

Eine Datenbank für den Mathe-Unterricht

Dr. Hubert Langlotz

Dr. Wilfried Zappe



¹ https://st6.cannypic.com/thumbs/21/211995_352_canny_pic.jpg

Gliederung:

Abschnitt	Seite
Vorbemerkungen	3
Stabilisierung von relativen Häufigkeiten	4
Statistische Auswertungen von Daten	6
Prognoseintervalle	10
Konfidenzintervalle	14
Daten auf Normalverteilung prüfen	20

Übersicht über begleitende Dateien für den TI-Nspire CX CAS

Daten_Neugeborene_ti_tns

Daten_Neugeborene_rH_Geschlecht.tns

Daten_Neugeborene_Prognose_Geschlecht.tns

Daten_Neugeborene_Konfidenz_einfach.tns

Daten_Neugeborene_Konfidenz_Erwartungswert.tns

Daten_Neugeborene_Konfidenz_grafisch_mit prog.tns

Daten_Neugeborene_Konfidenz_grafisch-ohne prog.tns

Eine Datenbank für den Mathe-Unterricht

Vorbemerkungen

Texas Instruments stellt hiermit eine Datenbank zur Verfügung, die interessierte Kolleginnen und Kollegen für ihren Mathe-Unterricht unentgeltlich nutzen können. Die Datenbank enthält die Angaben über Körpergröße, Körpergewicht und Geschlecht von über 2000 Neugeborenen. Über sechs Jahre lang wurden dazu Anzeigen aus Tageszeitungen ausgewertet. Die Anzeigen erscheinen etwa in wöchentlichem Rhythmus. Sicher kann man davon ausgehen, dass die Daten nur mit dem Einverständnis der Eltern veröffentlicht wurden. Sie enthalten mitunter auch Angaben über das Geburtsdatum und die Namen der Eltern. Aus Datenschutzgründen werden solche Angaben aber in der hier zur Verfügung gestellten Datenbank nicht weitergegeben. Ursprünglich erfasst wurden auch die Vornamen der Kinder. Sie wurden ebenfalls aus Datenschutzgründen entfernt. Vor der Entfernung der Vornamen wurde ermittelt, aus wie vielen Buchstaben und Zeichen der jeweilige Vorname zusammengesetzt ist. Ebenso wurden die Daten von Mehrlingsgeburten entfernt, um eine Unabhängigkeit zwischen den erhobenen Daten abzusichern.



Sie sehen im Folgenden einen Ausschnitt aus der Tabelle mit den Daten:

2

	A nr	B größe	C masse	D geschlecht	E anzbu	F
=	=seq(k,k,1					
1	1	50	3.29	0	11	
2	2	52	3.73	1	6	
3	3	57	4.295	0	4	
4	4	50	2.8	1	10	

nr: Jeder Datensatz ist mit einer Nummer versehen.

größe: Enthält die Körpergröße in cm.

masse: Zeigt das Körpergewicht in kg an.

geschlecht: Eine „0“ steht für einen neugeborenen Jungen, eine „1“ für ein Mädchen.

anzbu: Jede Zahl gibt an, aus wie vielen Buchstaben und Zeichen der Vorname des Kindes zusammengesetzt ist.

Zusammengetragen wurden im Zeitraum von 2014 bis 2020 die Daten von Dr. Wilfried Zappe (Ilmenau). Sie beziehen sich in der Mehrheit auf Neugeborene aus dem IIm-Kreis.

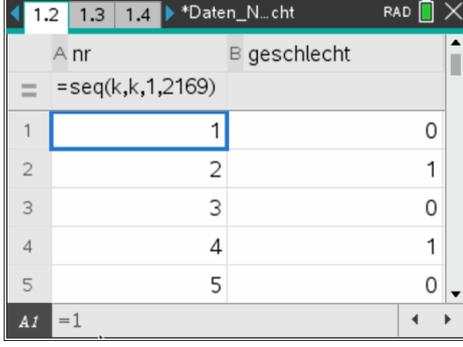
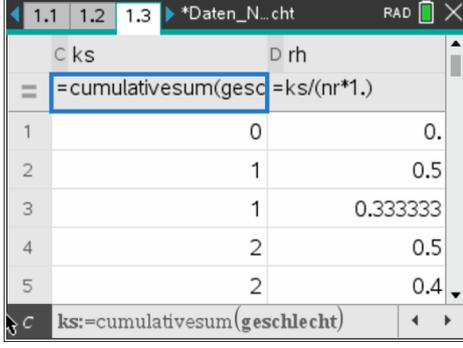
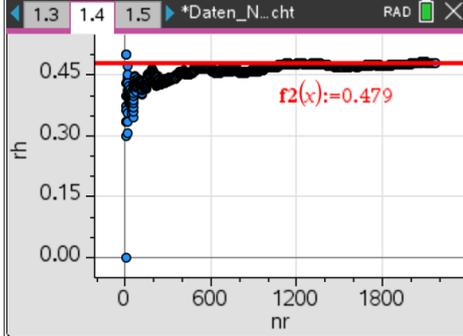
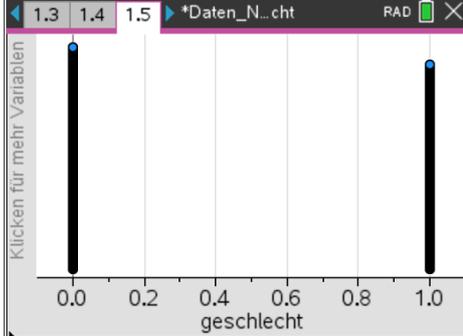
Die Kollegen Dr. Hubert Langlotz und Dr. Wilfried Zappe geben nun einige Anregungen dafür, was mit diesen Daten im Unterricht machbar ist. An dem Teil zur Statistik hat Kollege Tobias Kellner mitgearbeitet. Zur Programmierung bei Konfidenzintervallen konnten wir auf die Unterstützung des Kollegen Sebastian Rauh bauen, der außerdem den Text nochmal kritisch durchgesehen hat und schöne Ideen zum Thema beigesteuert hat. Zu besonderem Dank sind wir dem Kollegen Dr. Andreas Prömmel verpflichtet, der sehr sorgfältig alles vorhandene Material gesichtet und geprüft hat.

Hinweise oder Ergänzungen zu unseren Ausführungen nehmen wir gern entgegen. (info@t3deutschland.de)

² Quelle: „Freies Wort“ vom 04.07.2020

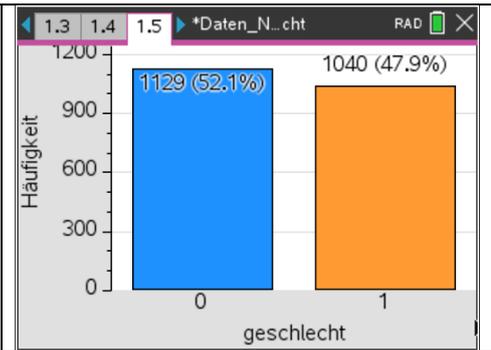
Stabilisierung von relativen Häufigkeiten und Durchschnittswerten

Mit zunehmender Anzahl von Einzeldaten ist eine Stabilisierung der relativen Häufigkeiten der Mädchen- bzw. Jungengeburten sowie der durchschnittlichen Größe und Masse neugeborener Kinder zu beobachten. Diese Prozesse können veranschaulicht werden. Sie tragen z. B. zu einem tieferen Verständnis des Wahrscheinlichkeitsbegriffes bei.

<p>Kopieren Sie die Datei „Daten Neugeborene_2“ und fügen Sie diese Datei in die Anwendung „Lists&Spreadsheet“ ein. Löschen Sie die Spalten „größe“ und „masse“, sodass zunächst nur die Spalten A: „nr“ und B: „geschlecht“ verbleiben.</p>	 <table border="1" data-bbox="922 421 1385 768"> <thead> <tr> <th>A nr</th> <th>B geschlecht</th> </tr> </thead> <tbody> <tr><td>1</td><td>0</td></tr> <tr><td>2</td><td>1</td></tr> <tr><td>3</td><td>0</td></tr> <tr><td>4</td><td>1</td></tr> <tr><td>5</td><td>0</td></tr> </tbody> </table>	A nr	B geschlecht	1	0	2	1	3	0	4	1	5	0
A nr	B geschlecht												
1	0												
2	1												
3	0												
4	1												
5	0												
<p>Bezeichnen Sie die Spalte C mit dem Namen „ks“ für kumulierte Summe und geben Sie in die Zelle darunter den Befehl „=cumsum(geschlecht)“ ein. Der Rechner ergänzt diese Anweisung automatisch mit dem vollständigen Befehl „cumulativesum(geschlecht)“ und berechnet diese kumulierten Summen. In der Spalte D tragen Sie den Namen „rh“ für relative Häufigkeit ein und fügen in die Zelle darunter die Anweisung „=ks/(nr*1.)“ ein, um die relativen Häufigkeiten zu berechnen. Der Faktor „1.“ wird verwendet, damit die relativen Häufigkeiten als Dezimalbrüche angezeigt werden.</p>	 <table border="1" data-bbox="922 801 1385 1149"> <thead> <tr> <th>C ks</th> <th>D rh</th> </tr> </thead> <tbody> <tr><td>0</td><td>0.</td></tr> <tr><td>1</td><td>0.5</td></tr> <tr><td>1</td><td>0.333333</td></tr> <tr><td>2</td><td>0.5</td></tr> <tr><td>2</td><td>0.4</td></tr> </tbody> </table>	C ks	D rh	0	0.	1	0.5	1	0.333333	2	0.5	2	0.4
C ks	D rh												
0	0.												
1	0.5												
1	0.333333												
2	0.5												
2	0.4												
<p>Fügen Sie die Anwendung „Data&Statistics“ ein. Weisen Sie der horizontalen Achse die Variable „nr“ zu und der vertikalen Achse die Variable „rh“. Sie können die Stabilisierung der relativen Häufigkeiten gut erkennen: Je höher n, desto weniger schwanken die relativen Häufigkeiten um einen gewissen festen Wert, der hier bei ca. 0,479 liegt, d. h. hier, dass der Anteil der Mädchengeburten in dieser Stichprobe bei ca. 47,9% liegt.</p>													
<p>Eine rasche Übersicht über den Anteil der Mädchen- und Jungengeburten unter allen erfassten Neugeborenen erhalten Sie über die Darstellung der Liste „geschlecht“ in der Anwendung „Data&Statistics“.</p>													

Die Anweisungen „Kategorisches X erzwingen“ und „Alle Bezeichnungen anzeigen“ ergeben den nebenstehenden Bildschirm.
Setzen Sie dazu den Cursor auf die Bezeichnung der horizontalen Achse und wählen Sie mit <ctrl><menü> die oben genannten Befehle.

Unter den 2169 erfassten Datensätzen sind 47,9% (1040) Mädchengeburten registriert.



Aufgaben:

Untersuchen Sie auf analogem Wege die Stabilisierung der durchschnittlichen Körpergröße bzw. des durchschnittlichen Geburtsgewichts Neugeborener.

Ergänzungen:

Um die etwas unübersichtlich wirkende Darstellung von über 2000 einzelnen relativen Häufigkeiten zu verbessern, lassen sich die Stabilisierungseffekte auch durch die Zusammenfassung von zunächst den ersten 100, dann den ersten 200 zu Partialsummen aufeinanderfolgenden Einzelwerten usw. veranschaulichen. Jede dieser Partialsummen wird dann durch die Anzahl ihrer Summanden geteilt.

Schieberegler für s auf Data&Statistics: | s ▶ 100

$$li := \text{seq} \left(s \cdot k, k, 1, \frac{2000}{s} \right)$$

▶ { 100,200,300,400,500,600,700,800,900,1000,1100,1200,1300,1400,1500 }

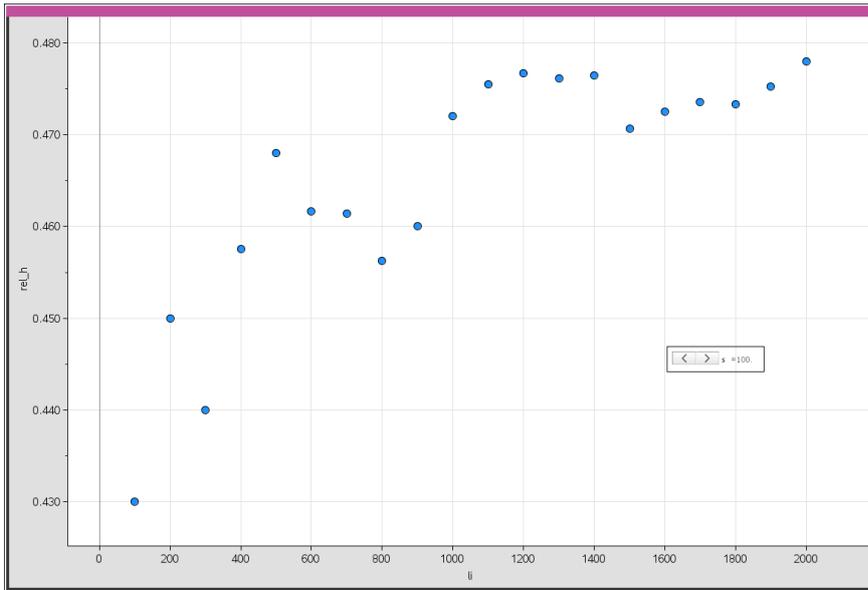
$$ma := \text{seq} \left(\sum_{k=1}^a (\text{geschlecht}[k]), a, s, 2000, s \right)$$

▶ { 43,90,132,183,234,277,323,365,414,472,523,572,619,667,706,756,805,854,902,950,1000 }

$$\text{rel_h} := \frac{ma}{li}$$

▶ { $\frac{43}{100}, \frac{9}{20}, \frac{11}{25}, \frac{183}{400}, \frac{117}{250}, \frac{277}{600}, \frac{323}{700}, \frac{73}{160}, \frac{23}{50}, \frac{59}{125}, \frac{523}{1100}, \frac{143}{300}, \frac{6}{100}, \frac{11}{250}, \frac{183}{400}, \frac{117}{250}, \frac{277}{600}, \frac{323}{700}, \frac{73}{160}, \frac{23}{50}, \frac{59}{125}, \frac{523}{1100}, \frac{143}{300}, \frac{6}{100}$ }

Die relativen Häufigkeiten sind hier nicht als Dezimalzahlen angegeben, weil sich durch die Darstellung in Brüchen ihre Herkunft besser erklären lässt.



Die grafische Darstellung muss ggf. durch *Zoom-Daten* angepasst werden.

Statistischen Auswertungen von Daten

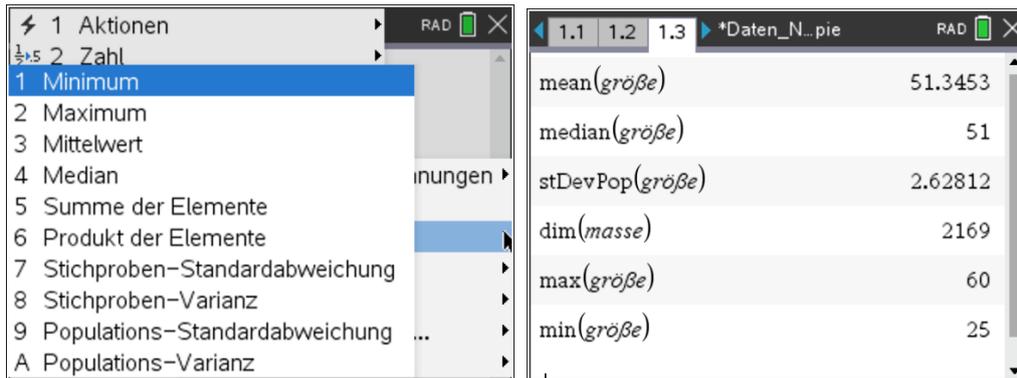
Das Merkmal „Geschlecht“ liegt in zwei Ausprägungen vor, die als „Klassenwerte“ bezeichnet werden:

Die Klassenwerte müssen als Zahlen eingegeben werden, um eine nach Klassen getrennte Auswertung vornehmen zu können.

Geschlecht	Klassenwert
männlich	0
weiblich	1

Ermittlung von statistischen Kenngrößen über Listenbefehle

Von großem Interesse für die Bewertung von Daten sind ihre Kenngrößen, wie das arithmetische Mittel, die Standardabweichung, der Median und andere. Der CAS-Rechner TI-Nspire bietet dazu verschiedene Möglichkeiten an. Ein einfacher Zugriff kann über die Anwendung *Calculator - Menü - Statistik - Listen Mathematik* erfolgen. Er bietet folgende Optionen, die eigentlich selbsterklärend sind:



Die Anweisung „dim“ gibt die Anzahl der Listenelemente zurück.

Statistik mit einer Variablen für alle Neugeborenen

Komfortabler gelingt die Anzeige von Kenngrößen mit der Anwendung *Calculator - Menü - Statistik - Statistische Berechnungen - Statistik mit einer Variablen*.

Es öffnet sich ein Fenster, das die Anzahl der Listen abfragt. Wir lassen es zunächst bei der voreingestellten „1“ und drücken „ok“. Nun öffnet sich ein weiteres Fenster (siehe Screenshot). Um die statistischen Werte für alle Neugeborenen bezüglich des Merkmals „Größe“ zu erhalten, wählen wir als x1-Liste die Variable „größe“ aus, ändern weiter nichts an den Voreinstellungen und drücken „ok“. In den Spalten E und F werden Kenngrößen und ihre aktuellen Werte angezeigt. Die farbige Unterlegung wurde nachträglich vorgenommen.



A nr	B größe	C masse	D geschlecht	E	F
=	seq(k,k,1,2169)				=OneVar('größe,1): CopyVar Stat
1	1	50	3.29	0 Titel	Statistik mit einer Variable
2	2	52	3.73	1 \bar{x}	51.3453
3	3	57	4.295	0 Σx	111368.
4	4	50	2.8	1 Σx^2	5.73321e6
5	5	57	4.35	0 $s_x := s_{n-1}x$	2.62873
6	6	52	3.69	1 $\sigma_x := \sigma_{n-1}x$	2.62812
7	7	32	3.25	0 n	2169.
8	8	52	3.612	0 MinX	25.
9	9	48	2.275	0 Q_1X	50.
10	10	50	3.2	0 MedianX	51.
11	11	51	3.15	1 Q_3X	53.
12	12	53	3.06	0 MaxX	60.
13	13	50	2.96	0 $SSX := \Sigma(x-\bar{x})^2$	14981.4

Bedeutung der angezeigten Statistikvariablen

\bar{x}	\bar{x} arithmetisches Mittel
Σx	Σx Summe der Werte
Σx^2	Σx^2 Summe der quadrierten Werte
$s_x := s_{n-1X}$	s_x Standardabweichung der Stichprobe (korrigiert)
$\sigma_x := \sigma_{nX}$	σ_x Standardabweichung der Grundgesamtheit (unkorrigierte)
n	n Anzahl der Werte
MinX	kleinster Wert
Q_1X	erstes Quartil
MedianX	Median der Werte
Q_3X	drittes Quartil
MaxX	größter Wert
$SSX := \Sigma(x-\bar{x})^2$	Summe der quadrierten Abweichungen der Einzelwerte vom arithmetischen Mittel

Aufgabe:

Bestimmen Sie auf analogem Wege die Kenngrößen der Listen „masse“ und „geschlecht“.

Statistik mit einer Variablen getrennt nach Kategorien, also wenn z. B. aus der Gesamtheit der Daten „größe“ nur die statistischen Angaben für die neugeborenen Jungen berechnet werden sollen.

Menü – Statistik – Statistische Berechnungen – Statistik mit einer Variablen wählen.

Es öffnet sich ein Fenster, das die Anzahl der Listen abfragt. Wir lassen es wieder bei der voreingestellten „1“ und drücken „ok“. Nun öffnet sich ein weiteres Fenster (siehe Bildschirmabdruck). Um die statistischen Werte für alle neugeborenen Jungen bezüglich des Merkmals „Größe“ zu erhalten, wählen wir als x1-Liste die Variable „größe“ aus. Die Häufigkeitsliste bleibt bei der voreingestellten „1“. Als Kategorielliste wählen wir die Variable „geschlecht“, tragen in das Feld Kategorien {0} für „männlich“ ein. Als Ergebnisspalte kann man eine der noch freien Spalten wählen. Wir nehmen hier wieder die Spalte E und lassen die vorigen Ergebnisse durch Drücken von „ok“ überschreiben.



A nr	B größe	C masse	D geschlecht	E	F
=	=seq(k,k,1,2169)				=OneVar('größe,1,'geschlecht,0):
1	1	50	3.29	0 Titel	Statistik mit einer Variable
2	2	52	3.73	1 \bar{x}	51.6594
3	3	57	4.295	0 Σx	58323.5
4	4	50	2.8	1 Σx^2	3.02138e6
5	5	57	4.35	0 $s_x := s_{n-1X}$	2.73169
6	6	52	3.69	1 $\sigma_x := \sigma_{nX}$	2.73048
7	7	32	3.25	0 n	1129.
8	8	52	3.612	0 MinX	25.
9	9	48	2.275	0 Q_1X	50.
10	10	50	3.2	0 MedianX	52.
11	11	51	3.15	1 Q_3X	53.
12	12	53	3.06	0 MaxX	60.
13	13	50	2.96	0 $SSX := \Sigma(x-\bar{x})^2$	8417.3

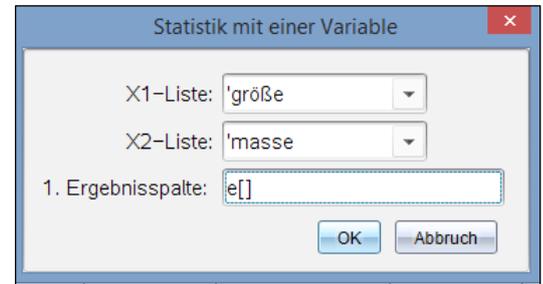
Aufgabe:

Ermitteln Sie analog die Statistikvariablen für die neugeborenen Mädchen.

Statistik mit einer Variablen für mehrere Listen

Menü – Statistik – Statistische Berechnungen – Statistik mit einer Variablen wählen.

Es öffnet sich ein Fenster, das die Anzahl der Listen abfragt. Wir wollen die Statistikvariablen für die Größe und die Masse gleichzeitig bestimmen lassen. Die voreingestellte „1“ wird durch eine „2“ ersetzt und „ok“ gedrückt. Nun öffnet sich ein weiteres Fenster (siehe Bildschirmabdruck). Um die statistischen Werte für alle Neugeborenen bezüglich der Merkmale „Größe“ und „Masse“ zu erhalten, wählen wir als x1-Liste die Variable „größe“ und als x2-Liste „masse“ aus. Für die 1. Ergebnisspalte wird Spalte E gewählt und „ok“ gedrückt.



Die Statistikvariablen für jede Liste werden in getrennten Spalten angezeigt.

	größe	masse	geschlecht	E	F	G
1	1	50	3.29	0	Titel	
2	2	52	3.73	1	\bar{x}	51.3453
3	3	57	4.295	0	Σx	111368.
4	4	50	2.8	1	Σx^2	5.73321e6
5	5	57	4.35	0	$s_x := s_{n-1}x$	2.62873
6	6	52	3.69	1	$\sigma_x := \sigma_n x$	2.62812
7	7	32	3.25	0	n	2169.
8	8	52	3.612	0	MinX	25.
9	9	48	2.275	0	Q ₁ X	50.
10	10	50	3.2	0	MedianX	51.
11	11	51	3.15	1	Q ₃ X	53.
12	12	53	3.06	0	MaxX	60.
13	13	50	2.96	0	$SSX := \Sigma(x-\bar{x})^2$	14981.4

Für eine Auswertung nach dem Klassenwert gibt es hier keine Abfrage, dies scheint für eine Anzahl von mehr als einer Liste nicht möglich zu sein.

Statistik mit zwei Variablen; gibt es einen Zusammenhang zwischen Größe und Geburtsgewicht?

Menü – Statistik – Statistische Berechnungen – Statistik mit zwei Variablen wählen.

Es öffnet sich ein Fenster, das die Anzahl der Listen abfragt. Wir wollen die Statistikvariablen für die Größe und die Masse sowie ihren Zusammenhang bestimmen lassen. Die voreingestellte „1“ wird belassen. Die Auswertung soll nur für alle erfassten Neugeborenen erfolgen, deshalb bleiben die Felder „Kategorieliste“ und „Mit Kategorien“ leer. Für die 1. Ergebnisspalte wird Spalte E gewählt und „ok“ gedrückt. (Eine Auswertung nach Kategorien wäre aber auch möglich.)



Die Statistikvariablen für jedes Merkmal werden in ein und derselben Spalte angezeigt. Es sind so viele Variablen, dass man den Bildschirm nach unten scrollen muss, um alle zu sehen.

A nr	B gröÙe	C masse	D geschlecht	E	F	G
=seq(k,k,1,2169)					=TwoVar('gröÙe','masse,1): Copy\ =OneVar('masse,1): C	
1	1	50	3.29	0 Titel	Statistiken mit zwei Variablen	Statistik mit einer Vari
2	2	52	3.73	1 \bar{x}	51.3453	3.4
3	3	57	4.295	0 Σx	111368.	739
4	4	50	2.8	1 Σx^2	5.73321e6	257
5	5	57	4.35	0 $s_x := s_{n-1}x$	2.62873	0.49
6	6	52	3.69	1 $\sigma_x := \sigma_n x$	2.62812	0.49
7	7	32	3.25	0 \bar{y}	2169.	2
8	8	52	3.612	0 \bar{y}	3.40803	
9	9	48	2.275	0 Σy	7392.01	
10	10	50	3.2	0 Σy^2	25719.4	
11	11	51	3.15	1 $s_y := s_{n-1}y$	0.493142	
12	12	53	3.06	0 $\sigma_y := \sigma_n y$	0.493029	5
13	13	50	2.96	0 Σxy	381681.	527
14	14	48	3.64	1 r	0.759838	
15	15	49	2.73	1 MinX	25.	
16	16	48	2.905	1 Q ₁ X	50.	
17	17	56	4.195	0 MedianX	51.	
18	18	47	2.56	1 Q ₃ X	53.	
19	19	52	3.59	1 MaxX	60.	
20	20	48	2.67	0 MinY	0.45	

Die Statistikvariable „r“ ist der „Korrelationskoeffizient“. Er ist definiert durch³

$$\text{Kor}_e(x, y) := \rho_e(x, y) := r_{xy} := \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2 \cdot \sum_{i=1}^n (y_i - \bar{y})^2}}$$

Er wird verwendet, um festzustellen, wie hoch ein Zusammenhang zwischen zwei numerischen oder quantitativen Merkmalen ist. Er bezeichnet das Maß für die Richtung und Stärke einer statistischen Beziehung zwischen zwei Merkmalen.

Der Korrelationskoeffizient nimmt Werte zwischen -1 und +1 an. Ein Wert von +1 bedeutet, dass bei steigenden Werten des ersten Merkmals die des zweiten Merkmals maximal steigen, und umgekehrt. Ein Wert von -1 bedeutet, dass bei steigenden Werten des ersten Merkmals die Werte des anderen maximal sinken.

Man kann sogar zeigen, dass genau dann, wenn $r = \pm 1$ ein fast sicherer affin linearer Zusammenhang zwischen x und y besteht.

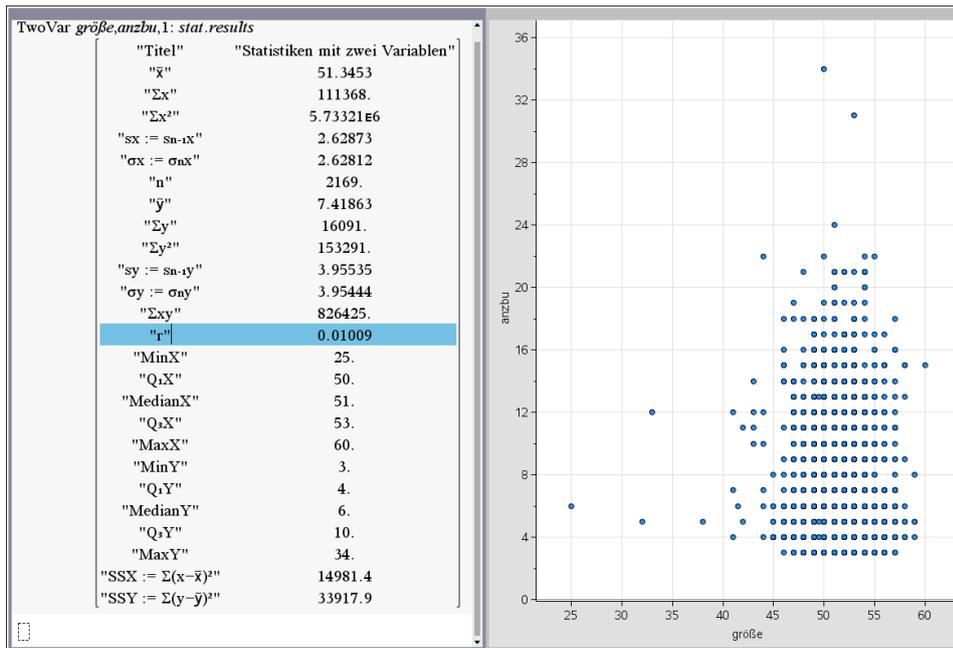
Ein Wert von $r = +0,6$ bedeutet, dass ein mittlerer positiver Zusammenhang besteht, ein Wert von +0,2, dass ein kleiner positiver Zusammenhang vorliegt. Hierbei ist zu erwähnen, dass der Koeffizient zwar etwas über die Korrelation aussagt, sich aus dem Ergebnis aber nicht der kausale Zusammenhang ableiten lässt. Als Beispiel ist hier die Besiedelung des österreichischen Südburgenlands durch Störche zu erwähnen. Diese korreliert zwar positiv mit der Geburtenrate, aber es lässt sich trotzdem kein ursächlicher Zusammenhang ableiten.

Der Korrelationskoeffizient r ist für Größe und Masse bei den hier untersuchten neugeborenen Kindern ca. 0,76. Beide Merkmale weisen also eine mittlere Korrelation auf. Man darf hier sicher auch einen kausalen Zusammenhang zwischen Größe und Körpergewicht vermuten.⁴

Auf analogem Wege lässt sich z. B. untersuchen, welche Korrelation zwischen der Größe und der Anzahl der Buchstaben des Vornamens besteht:

³ <http://de.wikipedia.org/wiki/Korrelationskoeffizient>

⁴ Vgl. z. B. Daten und Zufall im Mathematikunterricht, Cornelsen Verlag 2012, S. 95ff.



Es verwundert nicht, dass der Korrelationskoeffizient hier mit rund 0,01 nahe bei Null liegt.

Grafische Darstellung von Daten durch Boxplots

Durch Aufrufen der Anwendung *Data&Statistics* kann man die Daten auch in verschiedenen Varianten grafisch darstellen. Von besonderer Bedeutung sind dabei Boxplots.

Öffnen Sie *Data&Statistics*.

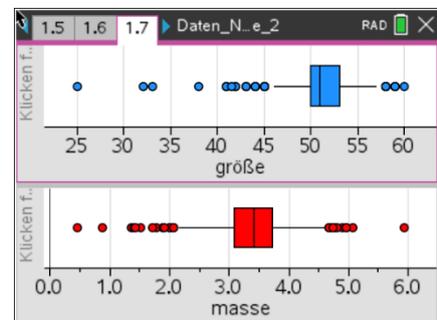
Drücken Sie <tab> und wählen Sie die Sie interessierende Liste aus.

Bestätigen Sie mit <enter>.

Drücken Sie <ctrl><menü> und wählen Sie *Box Plot*.

Sie erhalten die zugehörige *Kästchengrafik*.

Durch Überstreichen der Kästchengrafik mit dem Cursor werden das untere Quartil, der Median und das obere Quartil sowie die Werte der Ausreißer numerisch angezeigt.



Prognoseintervalle

In den Jahren 2014 bis 2019 wurden in Deutschland 2 356 752 Jungen und 2 238 378 Mädchen geboren.

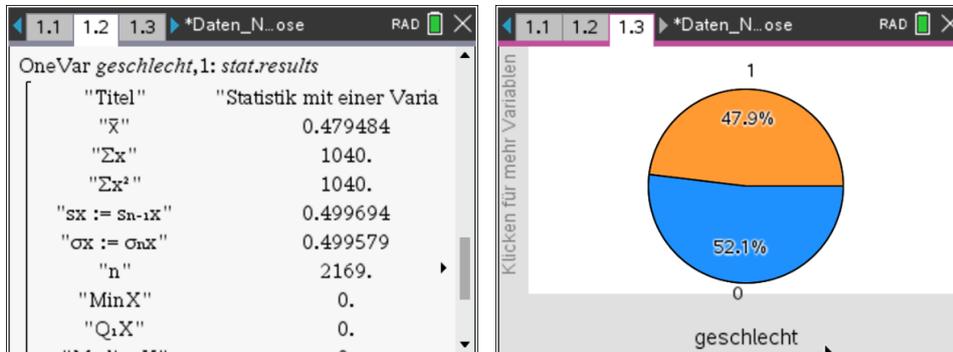
Jahr	Jungen	Mädchen
2019	399 292	378 798
2018	404 052	383 471
2017	402 510	382 374
2016	405 585	386 546
2015	378 478	359 097
2014	366 835	348 092
Summe	2 356 752	2 238 378

Quelle: <https://de.statista.com/statistik/daten/studie/880778/umfrage/anzahl-der-geburten-in-deutschland-nach-geschlecht/>

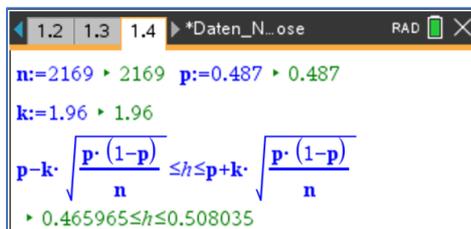
Der Anteil der Mädchengeburten in der Gesamtheit in Deutschland betrug danach in diesem Zeitraum

$$\frac{2\,238\,378}{2\,356\,752 + 2\,238\,378} \approx 0,487.$$

Der Anteil der Mädchengeburten in der in den Jahren 2014 bis Mitte 2020 erhobenen Stichprobe (n = 2169, Liste „geschlecht“) beträgt 0,479.



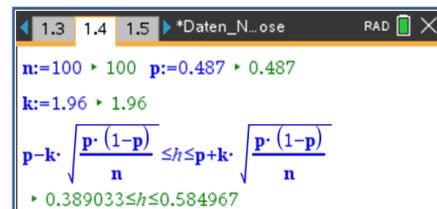
Wir prüfen, ob diese relative Häufigkeit im 95%-Prognoseintervall zu p = 0,487 liegt:



Das 95%-Prognoseintervall ist $0,466 \leq h \leq 0,508$. Die relative Häufigkeit $h = 0,479$ des Stichprobenergebnisses liegt in diesem Intervall. Das Stichprobenergebnis ist also statistisch verträglich mit $p = 0,487$.

Wir testen, wie viele von 20 (100) zufälligen Stichproben vom Umfang $n = 100$ aus der Liste „geschlecht“ statistisch verträglich mit dem 95%-Prognoseintervall sind:

Zunächst wird für $n = 100$ und $p = 0,487$ (Wahrscheinlichkeit für Mädchengeburten in Deutschland) das zugehörige 95%-Prognoseintervall $0,390 \leq h \leq 0,584$ bestimmt.



Genau eine zufällige Stichprobe vom Umfang $n = 100$ aus der Liste „geschlecht“ ohne Zurücklegen wird mit dem Befehl „randsamp (geschlecht, 100,1)“ bestimmt. Die relative Häufigkeit rh der Mädchengeburten in dieser Stichprobe kann durch die Summenbildung erfolgen, da die Mädchengeburten mit „1“ charakterisiert wurden. Im nebenstehenden Beispiel ergibt sich die Summe 52, sodass sich für die relative Häufigkeit rh feststellen lässt, dass sie im Prognoseintervall $0,390 \leq rh \leq 0,584$ liegt.



Eine Liste von 20 solcher Stichproben kann mit dem Befehl „seq“ erzeugt werden. Der Befehl „countif“ dient zum Abzählen derjenigen Listenelemente, die im 95%-Prognoseintervall

$$\text{liste} := \text{seq}(\text{sum}(\text{randsamp}(\text{geschlecht}, 100, 1)), i, 1, 20)$$

$$\text{countif}(\text{liste} \cdot 0.01, 0.389 \leq ? \leq 0.58) \rightarrow 19$$

liegen. Im nebenstehenden Beispiel wären das 19 von 20, also gerade 95% der Stichproben, die im 95%-Prognoseintervall liegen. Wird diese Berechnung wiederholt, können sich auch andere Werte ergeben, weil es sich ja um Zufallsversuche handelt. Aber eine Mittelwertbildung der Ergebnisse über eine größere Anzahl wird in der Nähe der 95%-Marke liegen. Um eine solche Wiederholung der Berechnung zu realisieren, genügt es,

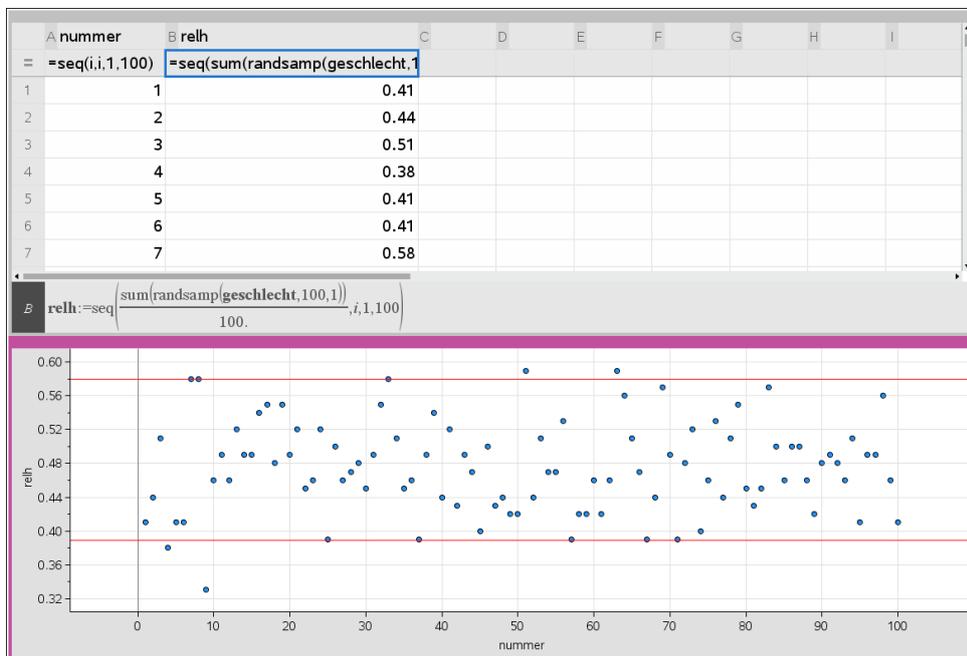
falls die Rechenschritte in der Anwendung „Notes“ formuliert wurden, den Cursor in die Anweisung „liste“ zu setzen und <enter> zu drücken.

Eine grafische Veranschaulichung für 100 solcher Stichproben kann über das nachfolgend skizzierte Vorgehen in der Anwendung „Lists&Spreadsheet“ realisiert werden. In der Spalte A werden die natürlichen Zahlen von 1 bis 100 erzeugt. Diese Liste erhält den Namen „nummer“. In der Spalte B werden 100 zufällige Stichproben ohne Wiederholung aus der Liste „geschlecht“ erzeugt und die relativen Häufigkeiten der Mädchengeburt in jeder Stichprobe berechnet. Diese Liste erhält den Namen „relh“. Die notwendigen Befehle lassen sich dem Screenshot entnehmen.

In der Anwendung „Data&Statistics“ wird jedem Element der Liste „nummer“ die zugehörige relative Häufigkeit aus der Liste „relh“ zugeordnet und das Wertepaar als Punkt dargestellt.

Zeichnet man noch die Grenzen des Prognoseintervalls als Graphen der konstanten Funktion $y = 0,390$ und $y = 0,584$ ein, lässt sich mit einem Blick erfassen, wie viele der Punkte außerhalb des Prognoseintervalls liegen.

Auch hier lassen sich in der Tabellenkalkulation mit <ctrl><R> beliebig viele Wiederholungen erzeugen.



Prognoseintervalle für den Mittelwert einer normalverteilten Zufallsgröße:

Für die annähernd normalverteilten Zufallsgrößen „masse“ und „größe“ lassen sich auf analogem Wege Prognoseintervalle für Mittelwerte bestimmen. Dazu werden der Mittelwert \bar{x} und die Standardabweichung σ_x von der Gesamtheit aller erfassten Daten verwendet. Das Prognoseintervall hat dann die Gestalt $\bar{x} - k \cdot \frac{\sigma_x}{\sqrt{n}} \leq h \leq \bar{x} + k \cdot \frac{\sigma_x}{\sqrt{n}}$ (vgl. Bigalke/ Köhler „Mathematik, Gymnasiale Oberstufe, Qualifikationsphase Leistungskurs Q3“, Cornelsen, 2018, Seite 268).

Die folgenden Bildschirme veranschaulichen das Vorgehen:

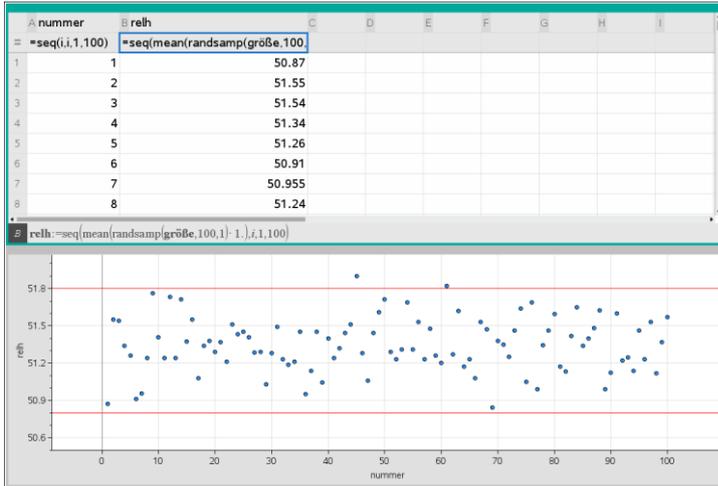
A nr	B gröÙe	C	D
=seq(k,k,1			
1	1	50	
2	2	52	
3	3	57	
4	4	50	
5	5	57	

OneVar gröÙe,1: stat.results	
"Titel"	"Statistik mit einer Varia
"x̄"	51.3453
"Σx"	111368.
"Σx²"	5.73321E6
"sx := Sn-1x"	2.62873
"sx := Onx"	2.62812
"n"	2169.
"MinX"	25.
"Q1X"	50.

```

n:=100 ▶ 100 mw:=51.3 ▶ 51.3
s:=2.63 ▶ 2.63 k:=1.96 ▶ 1.96
mw-k· $\frac{s}{\sqrt{n}}$  ≤ h ≤ mw+k· $\frac{s}{\sqrt{n}}$ 
▶ 50.7845 ≤ h ≤ 51.8155
liste:=randSamp(gröÙe,100,1)
▶ {49,51,51,51,47,50,50,49,49,50,50,53,52,4
mean(liste·1.) ▶ 51.345

```



2. Verwendung der Doppelungleichung

Auch hier bekommen wir ein Intervall, welches die „wirkliche Wahrscheinlichkeit“ überdeckt.

$$h - 1.96 \cdot \sqrt{\frac{h \cdot (1-h)}{n}} \leq h \leq h + 1.96 \cdot \sqrt{\frac{h \cdot (1-h)}{n}}$$

$$h - 1.96 \cdot \sqrt{\frac{h \cdot (1-h)}{n}} \quad | h=0.46 \text{ and } n=50$$

0.321851

$$h + 1.96 \cdot \sqrt{\frac{h \cdot (1-h)}{n}} \quad | h=0.46 \text{ and } n=50$$

0.598149

3. Nutzung der im CAS vorhandenen Näherungsformel

Die im Statistikmodul existierende Näherungsformel

liefert ebenfalls ein Sicherheitsintervall [0.33; 0.60].

Anmerkung:

Es wird hier das sogenannte Wald-Intervall berechnet, bei welchem zum einfacheren Berechnen p durch h abgeschätzt wird. Die Rundung erfolgt entsprechend der Ungleichheitszeichen in der Doppelungleichung.

```
Interval_1Prop 23,50,0.95: stat.results
┌──────────┬──────────┬──────────┐
│ "Titel"   │ "1-Prop z-Intervall" │           │
│ "CLower"  │                   │ 0.321854 │
│ "CUpper"  │                   │ 0.598146 │
│ "p̂"       │                   │ 0.46     │
│ "ME"      │                   │ 0.138146 │
│ "n"       │                   │ 50       │
└──────────┴──────────┴──────────┘
```

Erzeugung mehrerer Konfidenzintervalle unter Verwendung von „Notes“

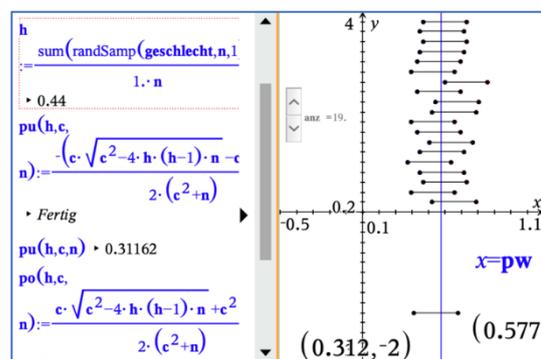
Variante 1:

Mit Hilfe der Applikation „Notes“ lassen sich schnell viele verschiedene Konfidenzintervalle erzeugen.

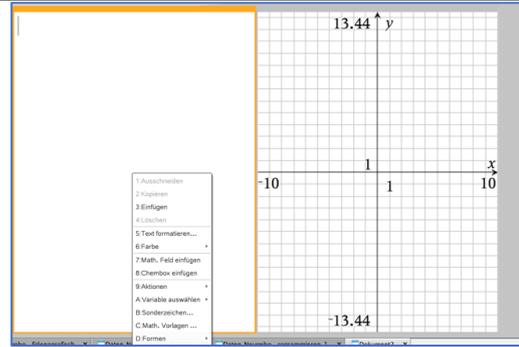
Im linken Fenster sind die notwendigen Funktionen definiert und im rechten lassen sich mittels Schieberegler mehrere Konfidenzintervalle darstellen.

Hier wurden 20 Konfidenzintervalle erzeugt, von denen 19 die „gesuchte“ Wahrscheinlichkeit überdecken.

Anleitung zur Erzeugung der Darstellung



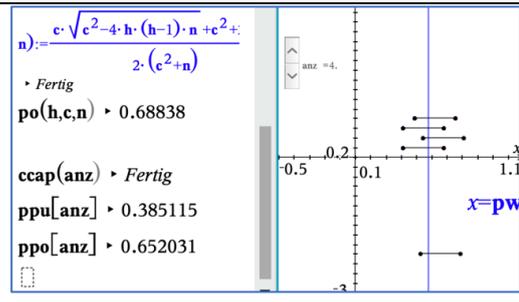
Zuerst erzeugt man sich ein zweigeteiltes Fenster, links fügt man eine Notes- und rechts eine Graphs-Applikation ein.
Beachten Sie, dass man im Notes-Fenster normalen Text, Bilder und mathematische Formeln (Math-Feld-einfügen) einfügen kann.



Die Stichprobengröße n ist ebenso wie der Faktor c (hier 1,96 für das 1,96-fache der Sigma-Umgebung) frei wählbar.
Mit den Variablen pu und po werden die Grenzen des jeweiligen Konfidenzintervalls berechnet.
(Man erhält die Terme durch Umformung der Doppelungleichung nach p .)
Befindet sich der Cursor im Definitionsfeld der Variable h , so kann durch <Enter> immer eine neue Punktschätzung erzeugt werden.

Konfidenzintervalle
Stichproben: $n:=50 \blacktriangleright 50$
 $c:=1,96 \blacktriangleright 1.96$
 $h:=\frac{\text{sum}(\text{randSamp}(\text{geschlecht},n,1))}{1 \cdot n} \blacktriangleright 0.56$
 $pu(h,c,n) := \frac{-(c \cdot \sqrt{c^2 - 4 \cdot h \cdot (h-1) \cdot n} - c^2 - 2 \cdot h \cdot n)}{2 \cdot (c^2 + n)} \blacktriangleright \text{Fertig}$
 $pu(h,c,n) \blacktriangleright 0.423058$
 $po(h,c,n) := \frac{c \cdot \sqrt{c^2 - 4 \cdot h \cdot (h-1) \cdot n} + c^2 + 2 \cdot h \cdot n}{2 \cdot (c^2 + n)} \blacktriangleright \text{Fertig}$

Möchte man mehrere Konfidenzintervalle gleichzeitig darstellen, so kann man dies z. B. mittels eines kleinen Programmes (ccap) unter Nutzung eines Schiebereglers (anz) erreichen.



Das Programm ccap leistet Folgendes:
Über den Parameter k wird die Nummer des zu erzeugenden Konfidenzintervalls übergeben, danach wird eine neue Stichprobe gezogen und die neue untere bzw. obere Grenze des Intervalls ermittelt (diese beiden Werte $ppu[k]$ bzw. $ppo[k]$ werden dann genutzt, um mit dem Wert $yy[k]$ die Endpunkte der Strecken darzustellen, welche das Konfidenzintervall im Grafikfenster darstellen sollen. Die Strecken im Grafikfenster müssen dann einmal manuell erzeugt werden.
Die For-Schleife sorgt nur dafür, dass immer nur die Intervalle von 1 bis k dargestellt werden, alle übrigen Strecken werden mittels des Befehls $yy[i]=-5$ in den nicht sichtbaren Bereich verschoben.
Im Graphikfenster ist zusätzlich noch die „Wahrscheinlichkeit pw in der Grundgesamtheit“ dargestellt. Diese Variable muss vorher definiert werden.

```

©Programm ccap
erstellt von S. Rauh

ccap
Define ccap(k)=
Prgm
Local probe,i
probe:=sum(randSamp(geschlecht,n,1))
1·n
ppu[k]:=pu(probe,c,n)
ppo[k]:=po(probe,c,n)
yy[k]:=k·0.2
For i,k+1,20
yy[i]:=-5
EndFor
EndPrgm

```

$$pw := \frac{\text{sum}(\text{geschlecht})}{\text{dim}(\text{geschlecht})} = \frac{1055}{2197}$$

Um auf die berechneten Werte zugreifen zu können, benötigt man noch ein Tabellenkalkulationsfenster, in welchem man in den drei dargestellten Spalten die durch das Programm erfassten Werte darstellt.

	A ppu	B ppo	C yy	D
=				
1	=0.366442	0.633557	0.2	
2	0.293748	0.557668	-5	
3	0.329694	0.596014	-5	
4	0.403986	0.670306	-5	
5	0.347969	0.614885	-5	
6	0.31162	0.576942	-5	
7	0.31162	0.576942	-5	

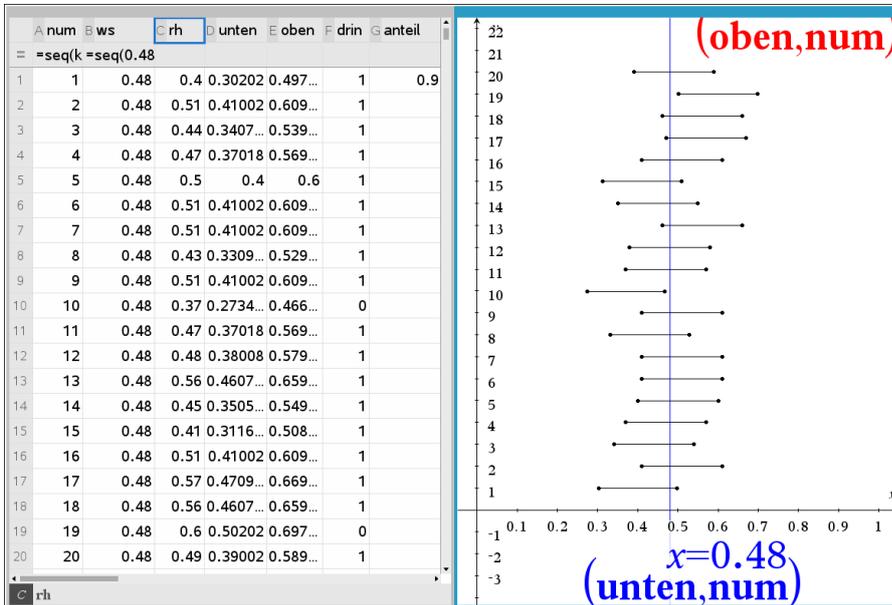
Aufgaben:

- Testen Sie das Programm Daten_Neugeborene_Konfidenz_grafisch_mit prog.tns
- Wählen Sie verschiedene Stichprobengrößen und untersuchen Sie den Zusammenhang zwischen Stichprobengröße und Intervalllänge.

Variante 2:

Einfaches Beispiel ohne Programmierung zur Veranschaulichung der Sicherheitswahrscheinlichkeit für Konfidenzintervalle

Es werden zwanzig 95,4%-Konfidenzintervalle für die Anteile der Mädchengeburt in Stichproben vom Umfang n = 50 aus der Liste „geschlecht“ erzeugt und grafisch veranschaulicht. Durch wiederholtes Durchführen lassen sich Mittelwerte für den Anteil der Konfidenzintervalle ermitteln, die die „unbekannte“ Wahrscheinlichkeit der Mädchengeburt in der Gesamtheit überdecken.



In diesem Beispiel erhält man zwei von zwanzig Konfidenzintervalle, die p = 0,48 nicht enthalten. 90% dieser Konfidenzintervalle überdecken also die Wahrscheinlichkeit p = 0,48.

Spalte	Bedeutung
A	Nummer der Simulation
B	angenommener Wert p = 0,48 der Wahrscheinlichkeit für eine Mädchengeburt
C	Anteil der Mädchengeburt in einer Stichprobe vom Umfang n = 50 aus der Liste $geschlecht: = \text{sum}(\frac{randsamp(geschlecht,50,1)}{50})$

D	Näherungswert für die untere Grenze des Konfidenzintervalls: $= C1 - 2 \cdot \sqrt{\frac{C1 \cdot (1-C1)}{50}}$
E	Näherungswert für die obere Grenze des Konfidenzintervalls: $= C1 + 2 \cdot \sqrt{\frac{C1 \cdot (1-C1)}{50}}$
F	Legt den Wert 1 fest, wenn die binomialverteilte Zufallszahl im Konfidenzintervall liegt, sonst den Wert 0: $= when(D1 \leq B1 \leq E1, 1, 0)$
G	Gibt den Anteil der Konfidenzintervalle an, die mit p verträglich sind: $= \frac{sum(drin)}{20}$

Hinweise:

Die Befehle für die Spalten C, D, E und F werden in Zeile 1 eingetragen und mit <Menü> <Daten> <Füllen> bis in die Zeile 20 als relative Zellbezüge kopiert. Die Spalten werden- wie oben zu sehen ist- bezeichnet. Die gewonnenen Daten werden als zwei Streudiagramme veranschaulicht. Die erhaltenen Punkte werden mit dem Geometriewerkzeug durch Strecken verbunden. Diese Strecken veranschaulichen die Konfidenzintervalle. Durch <CTRL> <R> in der Tabellenkalkulation können die Simulationen beliebig oft wiederholt werden. Zeichnet man noch die Gerade $x = 0,48$ ein, so lässt sich gut erkennen, welches der 20 Konfidenzintervalle $p = 0,48$ nicht überdeckt, falls ein solche Situation eintritt.

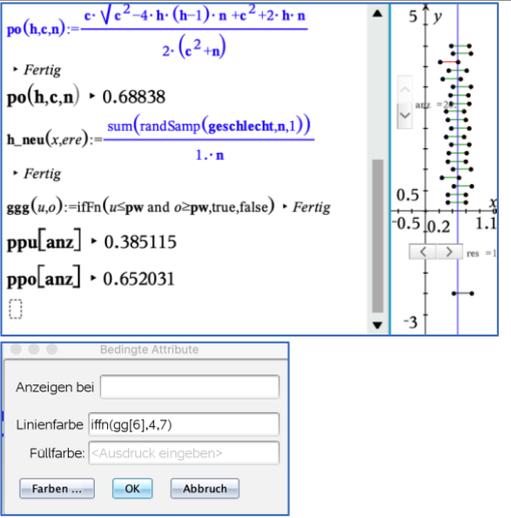
Aufgaben:

- Realisieren Sie die Simulationen auf Ihrem CAS-Rechner.
- Wiederholen Sie die Simulation mit <CTRL> <R> zehnmal.
- Ermitteln Sie für Ihre zehn Wiederholungen einen Durchschnittswert für den Anteil der Konfidenzintervalle, die mit $p = 0,48$ statistisch verträglich sind.
- Beschreiben Sie, wie die Simulation für andere Werte von p oder n angepasst werden kann.

Variante 3:

Konfidenzintervalle für Mädchengeburten ohne Programmierung

(Daten_Neugeborene_Konfidenz_grafisch-ohne prog.tns)

<p>Um ohne Programmierung auszukommen, benötigt man z. B. noch einen weiteren Schieberegler (res), der für die Aktualisierung der Daten im Lists&Spreadsheetfenster zuständig ist.</p> <p>Auf diesen wird dann im L&S-Fenster über die Variable <code>h_neu</code> zugegriffen.</p> <p>Die Variable <code>ggg</code> dient nur dazu, jene Konfidenzintervalle rot zu färben, die die wirkliche Wahrscheinlichkeit nicht überdecken. (Hierzu müssen bei allen Strecken die Bedingungen zum Zeichnen eingetragen werden.</p>	
--	--

<p>Erläuterung zu den einzelnen Spalten: aa: Berechnung von 20 neuen Punktschätzungen ppu, ppo: Intervallgrenzen ermitteln yy: y-Koordinate zuordnen ggg: boolesche Variable, um die Strecken rot färben zu können, falls die Prüfung den Wert false ergibt</p> <p>Die notwendigen Formeln kann man der tns-Datei entnehmen.</p>	<table border="1"> <thead> <tr> <th></th> <th>A aa</th> <th>B ppu</th> <th>C ppo</th> <th>D yy</th> <th>E gg</th> </tr> </thead> <tbody> <tr> <td>=</td> <td>seq(h_ne</td> <td>pu(aa,'c,</td> <td>po(aa,'c,</td> <td>seq('k*0.</td> <td>ggg(ppu,</td> </tr> <tr> <td>1</td> <td>0.48</td> <td>0.347969</td> <td>0.614885</td> <td>0.2</td> <td>true</td> </tr> <tr> <td>2</td> <td>0.54</td> <td>0.403986</td> <td>0.670306</td> <td>0.4</td> <td>true</td> </tr> <tr> <td>3</td> <td>0.5</td> <td>0.366443</td> <td>0.633557</td> <td>0.6</td> <td>true</td> </tr> <tr> <td>4</td> <td>0.5</td> <td>0.366443</td> <td>0.633557</td> <td>0.8</td> <td>true</td> </tr> <tr> <td>5</td> <td>0.5</td> <td>0.366443</td> <td>0.633557</td> <td>1.</td> <td>true</td> </tr> <tr> <td>6</td> <td>0.46</td> <td>0.329694</td> <td>0.596014</td> <td>1.2</td> <td>true</td> </tr> <tr> <td>7</td> <td>0.34</td> <td>0.224368</td> <td>0.478464</td> <td>1.4</td> <td>false</td> </tr> <tr> <td>8</td> <td>0.48</td> <td>0.347969</td> <td>0.614885</td> <td>1.6</td> <td>true</td> </tr> </tbody> </table> <p>ppu:=pu(aa,'c,'n) gg:=ggg(ppu,ppo)</p>		A aa	B ppu	C ppo	D yy	E gg	=	seq(h_ne	pu(aa,'c,	po(aa,'c,	seq('k*0.	ggg(ppu,	1	0.48	0.347969	0.614885	0.2	true	2	0.54	0.403986	0.670306	0.4	true	3	0.5	0.366443	0.633557	0.6	true	4	0.5	0.366443	0.633557	0.8	true	5	0.5	0.366443	0.633557	1.	true	6	0.46	0.329694	0.596014	1.2	true	7	0.34	0.224368	0.478464	1.4	false	8	0.48	0.347969	0.614885	1.6	true
	A aa	B ppu	C ppo	D yy	E gg																																																								
=	seq(h_ne	pu(aa,'c,	po(aa,'c,	seq('k*0.	ggg(ppu,																																																								
1	0.48	0.347969	0.614885	0.2	true																																																								
2	0.54	0.403986	0.670306	0.4	true																																																								
3	0.5	0.366443	0.633557	0.6	true																																																								
4	0.5	0.366443	0.633557	0.8	true																																																								
5	0.5	0.366443	0.633557	1.	true																																																								
6	0.46	0.329694	0.596014	1.2	true																																																								
7	0.34	0.224368	0.478464	1.4	false																																																								
8	0.48	0.347969	0.614885	1.6	true																																																								
<p>Ein Klick auf den Schieberegler „res“ erzeugt 20 neue Intervalle, welche man sich dann durch Klicken auf den Schieberegler „anz“ anzeigen lassen kann. Der Schieberegler „res“ wird benötigt, um in der Definition der Variable „aa“ eine Erneuerung zu erzeugen.</p> <p>Die notwendigen Formeln kann man der tns-Datei entnehmen.</p>																																																													

Konfidenzintervalle für den Erwartungswert einer Zufallsgröße:

Für die quantitativen Zufallsgrößen „masse“ und „größe“ lassen sich auf analogem Wege Konfidenzintervalle für Erwartungswerte bestimmen. Dazu werden der Mittelwert \bar{x} und die Standardabweichung σ_x der in einer Stichprobe erfassten Daten verwendet.

Der Parameter k beschreibt das angestrebte Sicherheitsniveau, z. B. gilt $k = 1,96$ für ein Sicherheitsniveau von 95%.

Das Konfidenzintervall hat dann die Gestalt $\mu - k \cdot \frac{\sigma_x}{\sqrt{n}} \leq \bar{x} \leq \mu + k \cdot \frac{\sigma_x}{\sqrt{n}}$

(vgl. Bigalke/ Köhler „Mathematik, Gymnasiale Oberstufe, Qualifikationsphase Leistungskurs Q3“, Cornelsen, 2018, Seite 268).

Der angegebenen Literatur entnimmt man noch folgenden Hinweis:

„Diese Näherungsverfahren darf nur angewendet werden, wenn der Stichprobenumfang n mindestens 30 beträgt und die Stichprobe maximal 5% der Grundgesamtheit umfasst.“ (ebenda, Fußnote auf Seite 268)

Betrachten wir den gegebenen Datensatz der Körpergröße Neugeborener als Stichprobe für alle Neugeborenen in Deutschland im gleichen Zeitraum, dann ergibt sich danach als Erwartungswert für die Körpergröße Neugeborener das 95%-Konfidenzintervall

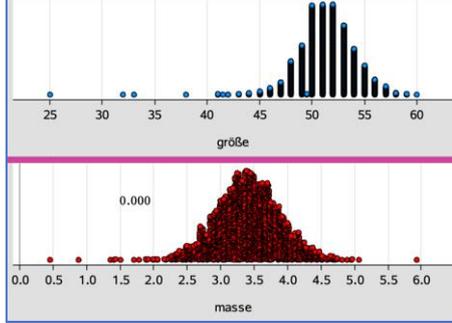
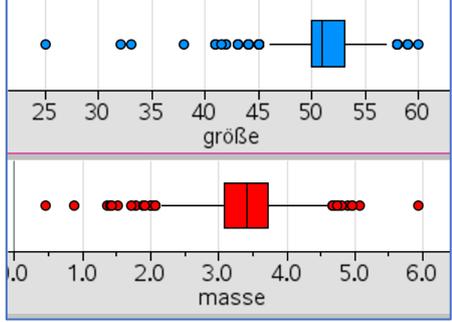
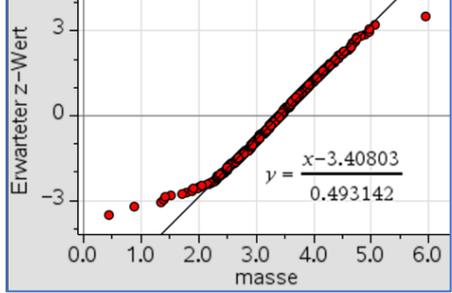
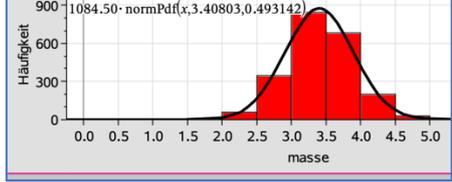
$$51,24 \leq \mu \leq 51,45.$$

Aufgabe:

Ermitteln Sie ein 95%-Konfidenzintervall für den Erwartungswert des Körpergewichtes Neugeborener.

Daten auf Normalverteilung prüfen

Mit schulischen Mitteln können wir nur qualitative Aussagen treffen, ob eine gegebene Verteilung annähernd normalverteilt ist oder nicht.

<p>Darstellung als Punktdiagramm</p> <p>Zunächst kann man die Daten in einem Punktdiagramm darstellen, um eine erste Vorstellung der Verteilung zu bekommen.</p> <p>Eine glockenförmige Verteilung liefert ein erstes Indiz für eine mögliche Normalverteilung.</p>	
<p>Nutzung der Boxplots</p> <p>Durch Veränderung des Plottyps auf „Boxplot“ können beide Verteilungen neu dargestellt werden.</p> <p>Umso symmetrischer ein Boxplot ist, desto besser kann man die Verteilung durch die Normalverteilung nähern. Man sieht hier z. B., dass die Masse anscheinend besser durch eine Normalverteilung angepasst werden kann, als die Größe.</p>	
<p>Normal Wahrscheinlichkeitsdiagramm</p> <p>Mit <code>ctrl</code> <code>menu</code> wird der Befehl Normal Wahrscheinlichkeitsdiagramm ausgelöst. Die Datenpunkte werden um eine Gerade angeordnet. Je dichter sie an der Geraden liegen, desto besser lassen sich die Daten durch eine Normalverteilung modellieren. Der gleichzeitig angezeigten Geradengleichung (hier: $y = \frac{x-3,40803}{0,493142}$) lassen sich der Erwartungswert (3,41) und die Standardabweichung (0,49) entnehmen.</p>	
<p>Normalverteilung anzeigen</p> <p>Zunächst ändert man den Plottyp auf „Histogramm“. Mit <code>Menü – Analysieren – Normal Pdf anzeigen</code> kann eine Glockenkurve mit zugehöriger Gleichung erzeugt werden. Die Werte für Erwartungswert und Standardabweichung stimmen mit den vorher erzeugten überein.</p>	
<p>Sigma-Regeln nutzen</p> <p>Gemäß der Sigma-Regeln müssen in der Ein-Sigma-Umgebung des Erwartungswertes ca. 68% der Werte liegen, für die Aufgabe wäre dies das Intervall $[3,41-0,49; 3,41+0,49] = [2,92; 3,9]$.</p>	<pre>dim(masse) 2169 countIf(masse,2.92<=:3.9) 1535 1535 0.707699 2169</pre>