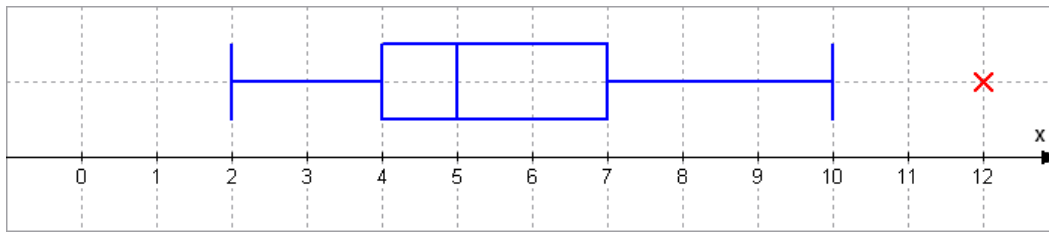


Beschreibende Statistik – Kenngrößen in der Übersicht (Ac)



Im folgenden wird die Berechnungsweise des **TI 83** (sowie von SPSS, s. unten) verwendet. Diese geht auf eine Festlegung von Moore und McCabe (2002) zurück. In der Literatur existieren insbesondere für die Berechnung der Quartile Q1 und Q3 (s.u.) noch ganz andere Formeln. Zum Beispiel rechnen EXCEL und das NRW-Schulportal Learn-line nach einer Methode von John W. Tukey (1983; Begründer der Explorativen Datenanalyse)!

Gegeben sei eine Liste $x(i)$ von n Daten. Eine ungeordnete Liste muss zunächst sortiert werden. U.a. werden jeweils 5 Kennzahlen ermittelt: x_{\min} , x_{\max} , med, Q1, Q3 (siehe Erläuterung unten).

Beispiel 1 mit $n=9$:

1	2		3	3	3	4	4		5	6
$x_{\min}=1$			Q1=2,5		med=3		Q3=4,5			$x_{\max}=6$

Beispiel 2 mit $n=10$:

0	1	2	4	4		6	7	8	9	9
$x_{\min}=0$		Q1=2		med=5		Q3=8				$x_{\max}=9$

Ersichtlich gibt es Unterschiede in der Behandlung bei geraden und ungeraden n !

1) Formeln zur Berechnung der Kenngrößen (mit dem TI 83):

Die Spannweite (Range): Differenz zwischen Minimum und Maximum der Liste

$$\text{range} = x_{\max} - x_{\min} \qquad \text{Im Beispiel 1: } \text{range} = 6 - 1 = 5$$

Bei den folgenden 3 Kenngrößen kann es vorkommen, dass der Listenindex zwischen 2 ganzen Zahlen liegt, z.B. $x(6,5)$. Es muss dann das arithm. Mittel verwendet werden. $x(6,5) \Rightarrow \frac{x(6) + x(7)}{2}$.

Der Zentralwert (Median oder med): Der Wert in der Mitte der sortierten Liste

$$\text{med} = x\left(\frac{n+1}{2}\right) \qquad \text{Im Beispiel 1: } \text{med} = x(5) = 3$$

Das erste (untere) Quartil Q1: Q1 gibt den oberen Bereich des ersten Viertels der Liste an.
 Andere Definition: Q1 ist der Median der links von **med** liegenden Liste.

$$n \text{ ungerade: } Q1 = x\left(\frac{n+1}{4}\right) \qquad n \text{ gerade: } Q1 = x\left(\frac{n+2}{4}\right)$$

$$\text{Im Beispiel 1: } Q1 = x(2,5) = \frac{2+3}{2} = 2,5$$

Das dritte (obere) Quartil Q3: Q3 gibt den oberen Bereich des dritten Viertels der Liste an.
 Andere Definition: Q3 ist der Median der rechts von **med** liegenden Liste:

$$n \text{ ungerade: } Q3 = x\left(\frac{3n+3}{4}\right) \qquad n \text{ gerade: } Q3 = x\left(\frac{3n+2}{4}\right)$$

$$\text{Im Beispiel 1: } Q3 = x(7,5) = \frac{4+5}{2} = 4,5$$

Der Interquartilsabstand IQR (interquartile range) :

$$IQR = Q3 - Q1$$

$$\text{Im Beispiel 1: } IQR = 4,5 - 2,5 = 2$$

Wichtig: Im IQR (zwischen Q1 und Q3) liegt genau die Hälfte aller Daten ! Begründung ?

Ausreißer: Ein Wert, der mehr als das 1,5-fache des IQR von den Quartilen abweicht.

Wie findet man Ausreißer ? Man definiert ein Intervall $[z_u ; z_o] = [Q1 - 1,5 \cdot IQR ; Q3 + 1,5 \cdot IQR]$
 Liegt ein Wert der Liste außerhalb dieses Bereichs, so ist er ein Ausreißer.
 Im Beispiel 1 ist das Intervall $[-0,5 ; 7,5]$. Es gibt dort also keine Ausreißer .

Die Quantile: Der Begriff Quantil ist ein Oberbegriff bzgl. Quartil und Median .

Quantile = Punkte einer nach Rang oder Größe geordneten Datenliste.
 Z.B. gibt das 0,35-Quantil die Obergrenze für 35% der unteren geordneten Liste an .

Beispiel für n=8 : 3 4 5 5 6 7 8 9
 x(1) x(2) x(3) x(4) x(5) x(6) x(7) x(8)

Das 0,2-Quantil ist x(2)=4 . Das 0,5-Quantil(Median) ist x(4,5)=5,5 .

Mögliche Formel (von Ac): Das p-Quantil (0<p<1) besitzt den Wert
 $\frac{x(\lfloor (n+1) \cdot p \rfloor) + x(\lfloor (n+1) \cdot p \rfloor + 1)}{2}$, falls $\text{frac}((n+1)p) = 0,5$ $\lfloor \rfloor$ = „Gaußklammerfunktion“
 $x(\text{round}((n+1) \cdot p))$, falls $\text{frac}((n+1)p) \neq 0,5$
 Diese Formeln sind nicht immer (aber häufig) kompatibel zur Quartilsdefinition (Q1,Q3) des TI 83 !

Markus Paul gibt für das p-Quantil folgende Berechnungen an (vermutlich nach Tukey bzw: EXCEL):

$$\frac{x(n \cdot p) + x(n \cdot p + 1)}{2} \qquad , \text{ falls } n \cdot p \text{ ganzzahlig}$$

$$x(\lfloor n \cdot p + 1 \rfloor) \qquad , \text{ falls } n \cdot p \text{ nicht ganzzahlig}$$

Die Perzentile:

Spezialfall der Quantile. Punkte, welche die Obergrenze für die Hundertstel (q%) der unteren geordneten Liste angeben.

Arithmetischer Mittelwert (mean):

$$\bar{x} = \frac{x(1) + x(2) + \dots + x(n)}{n}$$

Standardabweichung (standard deviation) und Varianz:

Erst wird die Varianz $V(x)$ als mittlere Abweichung der Quadrate vom Mittelwert gebildet.

$$V(x) = \frac{(x(1) - \bar{x})^2 + (x(2) - \bar{x})^2 + \dots + (x(n) - \bar{x})^2}{n}$$

Dann ist die Standardabweichung σ_n (σ_x beim TI 83) die Wurzel aus der Varianz:

$$\sigma_n = \sqrt{\frac{(x(1) - \bar{x})^2 + (x(2) - \bar{x})^2 + \dots + (x(n) - \bar{x})^2}{n}}$$

Achtung:

In der Praxis verwendet man bei Daten eher ein anderes Sigma, nämlich σ_{n-1} (S_x beim TI83).

Bei diesem wird in der Wurzel durch (n-1) statt durch n dividiert !

Dennoch hat auch σ_n z.B. bei theor. Verteilungen seine Berechtigung .

Vorteile des Medians gegenüber dem arith.Mittel sowie des IQRs geg. der Standardabweichung:

Median und IQR sind unempfindlich gegenüber Ausreißern und unzuverlässigen Messungen oder Übertragungsfehlern, weil sie **keine Gewichtung** der Daten vornehmen !!

2) Grafische Darstellung von Datenreihen:

A) Boxplot (Box-Whisker-Plot):

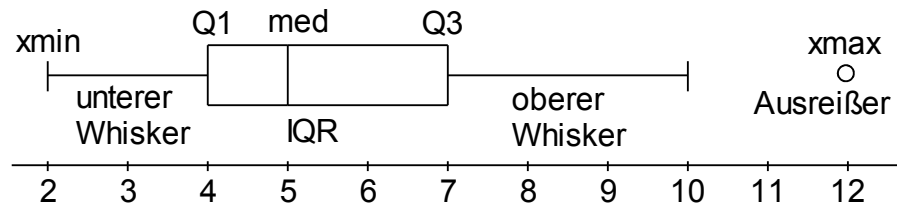
Der Boxplot stellt die Kenngrößen mittels einer Box dar.

Willkürliches Beispiel:

2 3 4 4 5 5 6 6 7 10 12

Es gelten dann: $x_{\min}=2$ $x_{\max}=12$ $med=5$ $Q1=4$ $Q3=7$ $IQR=3$

(„modifizierter Boxplot“)



Whisker (Schnurrbarthaare) sind beim **normalen Boxplot** die Verbindungslinien von Q1 zu x_{\min} sowie von Q3 zu x_{\max} .

Beim **modifizierten Boxplot** (siehe Grafik oben) kann es aber vorkommen, dass die Whisker einen der Randpunkte oder gar beide Randpunkte (x_{\min} , x_{\max}) nicht erreichen, weil Ausreißer immer außerhalb des Bereichs der Whisker gezeichnet werden.

Genauer:

In obiger Grafik gilt $IQR = 3$ und somit $1,5 \cdot IQR = 4,5$.

Das für den Ausschluss von Ausreißern zu betrachtende Intervall ist $[Q1 - 1,5 \cdot IQR ; 7 + 1,5 \cdot IQR]$.

Setzt man die entsprechenden Zahlen ein, so erhält man das Intervall $[-0,5 ; 11,5]$.

$x_{\min} = 2$ liegt innerhalb dieses Intervalls, ist also kein Ausreißer.

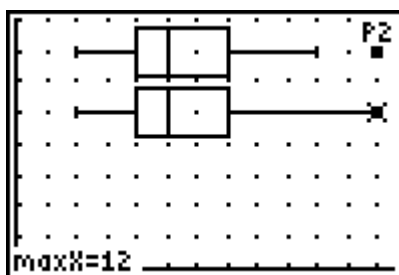
$x_{\max} = 12$ liegt außerhalb dieses Intervalls und ist demnach ein Ausreißer! Aus diesem Grunde schließt der obere Whisker die Zahl 12 nicht ein, sondern er erstreckt sich von Q3 bis zum letzten Wert, der noch innerhalb des Intervalls $[-0,5 ; 11,5]$ liegt. In diesem Fall ist das die Zahl 10.

Whisker und IQR sind Bereiche, keine Punkte. Im IQR liegen 50% aller Daten, im Bereich zwischen x_{\min} und Q1 sowie Q3 und x_{\max} liegen nochmals je 25% aller Daten.

Genauso wie bei den Quartilen gibt es für die Definition der Lage der Whisker in der Literatur verschiedene Möglichkeiten.

Darstellung mit dem TI83.

Modifizierter Boxplot (mit Ausreißer) und normaler Boxplot im Vergleich:



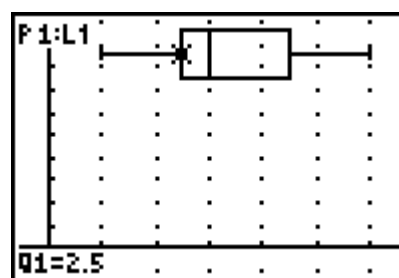
Boxplots mit dem TI 83

Zunächst nochmal das einführende Beispiel : 1 2 3 3 3 4 4 5 6

L1	L2	L3	1
1 2 3 3 3 4 4 5 6	-----	-----	1
L1(?) = 1			

```

2003 Plot2 Plot3
Off Off
Type: [L1] [L2] [L3]
      [Box] [Line] [Bar]
      [Normal] [Dot] [Line]
Xlist: L1
Freq: 1
    
```



Dies ist ein **normaler Boxplot** (Type 5) . Mögliche Ausreißer würden hier nicht gezeichnet !
Die Kennzahlen können mit **TRACE** abgefragt werden (siehe Bild 3 oben).

Alternativ können sie auch mittels **STAT CALC 1-Var Stats** ausgegeben werden (Bilder unten).

```

EDIT [2nd][DEL] TESTS
1: 1-Var Stats
2: 2-Var Stats
3: Med-Med
4: LinReg(ax+b)
5: QuadReg
6: CubicReg
7: QuartReg
    
```

```

1-Var Stats
x̄=3.444444444
Σx=31
Σx²=125
Sx=1.509230856
σx=1.422916497
n=9
    
```

```

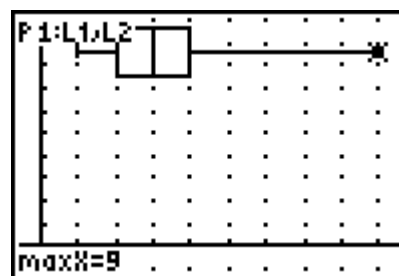
1-Var Stats
n=9
minX=1
Q1=2.5
Med=3
Q3=4.5
maxX=6
    
```

Ein weiteres Beispiel, diesmal werden absolute Häufigkeiten in L₂ mitverwendet . Im Statplot bei Freq den Wert L₂ (statt 1) eintragen ! Zuerst der normale Boxplot:

L1	L2	L3	2
1 2 3 3 3 4 4 5 9	2 1 2 2 1 1 1 1	-----	2
L2(?) =			

```

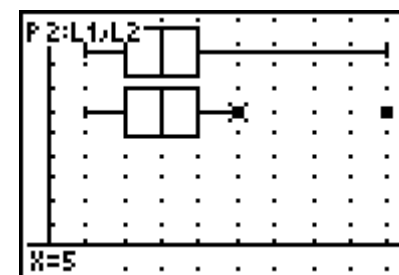
2003 Plot2 Plot3
Off Off
Type: [L1] [L2] [L3]
      [Box] [Line] [Bar]
      [Normal] [Dot] [Line]
Xlist: L1
Freq: L2
    
```



Verwendet man zusätzlich den modifizierten Boxplot (Type 4) , so können Ausreißer angezeigt werden. Hier ist das der Wert 9 . Der Whisker geht dann rechts nur noch bis zu x=5. Dies ist der letzte Wert der um den Ausreißer gekürzten Liste.

```

Plot1 Plot2 Plot3
Off Off
Type: [L1] [L2] [L3]
      [Box] [Line] [Bar]
      [Normal] [Dot] [Line]
Xlist: L1
Freq: L2
Mark: [Box] + .
    
```



Ohne GTR findet man bei diesem einfachen Beispiel die Kennzahlen ebenfalls mühelos:

Betrachte die sortierter Liste: 1 1 2 3 3 3 4 4 5 9

AbleSEN: x_{min}=1 x_{max}=9 med=3 Q₁=2 Q₂=4 IQR=2 1,5*IQR=3

Das durch 1,5*IQR definierte Intervall zum Ausschluss der Ausreißer ist dann [-1;7] . Da 9 nicht in diesem Bereich liegt ist es ein Ausreißer und somit geht der obere Whisker nur bis 5 .

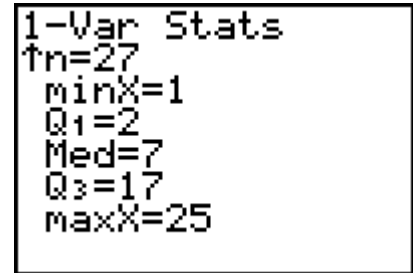
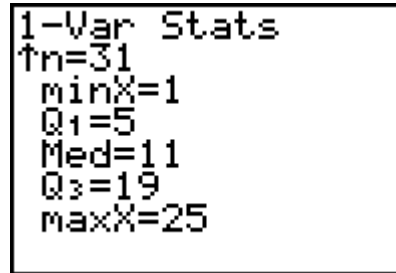
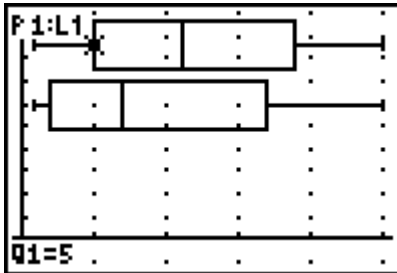
Wozu dienen Boxplots ? Vergleich zweier Datensätze:

Die Schüler der Klassen 9a und 9b geben die Entfernung (in km) ihres Wohnortes zur Schule an:

9a: 1 7 25 3 5 6 12 11 16 1 25 9 1 21 2 6 25 10 18 1 1 20 13 23 18 18 6 8 18 19 25

9b: 7 9 1 3 25 5 12 2 17 1 21 21 6 15 19 18 2 3 1 17 21 1 1 7 14 2 7

Aufgabe: Ordne die Listen und bestimme die Kennzahlen. Vergleiche mit den Ergebnissen des TI 83.



Hinweis: Die zweite Statistik erhält man mittels STAT CALC 1-Var Stats L₂ ENTER .

Schlussfolgerungen aus den Ergebnissen:

- Die 9b wohnt im Schnitt näher an der Schule
- In der 9b hat die **Hälfte** höchstens 7 km Schulweg, in der 9a höchstens 11 km
- Ein Viertel der 9b wohnt höchstens 2 km von der Schule entfernt, in der 9a sind es höchstens 5 km
- Keine Unterschiede gibt es beim kürzesten bzw. längsten Schulweg der beiden Klassen

B) Histogramme:

Das sind Rechtecke, deren **Flächen proportional zur klassenspezifischen Häufigkeit** sind.
Die Breite der Rechtecke (Klassenbreite!) kann variabel sein, was aber der TI83 nicht beherrscht.

Beispiel von oben: Wohnortentfernung von Schülern:

Man gibt am besten für jede Klasse 2 Listen ein, und zwar die jeweilige Entfernung und die dazugehörige Häufigkeit (Anzahl der Schüler mit dieser Entfernung).

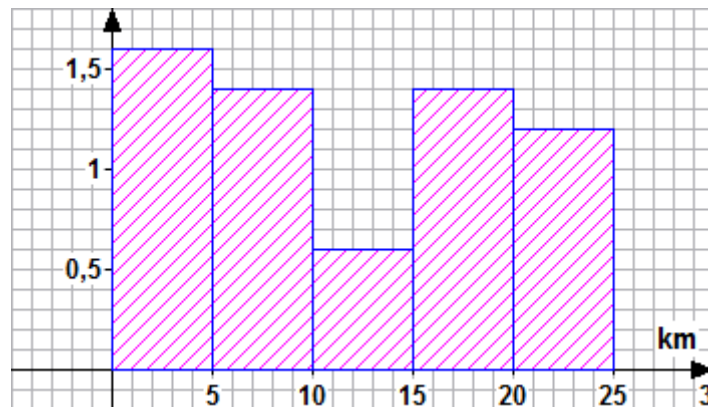
Für die 9a ist das z.B.

Entf / km	Anzahl
1	5
2	1
3	1
4	0
5	1
6	3
7	1
8	1
9	1
10	1
11	1
12	1
13	1
14	0
15	0
16	1
17	0
18	4
19	1
20	1
21	1
22	0
23	1
24	0
25	4

Teilt man nun die Entfernungen in Klassen der Breite 5km ein, so erhält man eine neue Liste mit der Klassenbreite 5, bei der die Rechteckflächen der Anzahl (Häufigkeit) entsprechen. Folglich ergibt sich die jeweilige Rechteckshöhe aus $\text{Fläche} / \text{Klassenbreite} = \text{Anzahl} / \text{Klassenbreite}$!

Entfernung in km [> a ; b]	Anzahl	Rechteckshöhe = Anzahl / 5
[0 ; 5]	8	1,6
[5 ; 10]	7	1,4
[10 ; 15]	3	0,6
[15 ; 20]	7	1,4
[20 ; 25]	6	1,2

Korrekt dargestelltes Histogramm (mit KarloPlot)



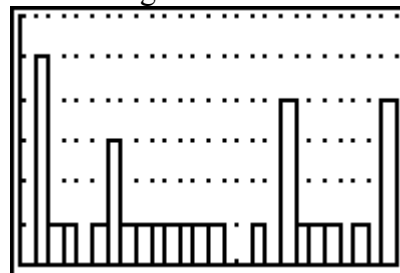
Histogramm mit TI83:

km in L1, Anzahl in L2

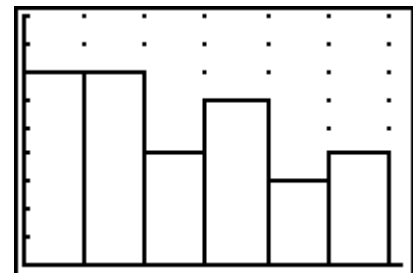
```

WINDOW
Xmin=0
Xmax=26
Xscl=1
Ymin=0
Ymax=6
Yscl=1
Xres=1
    
```

Einstellungen wie links:



mit Xscl=5 Xmax=31 Ymax=9



Man erkennt, dass

- die Rechteckshöhen nicht an die Flächen angepasst werden,
- insgesamt 6 Rechtecke statt 5 entstehen.

Offensichtlich ist der TI83 nicht für Histogramme geeignet.

C) Weitere gebräuchliche grafische Darstellungen

Außer den Histogrammen sind noch gebräuchlich:

Stängel-Blatt-Diagramm, Stabdiagramm, Häufigkeitspolygon, Kreisdiagramm (Torten-), Punktdiagramm(Scatter).

Der TI 83 bietet hiervon nur Histogramm, Scatter und Häufigkeitspolygon.

D) Speziellere grafische Möglichkeiten:

Normal-Quantil-Plot:

Sind die erhobenen Daten annähernd normalverteilt ?

Um dies zu entscheiden, kann über das Histogramm die Normalverteilungskurve mit entsprechendem Mittelwert und Standardabweichung gelegt werden.

In der explorativen Datenanalyse jedoch verwendet man Normal-Quantil-Plots.

Hierbei werden die Quantile der Häufigkeitsverteilung mit entsprechenden Quantilen der Standard-normalverteilung verglichen.

Liegen die Punkte auf einer Geraden, so spricht das für eine annähernde Normalverteilung .

Der TI 83 bietet hierfür den Plot-Type 6 .

Eine genauere Betrachtung ist nachzulesen bei

Markus Paul (T³ Europe): „Beschreibende Statistik und explorative Datenanalyse“

3) Anmerkungen zu anderer Software:

3.1) EXCEL u.a. berechnen nach der Tukey-Methode die Quartile folgendermaßen:

$$Q1 = x \left(\frac{\left[\frac{n+1}{2} \right] + 1}{2} \right) \quad \text{und} \quad Q3 = x \left(n+1 - \frac{\left[\frac{n+1}{2} \right] + 1}{2} \right)$$

[z] ist die sog. Gaußklammerfunktion
(größte ganze Zahl $\leq z$)

Bei dieser Methode wird bei ungeradem n der med-Wert in der Teiliste links (bzw. rechts) mitgezählt !
Bei geradem n ist die Methode identisch mit derjenigen des TI 83 .

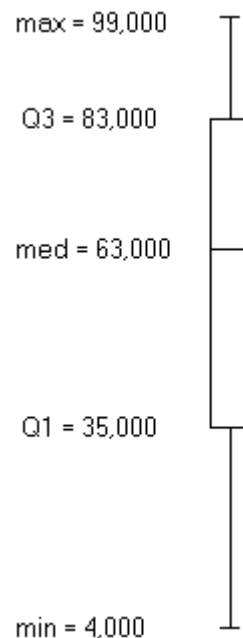
3.2) Mehrere Softwarepakete zeichnen Boxplots vertikal statt horizontal (siehe Grafik) oder sie bieten beide Darstellungsmöglichkeiten.

3.3) Die Länge der Whisker im modifizierten Boxplot wird sehr unterschiedlich gehandhabt:

a) maximal bis zum 1,5-fachen IQR-Abstand von der Box;
falls xmax bzw. xmin kleiner als dieser Abstand ist, dann bis zu xmax bzw. xmin

b) genau bis zum 0,05- bzw. 0,95-Quantil .

c) genau bis zum 0,025- bzw. 0,975-Quantil .



3.4) Gängige Statistik-Software-Pakete (kommerziell) sind:

- Fathom
- Minitab
- S-Plus
- SPSS

Meist erfordert diese Software eine nicht unbeträchtliche Einarbeitungszeit.

Es gibt aber auch freie (oder sehr preisgünstige) Pakete:

- Statistik-Labor(FU Berlin)
- VU-Statistik (Verlag Schroedel)
- Calc3D (außer Statistik noch weitere Themen)
- usw.
- GrafStat