

Wichtige statistische Koeffizienten und Formeln

Skalenniveau	Art des Koeffizienten	Koeffizient	Formel
(ab) Nominalskala: (Unterscheidung nach Gleichheit/Ungleichheit; jeder Ausprägung eine Zahl zugeordnet, die keine quantitative Aussage hat) Beispiele: Geschlecht, Religion, Nationalität dichotome Variablen: genau 2 Ausprägungen	Lagemaß (Beschreibung von univariaten Verteilungen)	Modalwert	am häufigsten vorkommende Ausprägung
	Streuungsmaß (Beschreibung von univariaten Verteilungen)	Spannweite/range	Maximum-Minimum
	Zusammenhangsmaße (Kennwerte für bivariate Verteilungen)	„Chi-Quadrat“ (χ) (der Wert steigt mit wachsendem N → Chi ² -Tabelle Freiheitsgrade/dF in einer Kreuzta- belle: $(r-1) \cdot (c-1)$)	$Residuum = n_{ij} - n_{0ij}$ $stand. Residuum = \frac{Residuum}{\sqrt{n_{0ij}}}$ Vergleichbarkeit: wenn Wert nicht innerhalb [-2;2] → überzufällige Abweichung $Chi^2 = \sum_{i=1}^r \sum_{j=1}^c \frac{(n_{ij} - n_{0ij})^2}{n_{0ij}}$ doppeltes Summenzeichen: jede Kombination von r = row = Zeilenanzahl und c = column = Spaltenan- zahl genau einmal
wichtige Werte für eine Kreuz- tabelle: beobachtete Häufigkeit (n_{ij} , wobei i = Zeile, j = Spalte), erwartete Häufigkeit (n_{0ij} : Zei- lensumme mal Spaltensumme geteilt durch Gesamtzahl, auch: Indifferenztafel), relative Häufigkeit (Anteil der Häufigkeit einer Zelle an der Gesamtanzahl), Spaltenprozent (Anteil der Häufigkeit einer Zelle an ihrer Spalte)	Prozentsatzdifferenz (d%) (hauptsächlich für 2x2-Tabellen) 0 = kein Zusammenhang 100 = vollständiger Zusammenhang	d% = Spaltenprozent 1 - Spaltenprozent 2	

		<p>phi (Φ) (für 2x2-Tabellen)</p>	<p>ungerichtet (0 bis 1): $\Phi = \sqrt{\frac{Chi^2}{N}}$ gerichtet (-1 bis 1): $\Phi = \frac{n_{11} \cdot n_{22} - n_{12} \cdot n_{21}}{\sqrt{n_{1.} \cdot n_{2.} \cdot n_{.1} \cdot n_{.2}}}$ $n_{1.}$: Randsumme der ersten Zeile (usw.)</p>
		<p>C (0 bis <1) (für beliebige Tabellengrößen, aber Werte von Tabellengröße abhängig! → C_{korr})</p>	$C = \sqrt{\frac{Chi^2}{N + Chi^2}}$ $C_{max} = \sqrt{\frac{\min(r-1; s-1)}{\min(r; s)}}$ $C_{korr} = \frac{C}{C_{max}}$
		<p>V (0 bis 1) (immer anwendbar)</p>	$V = \sqrt{\frac{Chi^2}{Chi_{max}^2}} = \sqrt{\frac{Chi^2}{N \cdot \min(r-1; c-1)}}$
<p>besonderes Zusammenhangsmaß: PRE („proportional reduction in error“) $PRE = \frac{E1 - E2}{E1}$ E1: Fehler nach erster Vorhersage (Kennwerte der AV), ohne Kenntnis der UV E2: Fehler mit Kenntnis der UV</p>	<p>lambda (λ) (Vorhersage anhand des Modus: E1 als alle Fälle, die vom Gesamt-Modus der AV nicht erfasst werden, E2 als alle Fälle, die von den partiellen Modi der AV in den einzelnen Ausprägungen der UV nicht erfasst werden)</p>	<p>Zeilenvariable (Randspalte) = AV: $\lambda_{r} = \frac{\sum_{j=1}^c \max n_j - \max n_i}{N - \max n_i}$ $\max n_j$ = Maxima der Spalten n_i = Randspalte symmetrisch: $\lambda_{s} = \frac{\sum_{j=1}^c \max n_j + \sum_{i=1}^r \max n_i - \max n_i - \max n_j}{2 \cdot N - \max n_i - \max n_j}$</p>	

<p>(ab) Ordinalskala: (Einordnung in Rangfolge; Rangfolge an Zahlen ersichtlich; Abstände zwischen Ausprägungen nicht interpretierbar)</p> <p>Beispiele: Altersgruppen, (Einkommens-)Schichten, Einstellungsskala („stimme voll zu“ bis „stimme gar nicht zu“)</p>	<p>Lagemaß (Beschreibung von univariaten Verteilungen)</p>	<p>Median (der Wert des „mittleren“ Falls; 50%-Teiler)</p>	<p><u>Urliste:</u> ungerades n: $\tilde{x} = x_{\frac{1}{2} \cdot (n+1)}$ gerades n: $\tilde{x} = \frac{1}{2} \cdot (x_{\frac{n}{2}} + x_{\frac{n}{2}+1})$ <u>gruppierte/klassierte Daten:</u> $\tilde{x} = Q_2 = U + \frac{\frac{1}{2} \cdot N - F_u}{F_m} \cdot K_b$ U: untere exakte Intervallgrenze (des Intervalls, indem die 50% der Fälle erreicht werden) N: Zahl der Fälle F_u: kumulierte Häufigkeit unterhalb unterer Intervallgrenze F_m: Häufigkeit im Intervall K_b: Intervallbreite (meistens 1)</p>
	<p>Streuungsmaß (Beschreibung von univariaten Verteilungen)</p>	<p>Interquartilsabstand (exakte Rechnung genau genommen nur bei > ordinalskalierten Variablen → Klassen)</p>	<p>$IQ = Q_3 - Q_1$ Q₁ ist der 25%-Teiler, Q₃ der 75%-Teiler für Q₃ bzw. Q₁ in obiger Formel das ½ durch ¾ bzw. ¼ ersetzen und bei der Auswahl des Intervalls die 50% durch 75% bzw. 25%</p>
	<p>Zusammenhangsmaße (Kennwerte für bivariate Verteilungen): Konkordanzmaße</p> <p>konkordante Paare (N_C): die Ausprägungen der zwei Fälle sind in beiden Variablen unterschiedlich, und zwar in die gleiche Richtung</p>	<p>tau_a (τ) (nur anwendbar ohne Ties, also fast nie)</p> <p>tau_b (-1 bis 1) (für quadratische Tabellen)</p> <p>tau_c (-1 bis 1) (für rechteckige Tabellen)</p>	<p>$tau_a = \frac{N_C - N_D}{N \cdot (N - 1) \cdot \frac{1}{2}}$</p> <p>$tau_b = \frac{N_C - N_D}{\sqrt{(N_C + N_D + T_X) \cdot (N_C + N_D + T_Y)}}$</p> <p>$tau_c = \frac{N_C - N_D}{\frac{1}{2} \cdot N^2 \cdot \frac{m-1}{m}}$ mit $m = \min(r; c)$</p>

	<p>diskordante Paare (N_D): die Ausprägungen der zwei Fälle sind in beiden Variablen unterschiedlich, und zwar in entgegengesetzten Richtungen</p> <p>Ties auf X (T_X): dieselbe Ausprägung auf der x-Achse (Spalten), unterschiedliche auf der y-Achse (Zeilen)</p> <p>Ties auf Y (T_Y): dieselbe Ausprägung auf der y-Achse (Zeilen), unterschiedliche auf der x-Achse (Spalten)</p> <p>Ties auf X und Y (T_{XY}): dieselbe Ausprägung auf beiden Achsen</p> <p>alles addiert: Gesamtanzahl aller möglichen Paare =</p> $N \cdot (N - 1) \cdot \frac{1}{2}$	<p>Somers' d (asymmetrisches Maß: hier wird gerichtet getestet, d.h. UV und AV bestimmen)</p>	$d_{yx} = \frac{N_C - N_D}{N_C + N_D + T_Y}$ <p>(Y=AV, X=UV)</p>
		<p>gamma (γ) (je weniger Ausprägungen, desto größer → steigt bei Recodierung, produziert Ungenauigkeiten)</p>	$gamma = \frac{N_C - N_D}{N_C + N_D}$
<p>PRE-Maß</p> $PRE = \frac{E1 - E2}{E1}$		<p>gamma als PRE-Maß ($E1 = 0,5 \cdot (N_C + N_D)$, $E2 = \min(N_C; N_D)$)</p>	$gamma(PRE) = \frac{N_C + N_D - 2 \cdot \min(N_C; N_D)}{N_C + N_D}$ $gamma(PRE) = \frac{ N_C - N_D }{N_C + N_D}$

<p>(ab) Metrische Skalen:</p> <p><i>Intervallskala:</i> („um wie viele Einheiten stärker ist eine Ausprägung?“; willkürlicher Nullpunkt, willkürliche Skaleneinheit; Abstände interpretierbar; gleiche Intervallgrößen zwischen den Ausprägungen)</p> <p>Beispiele: (additive) Indizes (z.B. Gewaltbereitschaft), Temperatur in °C</p> <p><i>Verhältnis-/Ratioskala:</i> („um wieviel größer ist eine Ausprägung?“; natürlicher Nullpunkt, willkürliche Skaleneinheit; Quotienten interpretierbar; Größenverhältnisse zwischen den Zahlen entsprechen Stärke der Merkmalsausprägungen)</p> <p>Beispiele: Körpergröße, Einkommen (in €), Temperatur in Kelvin</p>	<p>Lagemaß (Beschreibung von univariaten Verteilungen)</p>	<p>arithmetisches Mittel („Mittelwert“)</p>	<p>Grundformel (ungruppierte Daten):</p> $\bar{x} = \frac{1}{n} \cdot \sum_{i=1}^n x_i$	
	<p>Streuungsmaß (Beschreibung von univariaten Verteilungen)</p>	<p>Varianz/Standardabweichung</p>	<p>$Var = s^2$ Grundformel (in der Stichprobe):</p> $s = \sqrt{\frac{\sum_{i=1}^N (x_i - \bar{x})^2}{N - 1}}$	
		<p>Gestaltmaße</p>	<p>Schiefe (v) (größer 0 → linkssteil/rechtsschief, kleiner 0 → rechtssteil/linksschief)</p>	$v = \frac{1}{n} \cdot \sum_{i=1}^n \left(\frac{x_i - \bar{x}}{s} \right)^3$
			<p>Wölbung/Kurtosis (w) (größer 0 → steiler als Normalverteilung, kleiner 0 → flacher als NV)</p>	$w = \frac{1}{n} \cdot \sum_{i=1}^n \left(\frac{x_i - \bar{x}}{s} \right)^4$
		<p>Zusammenhangsmaße (Kennwerte für bivariate Verteilungen)</p>	<p>Korrelationskoeffizient r (-1 bis 1) (gibt den <u>linearen</u> Zusammenhang an) r² beschreibt den Anteil der Variation der AV, der durch die UV erklärt wird (siehe auch: r² als PRE-Maß)</p>	$r = \frac{\sum_{i=1}^N (x_i - \bar{x}) \cdot (y_i - \bar{y})}{\sqrt{\sum_{i=1}^N (x_i - \bar{x})^2 \cdot \sum_{i=1}^N (y_i - \bar{y})^2}}$ <p>aus z-Werten:</p> $r = \frac{1}{N} \cdot \sum_{i=1}^N z_{xi} \cdot z_{yi}$
		<p>Regression (Berechnung der Geraden, die den Verlauf einer Punktwolke möglichst gut beschreibt, Minimierung des Vorhersagefehlers; auch hier nur <u>linearer</u> Zusammenhang erfasst)</p>	<p>Allgemeine Formel für die Gerade: $y'_i = a_{yx} + b_{yx} \cdot x_i$</p> <p>Schnittpunkt mit der y-Achse: $a_{yx} = \bar{y} - b_{yx} \cdot \bar{x}$</p> <p>Steigung:</p>	

<p><i>Absolutskala:</i> (natürlicher Nullpunkt, natürliche Skaleneinheit)</p> <p>Beispiel: Anzahl von irgendetwas (z.B. Anzahl der Teilnehmer einer Studie)</p> <p>Z-Transformation (Standardisierung von Verteilungen):</p> $z_{xi} = \frac{x_i - \bar{x}}{s_x}$			$b_{yx} = \frac{\sum_{i=1}^N (y_i - \bar{y}) \cdot (x_i - \bar{x})}{\sum_{i=1}^N (x_i - \bar{x})^2}$
<p>Kombination aus nominalem und metrischem Skalenniveau</p>	<p>Zusammenhangsmaß</p>	<p>eta (η) (0 bis 1) (die Ausprägungen der nominalen UV werden Kolonnen genannt, es gibt k Kolonnen)</p>	$eta = \sqrt{\sum_{j=1}^k n_j \cdot (\bar{y}_j - \bar{y})^2}$ <p>n_j: Anzahl der Fälle in der j-ten Kolonne \bar{y}_j: Mittelwert der j-ten Kolonne</p>
	<p>PRE-Maß</p>	<p>eta² (dieselbe Aufteilung in erklärte und nicht erklärte Variation wie bei der Regression: E1 ist die Gesamtvariation, E2 die nicht erklärte Variation)</p> <p>erklärte Variation = externe Variation → zwischen Kolonnen nicht erklärte Variation = interne Variation → in den Kolonnen)</p>	<p>Gesamtvariation = erklärte Variation + n. e. Variat.</p> $\sum_{j=1}^k \sum_{i=1}^{n_j} (y_{ij} - \bar{y})^2 = \sum_{j=1}^k n_j \cdot (\bar{y}_j - \bar{y})^2 + \sum_{j=1}^k \sum_{i=1}^{n_j} (y_{ij} - \bar{y}_j)^2$ <p>doppeltes Summenzeichen: in jeder Kolonne wird jeder Wert einmal betrachtet (=alle Werte) $y_{ij} - \bar{y}_j$: Differenz jedes Wertes von seinem jeweiligen Kolonnenmittelwert</p> <p>es gilt: $\frac{E1 - E2}{E1} = eta^2$</p> <p>eta² ist also der Anteil der erklärten Variation an der Gesamtvariation und der Wert für das PRE-Maß</p>
	<p>PRE-Maß</p>	<p>Regression/r² als PRE-Maß (E1: Gesamtvariation der AV)</p> $E1 = \sum_{i=1}^N (y_i - \bar{y})^2$ <p>E2: Abweichungen von den Regressionswerten; nicht erklärte Variation</p> $E2 = \sum_{i=1}^N (y_i - y'_i)^2$	<p>Gesamtvariation = erklärte Variation + n. e. Variat.</p> $\sum_{i=1}^N (y_i - \bar{y})^2 = \sum_{i=1}^N (y'_i - \bar{y})^2 + \sum_{i=1}^N (y_i - y'_i)^2$ <p>y'_i: Wert bei x_i auf der Regressionsgeraden</p> <p>tatsächlich gilt: $\frac{E1 - E2}{E1} = r^2$</p> <p>r² ist also der Anteil der erklärten Variation an der Gesamtvariation und der Wert für das PRE-Maß</p>

Schließende Statistik Versuch, von den Werten einer Stichprobe auf die der ihr zugrunde liegenden Grundgesamtheit (die gesamte Population, über die eine Aussage gemacht wird) zu schließen	Standard-/Stichprobenfehler	Anteilswerte (p) (Bsp.: 30% aller Jugendlichen in Bayern haben Erfahrung mit Cannabis)	$SF = \sqrt{\frac{p \cdot (1-p)}{n}}$ n: Stichprobenumfang
		arithmetisches Mittel	$SF = \frac{s}{\sqrt{n-1}}$
	Konfidenzintervall (in diesem Intervall um den wahren Wert der GG liegen bei gegebener Irrtumswahrscheinlichkeit die Werte aller möglichen Stichproben; durch Umdrehen der Logik testet man rückwärts auf die GG)	Anteilswerte	$KI = p \pm t_{dF; 1-\frac{\alpha}{2}} \cdot \sqrt{\frac{p \cdot (1-p)}{n}}$ t: Wert aus der Tabelle der Quantile der t-Verteilung $dF = n - 1$; bei $n > 800 \rightarrow$ z-Werte α : Irrtumswahrscheinlichkeit bei Ablehnung von H_0 , standardmäßig 0,05, also 5%
		arithmetisches Mittel	$KI = \bar{x} \pm \frac{s}{\sqrt{n-1}} \cdot t_{dF; 1-\frac{\alpha}{2}}$
	notwendige Stichprobengröße von Zufallsstichproben	Variante 1: aus Anteilswert und Konfidenzintervall	$n = \frac{4 \cdot z_{1-\frac{\alpha}{2}}^2 \cdot p \cdot (1-p)}{d^2}$ $d = 2 \cdot e$, e: Stichprobenfehler, d: Breite des KI α ist standardmäßig 0,05
		Variante 2: aus Größe der GG	<u>Größe der GG > 100.000:</u> $n = \left(\frac{z_{\frac{\alpha}{2}}}{e}\right)^2 \cdot p \cdot (1-p)$ α wird standardmäßig auf 0,05 gesetzt, e auf 0,03, p auf 0,5 (ungünstigster Fall) $\rightarrow n = 1067,11$ als standardisierter Wert <u>Größe der GG < 100.000:</u> $n = \frac{N \cdot z_{\frac{\alpha}{2}}^2 \cdot p \cdot (1-p)}{z_{\frac{\alpha}{2}}^2 \cdot p \cdot (1-p) + N \cdot e^2}$

