

Ein bisschen Statistik

(orientiert an Hüsler/Zimmermann (2006) mit Umsetzung auf die linguistische Fragen)

1. Datentypen und Grafik

Grafische Darstellungen helfen Statistik zu verstehen. Grafiken helfen auch fehlerhafte Daten zu entdecken und Daten zu überblicken. Für jeden Datentyp gibt es passende Grafiken, auch deshalb ist es wichtig die Datentypen zu erkennen. Das wird später noch wichtiger, wenn die verschiedenen Berechnungen angestellt werden.

1.1 Qualitative / kategorielle Daten

qualitative Daten bestehen aus Werten, die Kategorien zuzuordnen sind.

- **nominale Daten** kategorisieren nach intern nicht geordneten Kategorien: Männlich/weiblich, maskulin/feminin/neutrum, Verb/Nomen/Adjektiv//a:/, /e:/...
- **ordinale Daten** zeigen eine innere Ordnung: langsam/normal/schnell, unbetont/ betont/ fokussiert, Anfänger/ Fortgeschrittene/ Profis, [a:]/[E:]/[e:]/[i:], ..

Beispiel 1

Beispiel 1 listet die Anzahl der Phoneme und Allophone auf, die im Text vorkommen. Die Daten sind in Tab. 1 gegeben.

Laut	abs. Häufigk.	rel. Häufigk.
a	15	9.5%
a:	9	5.7%
6	7	4.4%
E	26	16.5%
E:	3	1.9%
e:	7	4.4%
@	17	10.8%
e	1	0.6%
l	18	11.4%
i:	11	7.0%
i	1	0.6%
O	12	7.6%
o:	3	1.9%
o	1	0.6%
U	6	3.8%
u	1	0.6%
u:	1	0.6%
Y	4	2.5%
al	10	6.3%
aU	4	2.5%
OY	1	0.6%

Tab. 1: absolute und relative Häufigkeit der vorkommenden Laute

Die Darstellung erfolgt sinnvollerweise in Säulen- (Balken-)diagrammen.

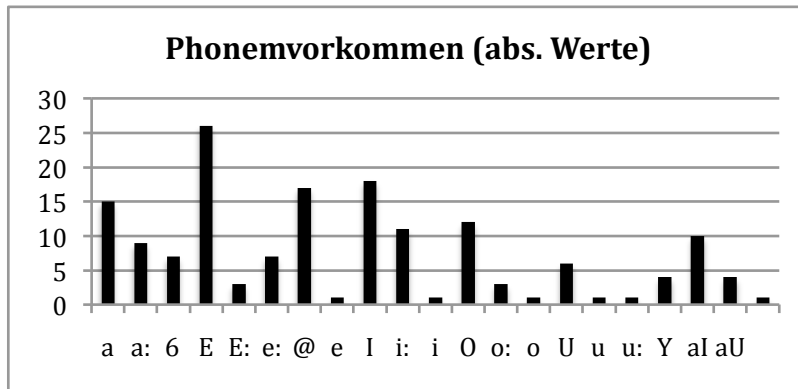


Abb. 1: Säulendiagramm der Vorkommenshäufigkeit der Vokalphoneme

Bei einer kleineren Anzahl Kategorien ist auch ein Kuchen-/Tortendiagramm möglich. Abb.2 zeigt, dass eine große Anzahl Kategorien die Grafik unübersichtlich werden lassen.

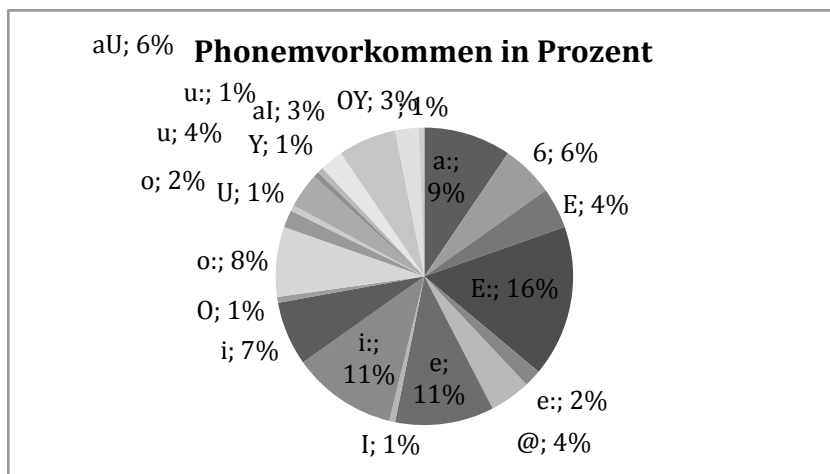


Abb. 2: Kuchendiagramm der Vorkommenshäufigkeit der Vokalphoneme

1.2 Quantitative Daten

Quantitative Daten zeigen Werte, die Rechenoperationen zulassen.

- **ordinale Daten** (mit quantitativer Bedeutung): Intelligenzquotient
- **intervallskalierte Daten** zeigen keinen absoluten Nullwert, die Verhältnisse sind bedeutungslos, Differenzen berechnungsfähig: Kalenderzeit,
- **rationale Daten** haben einen absoluten Nullwert, Verhältnisse sind bedeutsam: Länge, Dauer, Gewicht.

Beispiel 2

Beispiel 2 listet in Tab. 2 die Dauer von /a/ in Millisekunden auf. Dabei werden die beiden Geschwindigkeiten angegeben, sowie die Differenz.

a			
normal	schnell		Differenz
	48	86	-38
	65	60	5
	63	57	6
	37	22	15
	66	36	30
	52	50	2
	82	67	15
	115	85	30
	112	31	81
	69	47	22
	69	41	28
	71	65	6
	73	52	21
	62	42	20
	43	36	7

Tab. 2: Dauer von /a/ in ms in normaler und schneller Geschwindigkeit sowie die Differenz

Die Vorkommenshäufigkeit wird hier meist in Histogrammen dargestellt, wobei die Rechtecke die Klassen darstellen. die Wahl der Klassenbreite bestimmt oft das Bild (siehe Abb. 3):

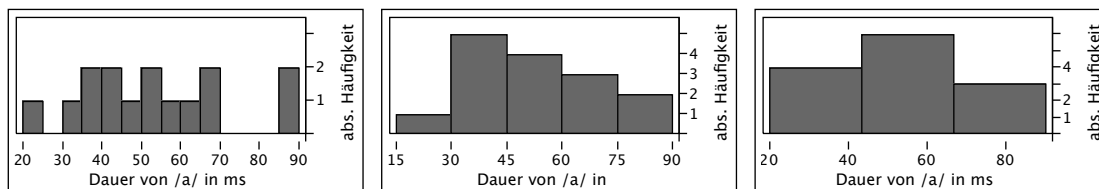


Abb. 3: Histogramm der Daten der schnellen Geschwindigkeit von Bsp. 2, man beachte die unterschiedlichen Klassengrenzen der drei Abb.

In Excel lassen sich keine Histogramme zeichnen, dafür lässt sich ein Flächendiagramm, das die Grenzen ausglättet darstellen. Allerdings müssen dafür zuerst die Kategorien mit der Form (=Zählenwenn()) errechnet werden.

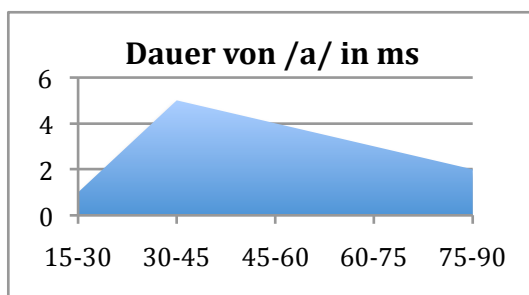


Abb. 4: Flächendiagramm der Daten der schnellen Geschwindigkeit von Bsp. 2

Das Kastendiagramm zeigt dazu schon eine Datenreduktion indem die wesentlichen Maßzahlen wie Median, Quartile und Maximum dargestellt werden.

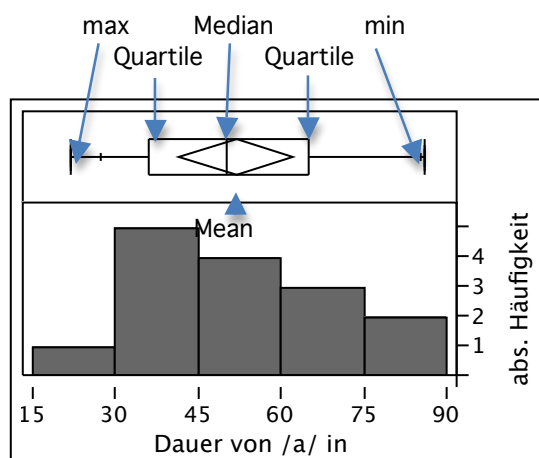


Abb. 5: Histogramm und Quantilen Box-Plot (Kastendiagramm) der Daten der schnellen Geschwindigkeit von Bsp. 2

1.3. Multivariate Daten

Häufig werden gleichzeitig mehrere Variablen in einem Versuch gemessen. Im Bsp. 2 sind das für jeden Laut die Dauer beim normalen Lesen und beim schnellen Lesen. (möglich wären aber z.B. auch Dauer und Grundfrequenz). Solche Daten werden als multivariate Daten bezeichnet, damit braucht man auch Darstellungen, die mehrere Ebenen darstellen. Bei zwei Variablen ist ein zweidimensionales Punktediagramm / Scatterplot üblich, eine dritte Dimension kann allenfalls mit unterschiedlichen Symbolen oder Farben dargestellt werden. Dreidimensionale Darstellungen wie Excel sie anbietet sind nicht besonders übersichtlich.

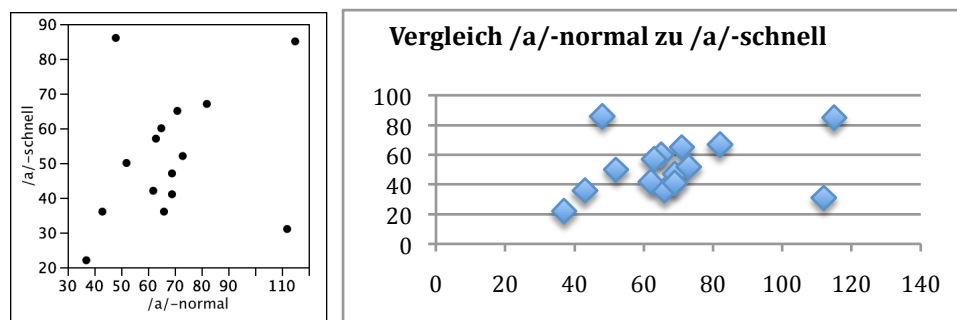


Abb. 6: Punktediagramm/Scatterplot der Dauer von /a/ in normal und schnell gelesenen Text (Bsp. 2) (zwei Darstellungen (Excel rechts))

2 Statistische Maßzahlen

Statistische Maßzahlen geben einen Datenüberblick und stellen schon eine Datenreduktion dar. Ein Wert allein (häufig der Mittelwert/mean) ist aber häufig zu stark reduzierend, deshalb müssen noch Lage- oder Streuparameter beigezogen werden. Diese Messwerte (außer dem Modus) sind nur bei intervallskalierten und rationalen Daten sinnvoll.

2.1 mittlere Lage

Die mittlere Lage wird durch verschiedene Werte angegeben, deren Verhältnis Aussagen zur Verteilung ermöglicht.

- **Mittelwert (mean)** \bar{x} : Der Mittelwert einer Summe von n Beobachtungen x_i , $i = 1, 2, \dots, n$ ist die Summe aller Beobachtungswerte x_i dividiert durch die Anzahl n der Beobachtungen:
$$\bar{x} = (x_1 + x_2 + \dots + x_n) / n$$
- **Median (median)** \tilde{x} : Der Median (Zentralwert) ist der Rangmittelpunkt der Beobachtungswerte, wenn diese der Größe nach geordnet sind. Bei n Beobachtungen ist es der $(n+1)/2$ -te Wert der nach Wert geordneten Beobachtungsreihe.
- **Modus (mode)** ist der am häufigsten vorkommende Wert.

Wenn *mean* und Median übereinstimmen, so ist das ein Hinweis auf eine symmetrische Verteilung, sind sie verschoben, so ist eine schiefe Verteilung der Daten zu erwarten (siehe Abb. 3 und Abb. 5)

2.2 Weitere Lageparameter

Andere Werte wie die Perzentile und Quantile beschreiben nicht zentrale Werte.

Perzentile: Die x -te Perzentile gibt den Anteil der Werte an die unter diesem Wert liegen: Die 50-te Perzentile ist der Median; die 90-te Perzentile gibt an, dass 90% aller Werte unter diesem liegen.

Quartile: Das erste Quartil entspricht dem 25. Perzentil, das zweite Quartil entspricht dem Median, das dritte Quartil entspricht dem 75. Perzentil. Quartilen werden im Boxplot verwendet (siehe Abb. 5)

Minimum: Der kleinste Wert der Daten.

Maximum. Der größte Wert der Daten.

2.3 Streuung

Für die Darstellung der Streuung werden Varianz und am häufigsten die Standardabweichung verwendet.

- Die Varianz s^2 entspricht (ungefähr) dem Mittelwert der quadrierten Abweichung $(x_i - \bar{x})^2$ der Beobachtung x_i vom Mittelwert \bar{x} . Der Nachteil der Varianz ist, dass sie durch die Quadrierung eine andere Einheit als die Daten besitzt (in unserem Fall wären das ms^2). Deshalb wird statt der Varianz meist die Standardabweichung angegeben.
$$s^2 = ((x_1 - \bar{x})^2 + (x_2 - \bar{x})^2 + \dots + (x_n - \bar{x})^2) / (n-1)$$
- Standardabweichung s (standard deviation, SD): Die Standardabweichung ist gleich der Quadratwurzel der Varianz. Sie hat damit dieselbe Einheit und lässt sich auf derselben Grafik darstellen.

Mittelwert und Standardabweichung werden häufig verwendet um die Verteilung der Daten zu beschreiben. In normalverteilten Daten gilt, dass

- 2/3 der Daten im Bereich $(\bar{x}-s, \bar{x}+s)$ liegen
- 95 % der Daten im Bereich $(\bar{x}-2s, \bar{x}+2s)$ liegen
- 99 % der Daten im Bereich $(\bar{x}-3s, \bar{x}+3s)$ liegen.

Bei nicht normalverteilten Daten

- 75 % der Daten im Bereich $(\bar{x}-2s, \bar{x}+2s)$ liegen
- 90 % der Daten im Bereich $(\bar{x}-3s, \bar{x}+3s)$ liegen.

Die Intervalle $\bar{x} \pm s$, $\bar{x} \pm 2s$, $\bar{x} \pm 3s$ werden als ein-, zwei-, dreifache Standardbereiche bezeichnet.

Tab. 3 zeigt die Daten der Werte für die schnelle Sprechgeschwindigkeit

	a		
		$x_i - \text{mean}$	$(x_i - \text{mean})^2$
x1	22	-29.8	888.04
x2	31	-20.8	432.64
x3	36	-15.8	249.64
x4	36	-15.8	249.64
x5	41	-10.8	116.64
x6	42	-9.8	96.04
x7	47	-4.8	23.04
x8	50	-1.8	3.24
x9	52	0.2	0.04
x10	57	5.2	27.04
x11	60	8.2	67.24
x12	65	13.2	174.24
x13	67	15.2	231.04
x14	85	33.2	1102.24
x15	86	34.2	1169.64
Mittelwert	51.8	0	322.0266667
Varianz	345.0285714		
Standardabweichung	18.57494472		

Tab. 3 Werte, Abweichungen vom Mittelwert und deren Quadrat als Veranschaulichung der Streuungswerte

Literaturhinweise

Hüsler, Jürg und Heinz Zimmermann (2006): *Statistische Prinzipien für medizinische Projekte*. 4. Auflage. Bern: Hans Huber.

Schlobinski, Peter (1996): *Empirische Sprachwissenschaft*. Opladen: Westdeutscher Verlag.

Woods, Anthony et al. (1986): *Statistics in language studies*. Cambridge: Cambridge University Press. (= Cambridge textbooks in linguistics) (versch. reprints)

3 Statistische Tests

Statistische Tests dienen dazu Hypothesen zu überprüfen, d.h. zu überprüfen, ob eine Behauptung wahr oder falsch ist. Die richtige Formulierung von Hypothesen ist deshalb ein wesentlicher Teil des wissenschaftlichen Arbeitens und stellt einen wichtigen Teil der Operationalisierung von Fragestellungen dar. Beispiel 3 zeigt diese Schritte vom Thema zur Fragestellung und zur Hypothese.

Beispiel 3

Das Thema des Kurses ist Sprechgeschwindigkeit. Unter diesem Themenbereich interessiert uns die Frage, wie Sprecher ihre Aussprache ändern, wenn sie einen schon gelesenen Text nochmals mit erhöhter Geschwindigkeit lesen. Diese allgemeine Frage lässt sich weiter eingrenzen, z.B. auf die Frage, wie sich die Pausen ändern, wie sich die Lautqualität oder die Länge der Laute ändert. Mit den hier schon gegebenen Daten grenzen wir die Frage ein, wie Sprecher die Vokale verändern. Aus dieser Frage lassen sich verschiedene Hypothesen ableiten:

- a) Vokale in normaler und schneller Geschwindigkeit sind unterschiedlich lang (zweiseitige Hypothese).
- b) Vokale in normaler Geschwindigkeit sind länger als solche in schneller Geschwindigkeit (einseitige Hypothese).

Für die statistische Berechnung wird zu diesen Hypothesen immer eine sog. Nullhypothese H_0 formuliert. Diese Nullhypothese entspricht dem, was man ablehnen möchte, hier also dass bei a) kein Unterschied der Vokallänge zwischen den Sprechgeschwindigkeiten besteht bzw. bei b) die Vokallänge bei normaler Sprechgeschwindigkeit nicht länger ist als bei schneller Geschwindigkeit. Die Alternativhypothese (häufig eben verkürzt Hypothese) entspricht dann dem, was bestätigt haben wollen.

In unserem Beispiel heißt die Nullhypothese also

$$H_0 \mu \text{ schnelle Geschwindigkeit} = \mu \text{ normale Geschwindigkeit}$$

Die Alternativhypothese lautet dann

$$a) H_1 \mu \text{ schnelle Geschwindigkeit} \neq \mu \text{ normale Geschwindigkeit}$$

$$b) H_1 \mu \text{ schnelle Geschwindigkeit} < \mu \text{ normale Geschwindigkeit}$$

Für die statistischen Tests verwendet man eine Testgröße, die aufgrund der Daten eine Entscheidung für die Null- oder Alternativhypothese ermöglicht. Diese Testgrößen sind in Tabellen zusammengefasst und in Statistikpaketen eingebaut, z. B. t-Test, z-Test, Wilcoxon-Test.

Zudem muss das Signifikanzniveau α festgelegt werden. Der kritische Bereich wird aus den unter der Nullhypothese nur mit geringer Wahrscheinlichkeit auftretenden Werten der Teststatistik gebildet. In der Sozialwissenschaft hat sich ein Signifikanzniveau von 0.05 als praktisch erwiesen, was aber nur ein willkürlich gewählter Wert ist. Ein Signifikanzniveau von 0.05 zeigt, dass die Wahrscheinlichkeit für einen Fehlentscheid 5 % beträgt.

Die Hypothesen lassen sich mittels verschiedener Testgrößen messen. Je nach Datenniveau muss ein anderer Test gewählt werden. Für den Vergleich zweier Gruppen wählt man meist den t-Test für intervallskalierte Daten und dem Chi-Quadratstest für nominalskalierte Daten. man spricht von parametrischen und nicht-parametrischen Tests, dabei stellen die parametrischen Tests höhere Anforderungen an die Daten als nicht-parametrische Tests. Parametrische Tests erwarten normalverteilte und zumindest intervallskalierte Daten, wohingegen beide Bedingungen bei nicht-

parametrischen Daten nicht gelten. Der am häufigsten verwendete Test bei intervallskalierten Daten ist der t-Test. Dieser ist relativ robust gegenüber Abweichungen von der Normalverteilung.

3.1 t-Test

Zu unterscheiden ist, ob korrelierende oder unabhängige Testgruppen vorliegen, d. h. ob wir zwei Messungen bei den selben Probanden, bzw. hier Lauten vornehmen, oder ob diese nichts miteinander zu tun haben.

3.1.1 t-Test für korrelierende Stichproben

Ausgangspunkt für unsere Überlegungen können die Daten aus Tabelle 2 sein. Sie stellen jeweils dieselben Laute im selben Lesetext bei normaler und erhöhter Sprechgeschwindigkeit dar. Die Daten sind also korreliert.

Für die Durchführung des t-Test muss die Differenz der Ergebnisse berechnet werden, sowie das Quadrat dieser Differenz. Dann wird die Standardabweichung (mit der vollen Formel) berechnet und der Standardfehler des Mittelwerts davon abgeleitet.

Belegstelle	a		D (Differenz)	D ² (Differenz ²)
	normal	schnell		
a1	48	86	-38	1444
a2	65	60	5	25
a3	63	57	6	36
a4	37	22	15	225
a5	66	36	30	900
a6	52	50	2	4
a7	82	67	15	225
a8	115	85	30	900
a9	112	31	81	6561
a10	69	47	22	484
a11	69	41	28	784
a12	71	65	6	36
a13	73	52	21	441
a14	62	42	20	400
a15	43	36	7	49
Summe			250	12514
Mittelwert	X = 68.466	Y = 51.8	16.6	
n	15			

Tab. 4: Dauer von /a/ in ms in normaler und schneller Geschwindigkeit sowie die für den t-Test notwendigen Berechnungen

Berechnung der **Standardabweichung s SD**:

$$s = SDD = \sqrt{\frac{\sum D^2 - \frac{1}{n}(\sum D)^2}{n-1}} = \sqrt{\frac{12514 - \frac{1}{15}(250)^2}{15-1}} = \sqrt{\frac{8347.44}{14}} = \sqrt{596.23} = 24.41$$

In Excel: STABW (Matrix)

Berechnung des **Standardfehlers** (Standard Error of Deviation)

$$SED = \frac{SDD}{\sqrt{n}} = \frac{24.41}{\sqrt{15}} \approx \frac{24.41}{3.87} \approx 6.3$$

Einsetzen der Ergebnisse in die Formel für den T-Test

$$t = \frac{\bar{X} - \bar{Y}}{SED} = \frac{68.466 - 51.8}{6.3} = \frac{16.666}{6.3} = 2.645$$

Um zu erkennen, ob der Wert signifikant ist, wird er in der t-Werte-Tabelle überprüft. Wir benötigen dazu den t-Wert und df (Freiheitsgrad). Der Freiheitsgrad entspricht in abhängigen Stichproben $n-1$, hier also 14. In der Tabelle sehen wir bei df 14 einen den nächstkleineren Wert zu 2.645 bei 2.624, davon abgeleitet ist der Wert der Wahrscheinlichkeit $=p$. In Excel lässt sich das mit der Formel $=TTEST(Matrix1; Matrix2;a;b)$ berechnen ($a=1$ für einen einseitigen Test, wir erwarten, dass schnellere Sprechgeschwindigkeit kürzere Laute liefert, $b=1$ für korrelierende Stichproben).

3.1.1 t-Test für unabhängige Stichproben

Entsprechend lässt sich ein t-Test auch für unabhängige Stichproben durchführen. Wir wollen wissen ob sich die Veränderung von a von der von o unterscheidet. Hypothese: beim Schnellen-Sprechen verkürzt sich /ɔ/ gleich wie /a/

Belegstelle	D von a	D ² von a	D von ɔ	D ² von ɔ
1	-38	1444	44	1936
2	5	25	2	4
3	6	36	47	2209
4	15	225	23	529
5	30	900	17	289
6	2	4	12	144
7	15	225	12	144
8	30	900	26	676
9	81	6561	27	729
10	22	484	33	1089
11	28	784	27	729
12	6	36	27	729
13	21	441	44	1936
14	20	400		
15	7	49		
Summe	250	12514	297	9207
Mittelwert	16.66		24.75	
n	15		13	

Tab. 5: Dauerdifferenz von /a/ und /ɔ/ in ms in normaler und schneller Geschwindigkeit sowie die für den t-Test notwendigen Berechnungen

Berechnung der Summe der Quadrate (Sum of squares)

$$SSx = \sum X^2 - \frac{(\sum X)^2}{n_1} \quad \text{und} \quad SSy = \sum Y^2 - \frac{(\sum Y)^2}{n_2}$$

$$t = \frac{\bar{X} - \bar{Y}}{\sqrt{\frac{SSx + SSy}{n_1 + n_2 - 2} \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}}$$

Wieder wird der Wert in der t-Tabelle nachgesehen, dabei gilt $df = n_1 + n_2 - 2$.

In Excel lässt sich ein T-Test mit folgender Formel erstellen $=TTEST(Matrix1; Matrix2; a;b)$, wobei $a=1$ einen einseitigen, $a=2$ einen zweiseitigen Test bildet, $b=1$ einen Test für korrelierende Stichproben, $b=2$ einen Test mit gleicher Varianz beider Stichproben und $b=3$ einen Test für unabhängige Stichproben darstellt. Das Ergebnis bieten den p-Wert, die Wahrscheinlichkeit für eine fälschliche Annahme der Alternativhypothese.

Tab. 9.1: Kritische Werte t_p des t -Testes mit ν Freiheitsgraden. Für $\nu = \infty$ gilt $t_p = z_p$, d.h. das p -Quantil der t -Verteilung mit ∞ Freiheitsgraden ist gleich dem p -Quantil der Normalverteilung.

ν	p									
	.60	.75	.90	.95	.975	.99	.995	.9975	.999	.9995
1	.325	1.000	3.078	6.314	12.706	31.821	63.657	127.32	318.31	636.62
2	.289	0.816	1.886	2.920	4.303	6.965	9.925	14.089	23.326	31.598
3	.277	0.765	1.638	2.353	3.182	4.541	5.841	7.453	10.213	12.924
4	.271	0.741	1.533	2.132	2.776	3.747	4.604	5.598	7.173	8.610
5	.267	0.727	1.476	2.015	2.571	3.365	4.032	5.773	5.893	6.869
6	.265	0.718	1.440	1.943	2.447	3.143	3.707	4.317	5.208	5.959
7	.263	0.711	1.415	1.895	2.365	2.998	3.499	4.029	4.785	5.408
8	.262	0.706	1.397	1.860	2.306	2.896	3.355	3.833	4.501	5.041
9	.261	0.703	1.383	1.833	2.262	2.821	3.250	3.690	4.297	4.781
10	.260	0.700	1.372	1.812	2.228	2.764	3.169	3.581	4.144	4.587
11	.260	0.697	1.363	1.796	2.201	2.718	3.106	3.497	4.025	4.437
12	.259	0.695	1.356	1.782	2.179	2.681	3.055	3.428	3.930	4.318
13	.259	0.694	1.350	1.771	2.160	2.650	3.012	3.372	3.852	4.221
14	.258	0.692	1.345	1.761	2.145	2.624	2.977	3.326	3.787	4.140
15	.258	0.691	1.341	1.753	2.131	2.602	2.947	3.286	3.733	4.073
16	.258	0.690	1.337	1.746	2.120	2.583	2.921	3.252	3.686	4.015
17	.257	0.689	1.333	1.740	2.110	2.567	2.898	3.222	3.646	3.965
18	.257	0.688	1.330	1.734	2.101	2.552	2.878	3.197	3.610	3.922
19	.257	0.688	1.328	1.729	2.093	2.539	2.861	3.174	3.579	3.883
20	.257	0.687	1.325	1.725	2.086	2.528	2.845	3.153	3.552	3.850
21	.257	0.686	1.323	1.721	2.080	2.518	2.831	3.135	3.527	3.819
22	.256	0.686	1.321	1.717	2.074	2.508	2.819	3.119	3.505	3.792
23	.256	0.685	1.319	1.714	2.069	2.500	2.807	3.104	3.485	3.767
24	.256	0.685	1.318	1.711	2.064	2.492	2.797	3.091	3.467	3.745
25	.256	0.684	1.316	1.708	2.060	2.485	2.787	3.078	3.450	3.725
26	.256	0.684	1.315	1.706	2.056	2.479	2.779	3.067	3.435	3.707
27	.256	0.684	1.314	1.703	2.052	2.473	2.771	3.057	3.421	3.690
28	.256	0.683	1.313	1.701	2.048	2.467	2.763	3.047	3.408	3.674
29	.256	0.683	1.311	1.699	2.045	2.462	2.756	3.038	3.396	3.659
30	.256	0.683	1.310	1.697	2.042	2.457	2.750	3.030	3.385	3.646
40	.255	0.681	1.303	1.684	2.021	2.423	2.704	2.971	3.307	3.551
60	.254	0.679	1.296	1.671	2.000	2.390	2.660	2.915	3.232	3.460
120	.254	0.677	1.289	1.658	1.980	2.358	2.617	2.860	3.160	3.373
200	.254	0.676	1.286	1.653	1.972	2.345	2.601	2.839	3.131	3.340
∞	.253	0.674	1.282	1.645	1.960	2.326	2.576	2.807	3.090	3.291