

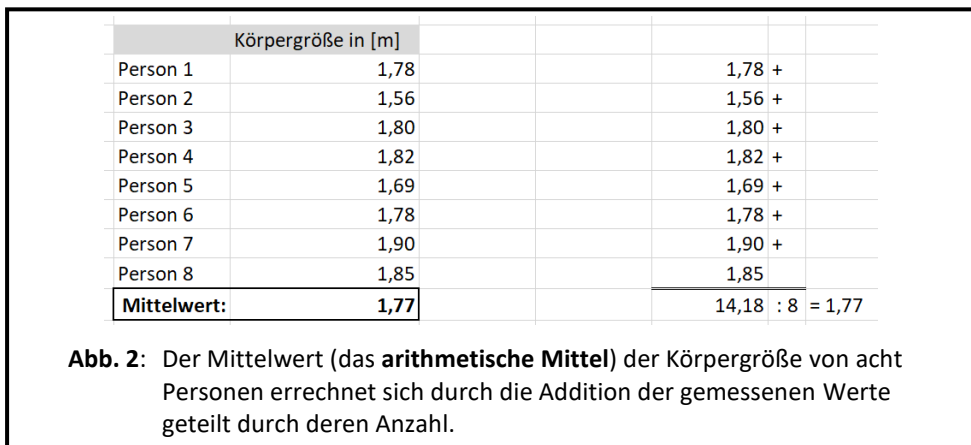
3. Deskriptive Statistik

„Die deskriptive (auch: beschreibende) Statistik hat zum Ziel, [...] Daten durch Tabellen, Kennzahlen [...] und Grafiken übersichtlich darzustellen und zu ordnen. Dies ist vor allem bei umfangreichem Datenmaterial sinnvoll, da dieses nicht leicht überblickt werden kann.“ – Wikipedia.

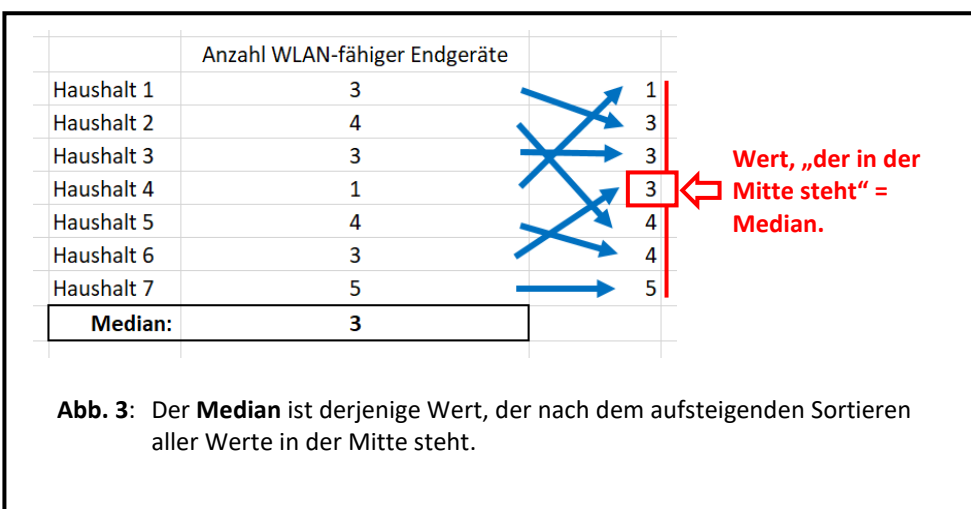
3.1 Wichtige Größen

3.1.1 Lagemaße

Häufig möchte man einen **Parameter** bei **zwei Gruppen** vergleichen, **beispielsweise die Körpergröße von Männern und Frauen**. Eine beliebte Methode hierfür ist die Bildung des (sicher jedem bekannten) **arithmetischen Mittelwerts**: Man addiert die Werte aller Messungen und teilt sie durch die Anzahl der Messungen (s. **Abb. 2**). Der arithmetische Mittelwert darf nur bei **intervallskalierten Daten** angewendet werden. Einfach ausgedrückt sind das Daten, bei denen ein z.B. doppelt so hoher Wert auch tatsächlich bedeutet, dass die Ausprägung des gemessenen Parameters doppelt so groß ist (vgl. mit **ordinal skalierten** Daten).



Weniger bekannt dürfte der **Median** sein: Beim Median werden alle vorhandenen Messwerte aufsteigend geordnet. Der Wert, der genau in der Mitte steht, ist der Median (s. **Abb. 3**). Falls eine gerade Anzahl an Werten vorliegt, berechnet sich der Median aus dem Mittelwert der beiden mittleren Werte.



Der Median kann auch bei **ordinal skalierten Daten** angewendet werden. Bei diesen Daten geht es nur darum, eine Reihenfolge bzw. Rangfolge festzustellen: Bei einem Autorennen kann ein **erster, zweiter und dritter Platz** vergeben werden. Dabei spielt es keine Rolle, ob der Erstplatzierte eine Stunde oder eine Minute schneller war als der Zweitplatzierte. Häufig kommen ordinal skalierte Daten auch bei Umfragen zum Einsatz:

Frage: **Wie häufig nutzen Sie Ihr Handy, um Ihren Kontostand abzurufen?**
 Antwortmöglichkeiten: **sehr oft (1), oft (2), gelegentlich (3), selten (4), nie (5).**

Über die Abstände zwischen den vorgegebenen Antwortmöglichkeiten lässt sich kaum eine Aussage machen: Ob jemand der „gelegentlich“ seinen Kontostand per Handy checkt, das doppelt so oft oder viermal so oft tut, wie jemand der „selten“ seinen Kontostand am Handy checkt, bleibt unklar. Auch **Schulnoten** sind ordinal skalierte Daten. Man kann nämlich z.B. nicht sagen, dass jemand der die Note 3 erhalten hat, dreimal so viele Fehler gemacht hat, wie jemand der die Note 1 erhalten hat. Insofern ist die Berechnung von arithmetischen Mittelwerten bei Schulnoten aus statistischer Sicht eigentlich unzulässig (s. **Abb. 4**).

	Schüler 1	Schüler 2	Schüler 3
Note 1	2	2	2
Note 2	3	2	3
Note 3	3	2	2
Note 4	4	4	4
Note 5	4	4	4
Note 6	2	2	2
arith. Mittelwert	3,0	2,7	2,8
Median	3,0	2,0	2,5

Abb. 4: Vergleich von **Median** und **arithmetischem Mittelwert** bei Schulnoten verschiedener Schüler.

Der Median ist sehr **robust** gegenüber **Ausreißern**. Das sind einzelne Werte, die sich deutlich von den anderen Messwerten einer Gruppe unterscheiden. Daher ist der Median in der Biologie oft beliebt, da gerade beim Umgang mit Lebewesen immer wieder Besonderheiten auftreten (s. **Abb.5**).

Häufigkeit des Gefiederputzens bei Vögeln pro Stunde		Häufigkeit des Gefiederputzens bei Vögeln pro Stunde	
Vogel 1	3	Vogel 1	3
Vogel 2	4	Vogel 2	4
Vogel 3	3	Vogel 3	3
Vogel 4	1	Vogel 4	1
Vogel 5	4	Vogel 5	4
Vogel 6	3	Vogel 6	3
Vogel 7	4	Ausreißer: Vogel 7	75
arith. Mittelwert	3,14	arith. Mittelwert	13,29
Median	3,00	Median	3,00

Abb. 5: Einzelne **Ausreißer** haben oft große Auswirkungen auf den Mittelwert, nicht jedoch auf den Median.

3.1.2 Streuungsmaße

Der Mittelwert oder Median alleine sagt noch wenig aus. Manchmal liegen die gemessenen Werte nämlich eng am tatsächlichen Mittelwert, manchmal **streuen** sie aber auch stark (s. **Abb. 6**).

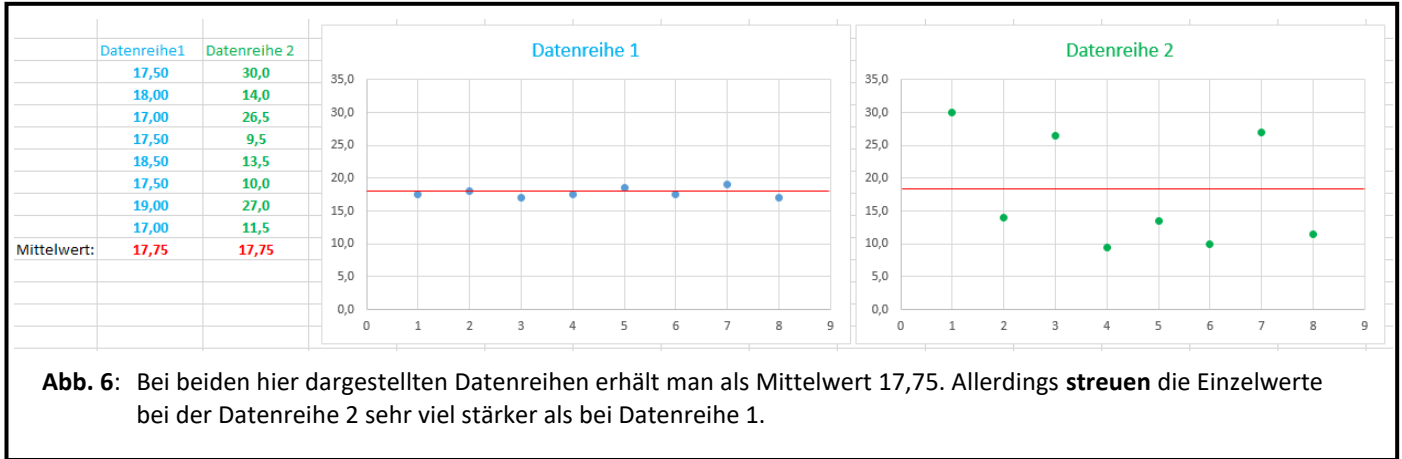


Abb. 6: Bei beiden hier dargestellten Datenreihen erhält man als Mittelwert 17,75. Allerdings **streuen** die Einzelwerte bei der Datenreihe 2 sehr viel stärker als bei Datenreihe 1.

Es gibt verschiedene Maße, die in der Lage sind diese Streuung in Zahlen auszudrücken: Zum Beispiel die **Varianz**, die vereinfacht ausgedrückt ungefähr den Mittelwert der Abweichungsquadrate darstellt. Aber auch die (empirische) **Standardabweichung**, die als Wurzel aus der Varianz definiert ist. Beide Parameter kann man in Tabellenkalkulationsprogrammen in der Regel einfach berechnen lassen (s. **Abb. 7**).

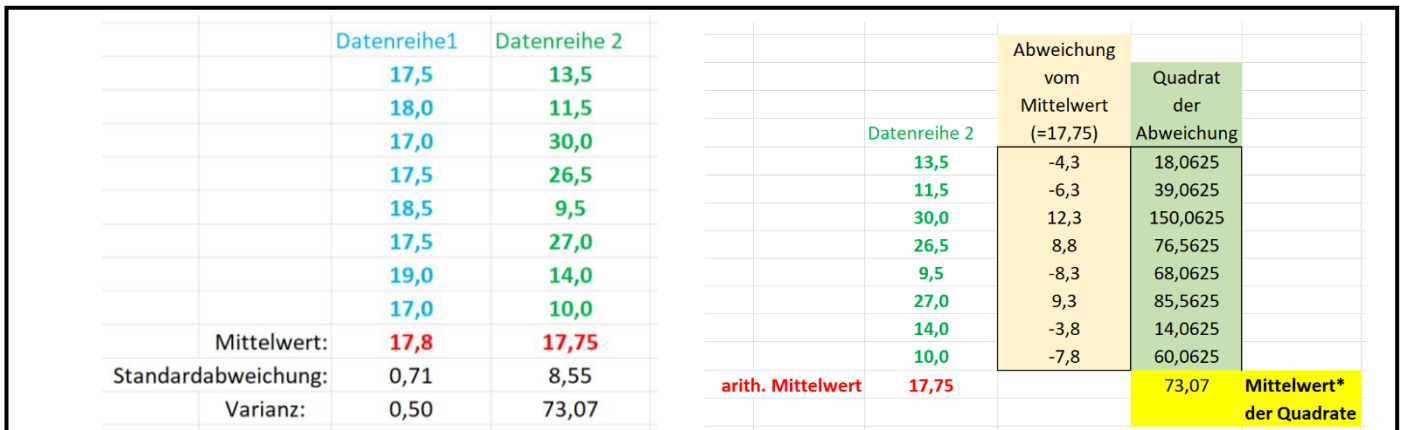
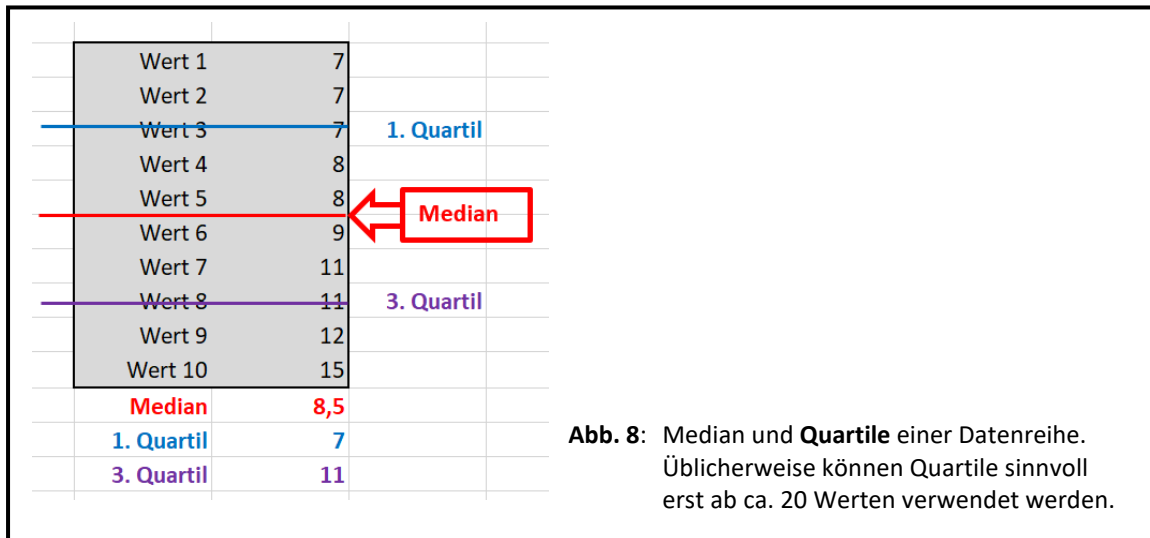


Abb. 7: links: Die Varianz und die Standardabweichung kann man von Tabellenkalkulationsprogrammen einfach über eine Formel in der entsprechenden Zelle berechnen lassen.

rechts: Veranschaulichung des mathematischen Verfahrens zur Berechnung der Varianz: Von jedem Messwert wird die Abweichung (Differenz) zum Mittelwert bestimmt. Dieser Wert wird quadriert. Der Mittelwert der Quadrate entspricht der Varianz.

* Wer genau nachrechnet, erhält für den Mittelwert hier eigentlich 63,94. Dies liegt daran, dass man bei der Berechnung dieses speziellen Mittelwerts die Summe nicht durch die Anzahl der Messungen teilt, sondern durch die um 1 verringerte Zahl (hier also: Summe geteilt durch 7, obwohl 8 Messungen vorliegen).

Auch bei ordinal skalierten Daten, kann man neben dem Median ein Streuungsmaß angeben: Den Interquartilsabstand. Die Quartile sind vereinfacht ausgedrückt die Werte unterhalb derer 25% (1. Quartil) bzw. 75% (3. Quartil) aller anderen Messwerte liegen.

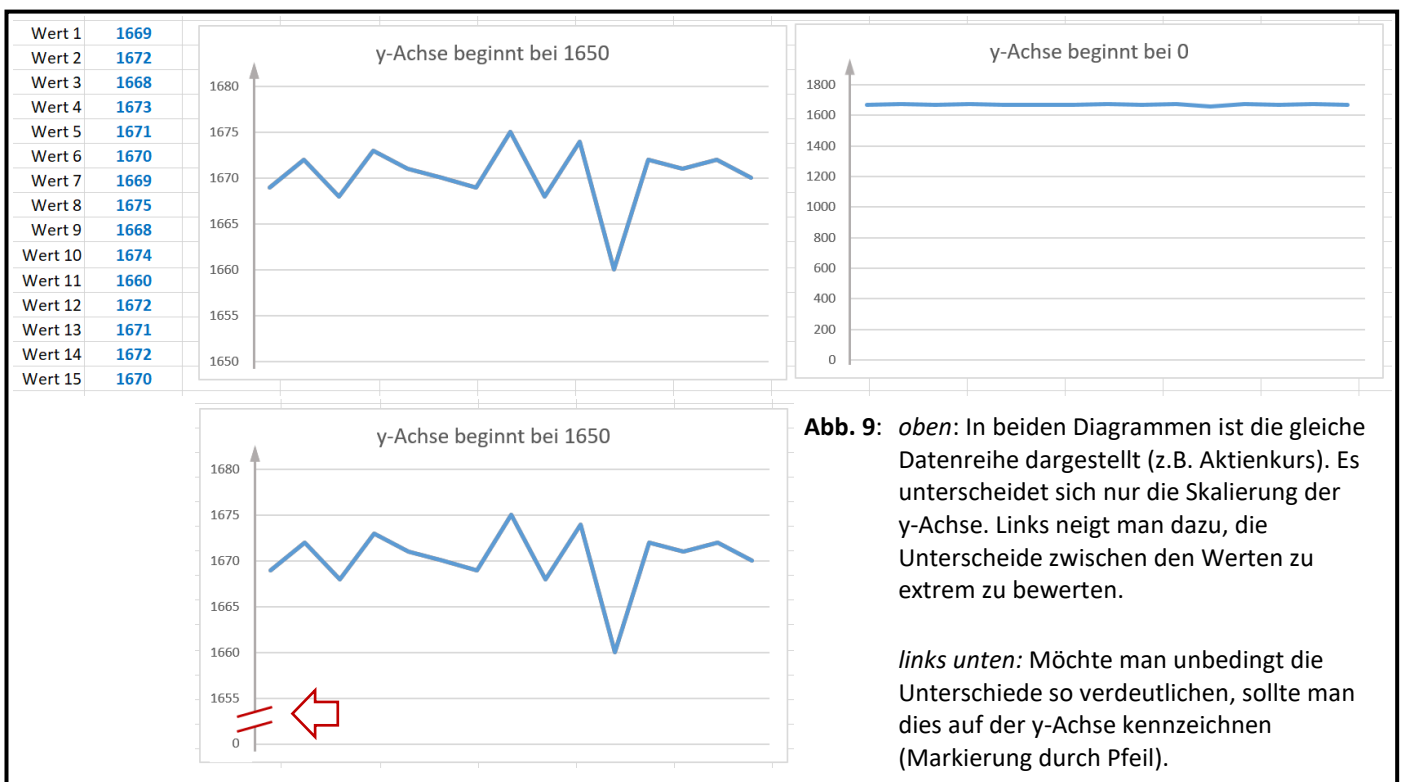


3.2 Diagrammtypen

Es gibt Standard-Diagrammtypen, die mit Tabellenkalkulationsprogrammen in der Regel leicht erstellt werden können. Jeder Diagrammtyp hat Vor- und Nachteile, die gegeneinander abgewogen werden müssen, eine detaillierte Betrachtung erfolgt hier jedoch nicht. Lediglich auf zwei typische Fehlerquellen soll hier hingewiesen werden.

Unterbrochene Skalierung der y-Achse

Gelegentlich beginnt in Diagrammen die y-Achse im Ursprung nicht mit dem Wert 0. Diese Darstellung wird dann gewählt, wenn man eigentlich kleine Unterschiede deutlich machen möchte. Es ist zwar nicht wirklich ein Fehler, sollte aber klar gekennzeichnet werden, z.B. durch eine Unterbrechung in der y-Achse (s. **Abb. 9**).



Falscher Diagrammtyp

Liniendiagramme bei denen Datenpunkte **verbunden** sind, dürfen nur verwendet werden, wenn zwischen den Datenpunkten auch tatsächlich eine **Verbindung** besteht, z.B. wenn sie zeitlich aufeinander folgen. Werden unterschiedliche Gruppen betrachtet, sind Säulendiagramme o.ä. zu verwenden.

Kreis- bzw. Tortendiagramme können nur verwendet werden, wenn die dargestellten Werte als Summe tatsächlich eine 100%-Gesamtheit ergeben.