann wurden so und so viel Prozent erreicht“, „wie viele Stücke liegen unterhalb der Toleranzgrenze“ usw.
 Bemerkung: Bei klassierten Daten erfolgen die Sprünge an den oberen Klassengrenzen.
- Die Aussagekraft einer graphischen Darstellung hängt stark vom gewählten Maßstab und der Lage der Koordinatenachsen ab! Es ist also Vorsicht geboten, sonst entsteht ein falscher Eindruck.

Beispiel: Notenverteilung



Um gleichartige Stichproben rasch miteinander vergleichen zu können (z. B. wöchentliche Umfrageergebnisse oder Proben einer laufenden Produktion), bedient man sich einiger Kenngrößen:

- Lagekennwerte** (spiegeln die Häufigkeiten und die Größenordnung der Merkmalsausprägungen wider); wichtig ist dabei die Idee der *Mittelwerte*: Diese, nur für quantitative Merkmale erklärten Größen, stellen sozusagen ein „Kondensat“ der Messwerte dar und bilden bei entsprechender Anzahl einen für die Stichprobe charakteristischen Ersatz der (unterschiedlichen) Merkmalsausprägungen.
- Streuungskennwerte** (beschreiben die Abweichungen der Merkmalsausprägungen voneinander)

Die wichtigsten Kennwerte der Lage: ($k \dots$ Anzahl der unterscheidbaren Merkmalsausprägungen)	Beispiel: Absolventenanzahl der E- und M-Klassen der letzten 6 Jahre	
<p>a) Der (<i>arithmetische</i>) <i>Mittelwert</i> \bar{x} (<i>arithmetisches Mittel</i>) wird gebildet, indem man die Summe aller Merkmalsausprägungen x_i durch den Stichprobenumfang n teilt:</p> $\bar{x} := \frac{1}{n} \cdot \sum_{i=1}^n x_i = \frac{1}{n} \cdot \sum_{i=1}^k n_i \cdot x_i = \sum_{i=1}^k h_i \cdot x_i$ <p>Der Betrag des Mittelwerts kann (sollte) eine Dezimalstelle mehr als die Messwerte haben.</p>	<p>16, 17, 20, 15, 18, 18, 17, 18, 14, 17, 18, 13, 19, 18, 14</p>	<p>22, 23, 23, 22, 18, 24, 19, 21, 27, 22, 24, 24</p>
<p>b) Der <i>Median</i> \tilde{x} (<i>Zentralwert</i>) ist der mittlere Wert bzw. das arithmetische Mittel der beiden mittleren Werte einer der Größe nach geordneten Stichprobe mit dem Umfang n:</p> $\tilde{x} := x_{\frac{n+1}{2}} \quad \text{für } n \text{ ungerade}$ $\tilde{x} := \frac{1}{2} \cdot (x_{\frac{n}{2}} + x_{\frac{n}{2}+1}) \quad \text{für } n \text{ gerade}$	<p>13, 14, 14, 15, 16, 17, 17, 17, 18, 18, 18, 18, 18, 19, 20</p>	<p>18, 19, 21, 22, 22, 22, 23, 23, 24, 24, 24, 27</p>
<p>c) Gibt es genau eine häufigste Ausprägung – gleichgültig, ob ein qualitatives oder quantitatives Merkmal vorliegt – so bezeichnet man diese als <i>Modus</i> \hat{x} (<i>Modalwert</i>) der Stichprobe.</p>	<p>$\hat{E} = 18$</p>	<p>\hat{M} existiert nicht, weil die Zahlen 22 und 24 gleich häufig sind.</p>

Bemerkungen:

- „Schwerpunkteigenschaft“ des arithmetischen Mittelwerts: Die Summe aller Differenzen zwischen den Merkmalswerten und dem Mittelwert ist null:

$$\sum_{i=1}^n (x_i - \bar{x}) = 0$$

- Vorteilhafte Ermittlung des Mittelwerts (wenn „händisch“ gerechnet werden muss): Durch die Annahme eines vorläufigen Mittelwerts x_v (meist der Modalwert) lassen sich große Zahlen vermeiden; für den Mittelwert gilt dann:

$$\bar{x} = x_v + \frac{1}{n} \cdot \sum_{i=1}^n (x_i - x_v)$$

- Der Mittelwert wird durch *Ausreißer* stark beeinflusst. In diesen Fällen wird die Lage durch den Median meist besser beschrieben.

- Die Anwendung der Mittelwertbildung auf qualitative Merkmale, die gereiht werden können (wie Schulnoten oder Güteklassen), ist eigentlich nicht zulässig: Die Zuordnung von verbalen Merkmalsbeschreibungen und Zahlen ist willkürlich, die Wahl der Zahlenskala beeinflusst aber die Mittelwertbildung entscheidend. Darüber hinaus muss der Mittelwert wieder rückinterpretiert werden.

- Bei klassierten Daten erfolgt die Bestimmung der Lagekennwerte sehr oft näherungsweise:

- Für die Mittelwertberechnung werden die Repräsentanten genommen,
- als Median dient der Wert, der die Gesamtfläche des Histogramms für die Häufigkeiten halbiert (der Median teilt ja jede (geordnete) Stichprobe in zwei gleich große Teilmengen) und
- als Modalwert nimmt man meistens den Repräsentanten jener Klasse, welche die größte Häufigkeit aufweist.

- Das arithmetische Mittel ist aber nicht immer zur Charakterisierung der Lage geeignet, wie nebenstehendes Beispiel jährlicher Umsatzsteigerungen zeigt (Umsatzsteigerungen werden immer auf das jeweils vergangene Jahr bezogen). Deshalb kommen auch andere Mittelwerte, wenn auch nicht sehr oft, zur Anwendung.

$$\sum_{i=1}^n (x_i - \bar{x}) = \sum_{i=1}^n x_i - \sum_{i=1}^n \bar{x} = n \cdot \bar{x} - \bar{x} \cdot \underbrace{\sum_{i=1}^n 1}_n$$

$$x_v + \frac{1}{n} \sum_{i=1}^n (x_i - x_v) = x_v + \frac{1}{n} \sum_{i=1}^n x_i - \frac{1}{n} \sum_{i=1}^n x_v$$

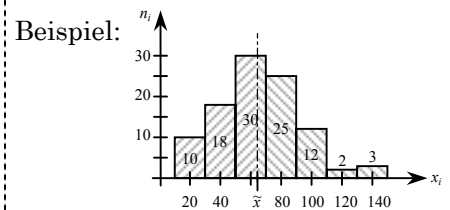
z. B. $x_i = 8; 11; 11; 15; 15; 15; 15; 17; 19$
 $\rightarrow \bar{x} = 15 + \frac{1}{9} \cdot (-7 + 2 \cdot (-4) + 2 + 4) = 14$

Ausreißer: einzelner Wert, der stark von den anderen Werten abweicht.

z. B. $x_i = 1; 21; 21; 22; 22; 22; 23; 23; 24$
 $\rightarrow \bar{x} = \frac{1}{9} \cdot (1 + 2 \cdot 21 + \dots + 24) \approx 19,9$

statt $\bar{x} \approx 22,3$ Median: $\tilde{x} = 22$

Die Beurteilungen „Nicht genügend – Gut – Gut“ würden bei der üblichen Notenskala ein „Befriedigend“ als Mischnote ergeben, nach der Skala: „Sehr gut“ (10) – „Gut“ (6) – „Befriedigend“ (5) – „Genügend“ (4) – „Nicht genügend“ (0) aber nur ein „Genügend“.



$$\bar{x} \approx \frac{10 \cdot 20 + 18 \cdot 40 + \dots + 3 \cdot 140}{100} = 65,8$$

$$\tilde{x} \approx 50 + 20 \cdot \frac{50 - (10 + 18)}{30} \approx 64,7$$

$$\hat{x} \approx 60$$

Beispiel: Der Umsatz erhöhte sich in drei aufeinander folgenden Jahren um 10%, 20% bzw. 30%, insgesamt um 71,6% ($1,1 \cdot 1,2 \cdot 1,3 - 1 = 0,716$).

Würde man mit einer mittleren Steigerung von 20% rechnen, so ergäbe sich ein Zuwachs von 72,8%.

Die wichtigsten Kennwerte der Streuung von Messwerten:

(n ... Stichprobenumfang, k ... Anzahl der unterscheidbaren Merkmalsausprägungen)

- a) Die *Spannweite* R ist die Differenz zwischen der größten und der kleinsten Merkmalsausprägung:

$$R := x_{\max} - x_{\min}$$

- b) Die *mittlere lineare Abweichung* \bar{s} ist der Mittelwert der Differenzbeträge aller Merkmalswerte x_i vom Mittelwert \bar{x} :

$$\bar{s} := \frac{1}{n} \cdot \sum_{i=1}^n |x_i - \bar{x}| = \sum_{i=1}^k h_i \cdot |x_i - \bar{x}|$$

Bemerkung: Die Beträge sind wegen der Schwerpunkteigenschaft des arithmetischen Mittels (Seite 5 oben) nötig.

- c) Die (*empirische*) *Varianz* s^2 (mittlere quadratische Abweichung) ist der Mittelwert der Abweichungsquadrate aller Merkmalswerte x_i vom Mittelwert \bar{x} :

$$s^2 := \frac{1}{n} \cdot \sum_{i=1}^n (x_i - \bar{x})^2 = \sum_{i=1}^k h_i \cdot (x_i - \bar{x})^2$$

Soll die Varianz der Stichprobe gleichzeitig als *Schätzwert* für die Varianz der Grundgesamtheit dienen, so benutzt man besser

$$s^2 := \frac{1}{n-1} \cdot \sum_{i=1}^n (x_i - \bar{x})^2 = \frac{n}{n-1} \cdot \sum_{i=1}^k h_i \cdot (x_i - \bar{x})^2$$

- d) Die (*empirische*) *Standardabweichung* s ist die Wurzel aus der Varianz (und hat damit die gleiche Einheit wie die Merkmalswerte und der Mittelwert; sie ist deshalb mit diesen vergleichbar):

$$s := \sqrt{s^2}$$

- e) Der *Variationskoeffizient* V (*relative Streuung*) ist der Quotient aus der Standardabweichung und dem Mittelwert:

$$V := \frac{s}{\bar{x}}$$

Der Variationskoeffizient relativiert die Streuung: Stichproben mit gleicher Streuung aber unterschiedlichen Mittelwerten können hinsichtlich des Streuverhaltens der Merkmalswerte besser miteinander verglichen werden.

Beispiel: Vergleich zweier Stichproben mit dem gleichen Mittelwert

$$x_i = 3; 5; 3; 4; 4; 8 \quad y_i = 6; 1; 3; 3; 5; 9$$

$$R = 8 - 3 = 5 \quad R = 8$$

$$\bar{x} = \frac{27}{6} = 4,5 \quad \bar{y} = 4,5$$

$$\bar{s}_x = \frac{8}{6} \approx 1,3 \quad \bar{s}_y \approx 2,2$$

Der Grund für die Verwendung der Quadrate an Stelle der Beträge liegt in der bequemerer Handhabung bei Abschätzungen und Beweisen.

$$s_x^2 = \frac{17,5}{6} \approx 2,9 \quad s_y^2 \approx 6,6$$

Der Grund für die Division durch den um eins kleineren Stichprobenumfang liegt in der Forderung nach Erwartungstreue eines Schätzwerts.

$$s_x^2 = \frac{17,5}{5} = 3,5 \quad s_y^2 = 7,9$$

Für die Grundgesamtheit ist daher die Standardabweichung das quadratische Mittel der Abweichungen aller Merkmalswerte vom arithmetischen Mittelwert.

$$s_x = \sqrt{\frac{17,5}{6}} \approx 1,7 \quad s_y \approx 2,6 \quad \text{bzw.}$$

$$\text{bzw. } s_x \approx 1,9 \quad s_y \approx 2,8$$

$$V_x \approx \frac{1,7}{4,5} \approx 0,38 \quad V_y \approx 0,57 \quad \text{bzw.}$$

$$\text{bzw. } V_x \approx 0,42 \quad V_y \approx 0,62$$

Bemerkung: Die Varianz und die Standardabweichung sind die in der Praxis am häufigsten verwendeten Streumaße.

Einige wichtige Eigenschaften der Varianz bzw. der Standardabweichung:

(n ... Stichprobenumfang)

- a) Die Varianz ist die Differenz des Mittelwerts der quadrierten Daten und dem Quadrat des Mittelwerts.

$$s^2 = \frac{1}{n} \cdot \sum_{i=1}^n x_i^2 - \bar{x}^2 = (\bar{x}^q)^2 - \bar{x}^2$$

- b) Minimaleigenschaft der Varianz: Die Varianz ist die kleinste quadratische Abweichung der Stichprobenwerte, d. h. es gilt für alle Werte $a \neq \bar{x}$:

$$s^2 < \frac{1}{n} \cdot \sum_{i=1}^n (x_i - a)^2$$

- c) Ungleichung von ČEBYŠEV

Für alle $k \in \mathbb{R}$ gilt: Im Intervall $]\bar{x} - k \cdot s ; \bar{x} + k \cdot s[$ liegen mindestens $n \cdot (1 - \frac{1}{k^2})$ Merkmalswerte.

Bemerkung: Diese Ungleichung erlaubt es, eine Abschätzung bezüglich der Streuung und damit der *Messungenauigkeit* für eine bestimmte Anzahl von Messwerten (z. B. 90% der Stichprobe) zu machen. Die Schätzung fällt im Allgemeinen recht „pessimistisch“ aus (der Messfehler wird viel größer vermutet als er tatsächlich ist), bei bestimmten Häufigkeitsverteilungen (z. B. Normalverteilung) lassen sich die Messungenauigkeiten enger eingrenzen.

$$\begin{aligned} \sum_{i=1}^n (x_i - \bar{x})^2 &= \sum_{i=1}^n x_i^2 - 2\bar{x} \sum_{i=1}^n x_i + \sum_{i=1}^n \bar{x}^2 = \\ &= \sum_{i=1}^n x_i^2 - 2n \cdot \bar{x}^2 + n \cdot \bar{x}^2 \end{aligned}$$

$$\begin{aligned} \sum_{i=1}^n (x_i - (\bar{x} \pm \delta))^2 &= \sum_{i=1}^n ((x_i - \bar{x}) \mp \delta)^2 = \\ \sum_{i=1}^n (x_i - \bar{x})^2 \mp 2\delta \cdot \sum_{i=1}^n (x_i - \bar{x}) + \sum_{i=1}^n \delta^2 & \\ \underbrace{\hspace{1.5cm}}_{n \cdot s^2} \quad \underbrace{\hspace{1.5cm}}_0 \quad \underbrace{\hspace{1.5cm}}_{> 0} & \end{aligned}$$

Sei n^* die Anzahl aller Messwerte x_i^* , die in $]\bar{x} - k \cdot s ; \bar{x} + k \cdot s[$ liegen und $n - n^*$ die Anzahl der restlichen Werte x_j^{**} ; dann gilt:

$$\begin{aligned} \sum_{i=1}^n (x_i - \bar{x})^2 &= \sum_{i=1}^{n^*} (x_i^* - \bar{x})^2 + \sum_{j=1}^{n-n^*} (x_j^{**} - \bar{x})^2 \\ \underbrace{\hspace{1.5cm}}_{n \cdot s^2} \quad \underbrace{\hspace{1.5cm}}_{\geq 0} \quad \underbrace{\hspace{1.5cm}}_{\geq (k \cdot s)^2} & \\ \rightarrow n \cdot s^2 &\geq (n - n^*) \cdot k^2 \cdot s^2 \end{aligned}$$