

SE Geographie und Ökonomie

Einheit 4: Bivariate Datenanalyse

Bernhard Schmidpeter
bernhard.schmidpeter@jku.at

Institut für Volkswirtschaftslehre

SoSe 2022

Motivation

- Oft ist man nicht an einem Merkmal interessiert, sondern an mehreren Merkmalen bzw. **Zusammenhängen zwischen verschiedenen Merkmalen**
- Wir untersuchen in dieser Einheit den **Zusammenhang zwischen zwei Merkmalen**. Zusammenhänge zwischen zwei Merkmalen (etwa: Unterschiede zwischen Gruppen) sind auch von wirtschaftspolitischem Interesse.
- **Anwendungsbeispiele:**
 - ▶ Arbeiten Frauen häufiger Teilzeit als Männer?
 - ▶ Verdienen auch Vollzeit beschäftigte Frauen weniger als (Vollzeit beschäftigte) Männer?
 - ▶ Wohnen reichere Leute in größeren Wohnungen?

Lernziele der Einheit 4

- Sie können bestimmen, **ob es einen Zusammenhang** zwischen zwei Merkmalen gibt.
- Abhängig davon, ob es sich um **diskrete oder stetige Merkmale** handelt, wissen Sie, **welche Maßzahlen geeignet** sind, einen möglichen Zusammenhang zwischen zwei Merkmalen abzubilden.
- Die können diese Maßzahlen **in EXCEL berechnen** und richtig **interpretieren**.
- Sie können beurteilen, ob dieser **Zusammenhang statistisch signifikant** ist. Das bedeutet, dass es diesen Zusammenhang mit hoher Wahrscheinlichkeit **auch in der Grundgesamtheit** gibt.

Zusammenhänge zwischen zwei Variablen

Was bedeutet Zusammenhang?

- Merkmal x **beeinflusst** Merkmal y : Wenn ich meine Arbeitszeit erhöhe, dann steigt mein Jahreseinkommen.
- Merkmal x und Merkmal y hängen zusammen (d.h. sie sind **nicht unabhängig voneinander**): Personen mit höheren Einkommen wohnen in größeren Wohnungen. Kausalität in beide Richtungen denkbar. Möglich, dass eine Dritte Variable (z.B.: Vermögen der Eltern) beide Variablen beeinflusst.
- Merkmal x beinhaltet Informationen über Merkmal y : Leute mit längeren Beinen sind üblicherweise größer.
- Üblicherweise sind wir an **kausalen Wirkungszusammenhängen** interessiert. Das ist aber oft sehr schwierig festzustellen (wir werden später darauf zurückkommen).

Zusammenhänge zwischen zwei Variablen

Wie der Zusammenhang zwischen zwei Variablen untersucht werden kann, hängt davon ab, ob es sich um **diskrete oder stetige metrische Merkmale** handelt:

		Variable 2	
		diskret	stetig
Variable 1	diskret	Kreuztabelle	Mittelwertvergleich
	stetig	Mittelwertvergleich	Korrelation

Mittelwertvergleich: Beispiel

Haben Frauen und Männer (diskretes Merkmal) im Durchschnitt unterschiedliche (Haushalts-)Einkommen (stetiges Merkmal)?

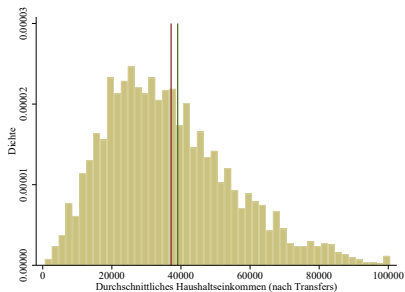
- **Nullhypothese (H_0): In der Grundgesamtheit** gibt es **keinen Unterschied** im Haushaltseinkommen zwischen den beiden Gruppen: $\mu_M = \mu_F$ (allgemein: $\mu_1 = \mu_2$)
- **Alternativhypothese (H_1): In der Grundgesamtheit** gibt es **einen Unterschied** im Haushaltseinkommen zwischen den beiden Gruppen: $\mu_M \neq \mu_F$ (allgemein: $\mu_1 \neq \mu_2$)
- Wir führen dazu einen **zweiseitigen Zweistichproben-t-Test für unabhängige Stichproben** durch.

Die Nullhypothese wird daher verworfen, wenn Männer ein signifikant höheres Haushaltseinkommen haben als Frauen **oder** wenn Männer ein signifikant niedrigeres Haushaltseinkommen haben als Frauen.

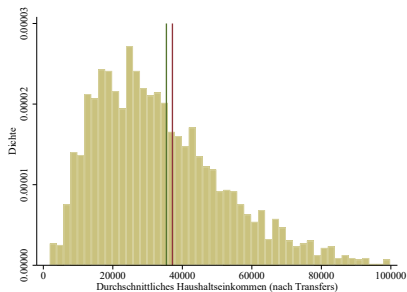
(Haushalts-)Einkommensverteilung von Männern und Frauen

	n	\bar{x}	s	x_{min}	x_{max}
HH-Einkommen (alle)	5,407	37,149.97	26,727.97	583	507,369
HH-Einkommen Männer	2,583	39,044.55	28,397.41	583	507,369
HH-Einkommen Frauen	2,824	35,417.08	24,983.54	1,809	507,369

Männer



Frauen



Zweiseitiger Zweistichproben-t-Test

Nullhypothese (H_0): $\mu_M = \mu_F \Rightarrow \mu_M - \mu_F = 0$

Alternativhypothese (H_1): $\mu_M \neq \mu_F \Rightarrow \mu_M - \mu_F \neq 0$

- Wir **wissen von letzter Einheit**, dass:

- ▶ $\bar{x}_M \sim N(\mu_M, \sigma_M^2/n_M)$
- ▶ $\bar{x}_F \sim N(\mu_F, \sigma_F^2/n_F)$

- Neu: Summen und Differenzen von normalverteilten Zufallsvariablen folgen ebenfalls einer Normalverteilung:

- ▶ $\bar{x}_M + \bar{x}_F \sim N(\mu_M + \mu_F, \sigma_M^2/n_M + \sigma_F^2/n_F)$
- ▶ $\bar{x}_M - \bar{x}_F \sim N(\mu_M - \mu_F, \sigma_M^2/n_M + \sigma_F^2/n_F)$

Zweiseitiger Zweistichproben-t-Test

- Wir können die Konfidenzintervalle von Differenzen dann ähnlich berechnen wie in unserer letzten Einheit
- Für ein $(1 - \alpha)$ Konfidenzintervall haben wir, wenn σ_M^2 und σ_F^2 bekannt sind:

$$\bar{x}_M - \bar{x}_F - u_{1-\frac{\alpha}{2}} \cdot \sqrt{\frac{\sigma_M^2}{n_M} + \frac{\sigma_F^2}{n_F}} \leq \underbrace{\mu_M - \mu_F}_{=0 \text{ unter } H_0} \leq \bar{x}_M - \bar{x}_F + u_{1-\frac{\alpha}{2}} \cdot \sqrt{\frac{\sigma_M^2}{n_M} + \frac{\sigma_F^2}{n_F}}$$

- Wenn σ_M^2 und σ_F^2 unbekannt sind, ersetzen wir diese durch die geschätzte Varianz $\hat{\sigma}_M^2$ und $\hat{\sigma}_F^2$ (Standardfall):

$$\bar{x}_M - \bar{x}_F - u_{1-\frac{\alpha}{2}} \cdot \sqrt{\frac{\hat{\sigma}_M^2}{n_M} + \frac{\hat{\sigma}_F^2}{n_F}} \leq \underbrace{\mu_M - \mu_F}_{=0 \text{ unter } H_0} \leq \bar{x}_M - \bar{x}_F + u_{1-\frac{\alpha}{2}} \cdot \sqrt{\frac{\hat{\sigma}_M^2}{n_M} + \frac{\hat{\sigma}_F^2}{n_F}}$$

Zweiseitiger Zweistichproben-t-Test

Standardisierte Teststatistik:

$$\bar{x}_M - \bar{x}_F - u_{1-\frac{\alpha}{2}} \cdot \sqrt{\frac{\hat{s}_M^2}{n_M} + \frac{\hat{s}_F^2}{n_F}} \leq \underbrace{0}_{\mu_M - \mu_F} \leq \bar{x}_M - \bar{x}_F + u_{1-\frac{\alpha}{2}} \cdot \sqrt{\frac{\hat{s}_M^2}{n_M} + \frac{\hat{s}_F^2}{n_F}}$$

$$u_{1-\frac{\alpha}{2}} \cdot \sqrt{\frac{\hat{s}_M^2}{n_M} + \frac{\hat{s}_F^2}{n_F}} \leq \bar{x}_M - \bar{x}_F \leq u_{1-\frac{\alpha}{2}} \cdot \sqrt{\frac{\hat{s}_M^2}{n_M} + \frac{\hat{s}_F^2}{n_F}}$$

$$-u_{1-\frac{\alpha}{2}} \leq \frac{\bar{x}_M - \bar{x}_F}{\sqrt{\frac{\hat{s}_M^2}{n_M} + \frac{\hat{s}_F^2}{n_F}}} \leq u_{1-\frac{\alpha}{2}}$$

$$\left| \frac{\bar{x}_M - \bar{x}_F}{\sqrt{\frac{\hat{s}_M^2}{n_M} + \frac{\hat{s}_F^2}{n_F}}} \right| \leq u_{1-\frac{\alpha}{2}}$$

Hinweis: Normalerweise verwenden wir den kritischen Wert der t-Verteilung

$$\text{t-Wert} = \left| \frac{\bar{x}_M - \bar{x}_F}{\sqrt{\frac{\hat{s}_M^2}{n_M} + \frac{\hat{s}_F^2}{n_F}}} \right| \leq t_{n_M+n_F-2; 1-\frac{\alpha}{2}}$$

Beispiel

Standardisierte Teststatistik:

H_0 wird nicht verworfen, wenn:

$$\text{t-Wert} = \left| \frac{\bar{X}_M - \bar{X}_F}{\sqrt{\frac{\hat{s}_M^2}{n_M} + \frac{\hat{s}_F^2}{n_F}}} \right| \leq t_{n_M+n_F-2; 1-\frac{\alpha}{2}} = \text{kritischer Wert } c$$

$$\text{t-Wert} = \left| \frac{39,045 - 35,417}{\sqrt{\frac{28,397^2}{2,583} + \frac{24,984^2}{2,824}}} \right| = \frac{3,627}{730} = 4.97$$

$$\text{kritischer Wert } c = t_{n_M+n_F-2; 1-\frac{\alpha}{2}} = t_{5,405; 0.975} = 1.96$$

$$\text{t-Wert} = 4.97 > 1.96 = \text{kritischer Wert } c$$

- **Nullhypothese wird** zugunsten der Alternativhypothese **verworfen**
- Es ist sehr wahrscheinlich, dass sich in der Grundgesamtheit die Haushaltseinkommen von Männern und Frauen unterscheiden

Einseitiger Zweistichproben-t-Test: Vorgehensweise

1. Berechnung einer standardisierten Teststatistik (t-Wert) und krit. Wert

- ▶ Teststatistik:

$$t\text{-Wert} = \left| \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{\frac{\hat{s}_1^2}{n_1} + \frac{\hat{s}_2^2}{n_2}}} \right|$$

- ▶ Ermittlung des kritischen Werts c der Verteilung: $c = t_{n_1+n_2-2;1-\alpha}$

2. Bewertung der Teststatistik

- ▶ Fall 1: $H_0 : \mu_1 \leq \mu_2$; $H_1 : \mu_1 > \mu_2$

- ★ H_0 wird (zugunsten von H_1) verworfen, wenn

$$t\text{-Wert} = \left| \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{\frac{\hat{s}_1^2}{n_1} + \frac{\hat{s}_2^2}{n_2}}} \right| \geq t_{n_1+n_2-2;1-\alpha} \text{ und } \bar{x}_1 > \bar{x}_2$$

- ▶ Fall 2: $H_0 : \mu_1 \geq \mu_2$; $H_1 : \mu_1 < \mu_2$

- ★ H_0 wird (zugunsten von H_1) verworfen, wenn

$$t\text{-Wert} = \left| \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{\frac{\hat{s}_1^2}{n_1} + \frac{\hat{s}_2^2}{n_2}}} \right| \geq t_{n_1+n_2-2;1-\alpha} \text{ und } \bar{x}_1 < \bar{x}_2$$

Umsetzung in EXCEL

Es bietet sich an, die Variablen der beiden Gruppen in getrennte Spalten zu kopieren. Alle notwendigen Befehle habe wir in Einheit 3 schon besprochen.

Alternativ: **Daten** → **Analyse** → **Datenanalyse** → **Zweistichproben t-Test: Unterschiedlicher Varianzen**

	hhinc Männer	hhinc Frauen
Mittelwert	39044.5451	35417.079
Varianz	806412884	624177386
Beobachtungen	2583	2824
Hypothetische Differenz der Mittelwerte	0	
Freiheitsgrade (df)	5165	
t-Statistik	4.96761202	
P(T<=t) einseitig	3.4984E-07	
Kritischer t-Wert bei einseitigem t-Test	1.6451487	
P(T<=t) zweiseitig	6.9968E-07	
Kritischer t-Wert bei zweiseitigem t-Test	1.96042339	

Anmerkung: Wenn sich die Varianzen stark unterscheiden (Faustregel: die Varianz einer Gruppe ist mehr als doppelt so groß wie die Varianz einer anderen Gruppe), dann sollten die Freiheitsgrade angepasst werden. Das Datenanalyse-Tool in Excel macht das automatisch, wodurch ein geringfügig anderer kritischer t-Wert berechnet wird. Für unsere Zwecke spielt das aber keine Rolle.

Zusammenhang zwischen zwei diskreten Merkmalen: Kreuztabelle

		Variable 2	
		diskret	stetig
Variable 1	diskret	Kreuztabelle	Mittelwertvergleich
	stetig	Mittelwertvergleich	Korrelation

Kreuztabelle

Eine **Kreuztabelle** (auch Kontingenztabelle oder Kontingenztafel) weist die absoluten oder die relativen Häufigkeiten aller Kombinationen der Merkmalsausprägungen von zwei diskreten Merkmalen aus (Merkmale müssen nicht notwendigerweise dichotom sein):

		Variable 2		Summe
		$y = 1$	$y = 2$	
Variable 1	$x = 1$	h_{11}	h_{12}	h_{1+}
	$x = 2$	h_{21}	h_{22}	h_{2+}
Summe		h_{+1}	h_{+2}	n

Bezeichnungen:

h_{ij}	absolute Häufigkeit der Kombination $x = x_i$ und $y = y_j$
n	Stichprobenumfang, $n = \sum_i \sum_j h_{ij}$
$p_{ij} = h_{ij}/n$	relative Häufigkeit der Kombination $x = x_i$ und $y = y_j$
$P_{ij} = p_{ij} \cdot 100$	relative Häufigkeit der Kombination $x = x_i$ und $y = y_j$ in %
$h_{i+}(p_{i+})$	Zeilensummen, Randhäufigkeiten des Merkmals x
$h_{+j}(p_{+j})$	Spaltensummen, Randhäufigkeiten des Merkmals y

Beispiel: Erwerbsstatus nach Geschlecht

Grundgesamtheit: Alle regulär Erwerbstätigen in Deutschland

Stichprobe: alle regulär Erwerbstätigen im GSOEP

Kreuztabelle: absolute Häufigkeiten

		Erwerbsstatus		Summe
		Vollzeit	Teilzeit	
Geschlecht	Männer	1,346	59	1,405
	Frauen	695	540	1,235
Summe		2,041	599	2,640

Kreuztabelle: relative Häufigkeiten

		Erwerbsstatus		Summe
		Vollzeit	Teilzeit	
Geschlecht	Männer	0.510	0.022	0.532
	Frauen	0.263	0.205	0.468
Summe		0.773	0.227	1.000

Randverteilungen: Gibt Auskunft über die Verteilung eines Merkmals, ohne die Verteilung des anderen Merkmals zu berücksichtigen.

Bedingte Wahrscheinlichkeiten

Unterscheidet sich der Anteil der Vollzeitbeschäftigten zwischen Männern und Frauen?

Bezeichnung:

$h_{ij}/h_{i+} = p_{ij}/p_{i+}$ bedingte relative Häufigkeit der Ausprägung y_j des Merkmals y bei gegebener Ausprägung x_i des Merkmals x

Bedingte relative Häufigkeiten

		Erwerbsstatus		Summe
		Vollzeit	Teilzeit	
Geschlecht	Männer	0.510/0.532	0.022/0.532	0.532/0.532
	Frauen	0.263/0.468	0.205/0.468	0.468/0.468

Bedingte relative Häufigkeiten

		Erwerbsstatus		Summe
		Vollzeit	Teilzeit	
Geschlecht	Männer	0.958	0.042	1.000
	Frauen	0.563	0.437	1.000

Stärke des Zusammenhangs

Idee: Um die Stärke des Zusammenhangs zu beurteilen, soll die **beobachtete Verteilung** mit jener Verteilung verglichen werden, die ich **erwarten** würden, wenn die beiden Merkmale keinen Zusammenhang aufweisen.

Bezeichnungen:

$p_{ij}^e = p_{i+} \cdot p_{+j}$	erwartete (e für expected) relative Häufigkeit von $x = x_i$ und $y = y_j$ bei Unabhängigkeit von x und y
$h_{ij}^e = p_{ij}^e \cdot n$ $= (h_{i+}^o \cdot h_{+j}^o) / n$	erwartete absolute Häufigkeit dieser Kombination bei Unabhängigkeit von x und y
h_{ij}^o	beobachtete (o für observed) absolute Häufigkeit dieser Kombination

Erwartete und beobachtete Häufigkeiten

Kreuztabelle: relative Häufigkeiten

		Erwerbsstatus		Summe
		Vollzeit	Teilzeit	
Geschlecht	Männer	$p_{11}^o = 0.510$ $p_{11}^e = 0.411$	$p_{12}^o = 0.022$ $p_{12}^e = 0.121$	$p_{1+}^o = 0.532$
	Frauen	$p_{21}^o = 0.263$ $p_{21}^e = 0.362$	$p_{22}^o = 0.205$ $p_{22}^e = 0.106$	$p_{2+}^o = 0.468$
Summe		$p_{+1}^o = 0.773$	$p_{+1}^e = 0.227$	1.000

Kreuztabelle: absolute Häufigkeiten

		Erwerbsstatus		Summe
		Vollzeit	Teilzeit	
Geschlecht	Männer	$h_{11}^o = 1,346$ $h_{11}^e = 1,086$	$h_{12}^o = 59$ $h_{12}^e = 319$	$h_{1+}^o = 1,405$
	Frauen	$h_{21}^o = 695$ $h_{21}^e = 955$	$h_{22}^o = 540$ $h_{22}^e = 280$	$h_{2+}^o = 1,235$
Summe		$h_{+1}^o = 2,041$	$h_{+1}^e = 599$	2,640

Maßzahlen

Das **Assoziationsmaß Chi-Quadrat** χ_{err}^2 (auch: Pearson's χ^2) mit

$$\chi_{err}^2 = \sum_i \sum_j \frac{(h_{ij}^o - h_{ij}^e)^2}{h_{ij}^e}$$

misst den Zusammenhang zwischen zwei diskreten Merkmalen.

Da das Assoziationsmaß χ_{err}^2 mit dem Stichprobenumfang steigt, bietet sich das **Cramersche Assoziationsmaß V** an:

$$V = \sqrt{\frac{\chi_{err}^2}{n \cdot (\min(r, s) - 1)}}$$

r gibt die Anzahl der Merkmalsausprägungen des Merkmals x an.

s gibt die Anzahl der Merkmalsausprägungen des Merkmals y an.

Es gilt $0 \leq V \leq 1$

Maßzahlen: Beispiel

$$\chi_{err}^2 = \sum_i \sum_j \frac{(h_{ij}^o - h_{ij}^e)^2}{h_{ij}^e} = \frac{(1,346 - 1,086)^2}{1,086} + \frac{(59 - 319)^2}{319} + \frac{(695 - 955)^2}{955} + \frac{(540 - 280)^2}{280} \approx 585$$

$$V = \sqrt{\frac{\chi^2}{n \cdot (\min(r,s) - 1)}} = \sqrt{\frac{585}{2,640 \cdot (\min(2,2) - 1)}} = \sqrt{\frac{585}{2,640}} \approx 0.471$$

Interpretationshilfe für das Cramersche Assoziationsmaß V:

$V = 0$	kein Zusammenhang
$0 < V \leq 0.3$	schwacher Zusammenhang
$0.3 < V \leq 0.7$	mittlerer Zusammenhang
$0.7 < V < 1$	starker Zusammenhang
$V = 1$	vollständiger Zusammenhang

Zusammenhang oder Zufall? χ^2 -Test auf Unabhängigkeit

Nullhypothese (H_0): Es gibt **keinen Zusammenhang** zwischen Geschlecht und Erwerbsstatus **in der Grundgesamtheit**

Alternativhypothese (H_1): Es gibt **einen Zusammenhang** zwischen Geschlecht und Erwerbsstatus **in der Grundgesamtheit**

Intuition:

Ein Zusammenhang in der Grundgesamtheit ist dann wahrscheinlich, wenn der Zusammenhang in der Stichprobe groß ist (V liegt deutlich über 0) und wenn der Stichprobenumfang groß ist (kleine Unsicherheit).

Teststrategie zielt auf χ_{err}^2 ab:

H_0 wird nicht verworfen, wenn gilt:

$$\chi_{(r-1)(s-1);1-\alpha}^2 \geq \chi_{err}^2 = \sum_i \sum_j \frac{(h_{ij}^o - h_{ij}^e)^2}{h_{ij}^e},$$

H_0 wird mit einer Irrtumswahrscheinlichkeit α verworfen, wenn gilt:

$$\chi_{(r-1)(s-1);1-\alpha}^2 < \chi_{err}^2 = \sum_i \sum_j \frac{(h_{ij}^o - h_{ij}^e)^2}{h_{ij}^e},$$

wobei $(r - 1) \cdot (s - 1)$ die Zahl der Freiheitsgrade bezeichnet.

χ^2 -Test: Beispiel

Teststrategie zielt auf χ_{err}^2 ab:

H_0 wird nicht verworfen, wenn gilt:

$$\chi_{err}^2 = \sum_i \sum_j \frac{(h_{ij}^o - h_{ij}^e)^2}{h_{ij}^e} \leq \chi_{(r-1)(s-1); 1-\alpha}^2$$

H_0 wird mit einer Irrtumswahrscheinlichkeit α verworfen, wenn gilt:

$$\chi_{err}^2 = \sum_i \sum_j \frac{(h_{ij}^o - h_{ij}^e)^2}{h_{ij}^e} > \chi_{(r-1)(s-1); 1-\alpha}^2$$

→ auch hier wird eine standardisierte Teststatistik mit einem kritischen Wert der Verteilung verglichen!

Beispiel:

$$\chi_{err}^2 = 585 > 3.84 = \chi_{1; 0.95}^2 = \chi_{(2-1)(2-1); 1-0.05}^2$$

→ **Die Nullhypothese** wird sehr deutlich **verworfen**: In der Grundgesamtheit gibt es (höchstwahrscheinlich) einen Zusammenhang zwischen Geschlecht und Erwerbsstatus.

Umsetzung in EXCEL

- 1 Erstellung der **Kreuztabelle der beobachteten absoluten Häufigkeiten**:
 - ▶ Hier bietet sich der **EXCEL-Befehl** ZÄHLENWENNS (*S* am Ende beachten!) an, wo die Zahl an Beobachtungen gezählt wird, die mehrere Bedingungen erfüllt (etwa *Männer* und *Vollzeitbeschäftigt*)
 - ▶ Nachdem alle absoluten Häufigkeiten h_{ij}^o erstellt wurden, können die Randverteilungen ermittelt werden.
- 2 Erstellung der **Kreuztabelle der erwarteten absoluten Häufigkeiten**:
 - ▶ Mit den Randverteilungen können die erwarteten absoluten Häufigkeiten h_{ij}^e berechnet werden, und in eine neue Kreuztabelle eingetragen werden.
- 3 Die einfachste Möglichkeit besteht darin, mit dem **EXCEL-Befehl** CHIUQ.TEST den p-Wert von χ_{err}^2 zu bestimmen. Dazu muss χ_{err}^2 gar nicht berechnet werden, sondern es ist ausreichend, die beobachteten und die erwarteten Häufigkeiten zu markieren. Ist dieser **p-Wert kleiner als** die angenommene Irrtumswahrscheinlichkeit α , dann wird die **Nullhypothese H_0 verworfen**.

Alternativ kann χ_{err}^2 berechnet werden und mit dem kritischen $\chi_{(r-1)(s-1);1-\alpha}^2$ -Wert verglichen werden. Den kritischen Wert erhält man mit **EXCEL-Befehl** CHIUQ.INV($1 - \alpha$; $(r - 1)(s - 1)$).

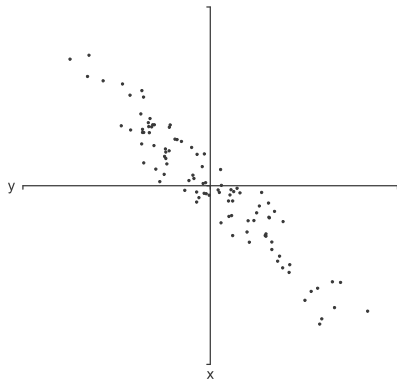
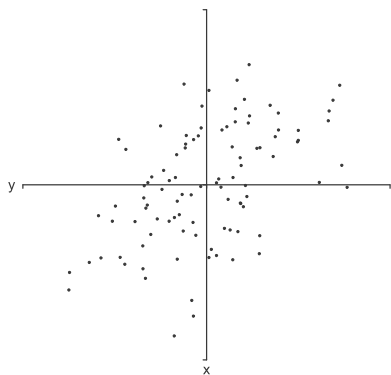
Zusammenhang zwischen zwei metrischen Variablen: Korrelation

		Variable 2	
		diskret	stetig
Variable 1	diskret	Kreuztabelle	Mittelwertvergleich
	stetig	Mittelwertvergleich	Korrelation

Der **(Bravais-Pearson-)Korrelationskoeffizient** ρ gibt an, (i) ob der Zusammenhang zwischen zwei metrischen Variablen x und y **positiv oder negativ** ist, und (ii) wie ähnlich dieser Zusammenhang einem linearen Zusammenhang ist.

Grafische Darstellung: Streudiagramm

Streudiagramm: Ein Streudiagramm ist eine grafische Darstellung eines zweidimensionalen metrischen Merkmals. Dabei wird jeder Erhebungseinheit der zugehörige Datenpunkt in einem Koordinatensystem zugeordnet. Streudiagramme erleichtern das Auffinden von Zusammenhängen.



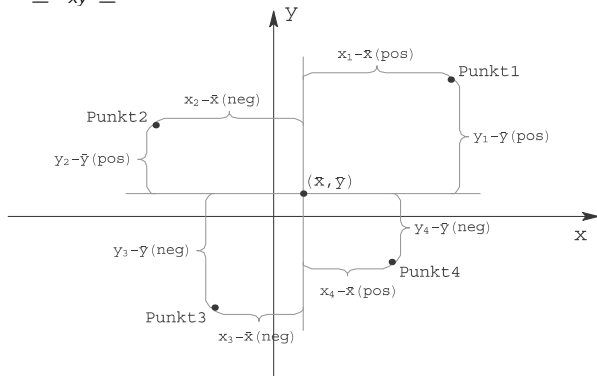
(siehe Duller, Abbildung 7.2)

Kovarianz als Ausgangspunkt

Die **Kovarianz** zu den Merkmalen x und y einer Stichprobe ist gegeben durch:

$$\hat{s}_{xy} = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x}) \cdot (y_i - \bar{y}) = \frac{1}{n-1} \sum_{i=1}^n (x_i \cdot y_i - \bar{x} \cdot \bar{y})$$

Wobei gilt: $-\infty \leq \hat{s}_{xy} \leq \infty$



(siehe Duller, Abbildung 7.1)

Korrelationskoeffizient: Berechnung

Der **Korrelationskoeffizient** stellt ein standardisiertes Maß zur Messung eines linearen Zusammenhangs zwischen zwei metrischen Merkmalen x und y dar:

$$\rho = \frac{\hat{s}_{xy}}{\hat{s}_x \cdot \hat{s}_y}$$

wobei \hat{s}_x und \hat{s}_y die Standardabweichung der Merkmale x und y darstellt.

Es gilt: $-1 \leq \rho \leq 1$

Anmerkungen:

- Es ist auch möglich, ein Konfidenzintervall von ρ zu berechnen und zu beurteilen, ob ρ statistisch signifikant von 0 verschieden ist. Da die statistische Signifikanz selten ausgewiesen wird und relativ aufwendig zu berechnen ist, wird dieses Thema nicht besprochen.
- Wenn die Grundgesamtheit (statt einer Stichprobe) vorliegt, ist die Korrektur (indem durch $n - 1$ statt n dividiert wird) bei Berechnung von Kovarianz und Standardabweichung nicht notwendig. In der Praxis ist der Unterschied aber meist gering und von geringer Relevanz.
- Bei der Berechnung des Korrelationskoeffizienten ρ ist es ohnehin unerheblich, ob diese Korrektur vorgenommen wird.

Korrelationskoeffizient: Interpretation

Interpretation:

$\rho > 0$ gleichsinniger (positiver) linearer Zusammenhang

$\rho = 0$ kein linearer Zusammenhang

$\rho < 0$ gegensinniger (negativer) linearer Zusammenhang

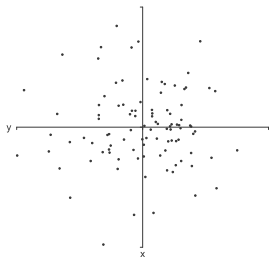
Die Richtung des Zusammenhanges ergibt sich aus dem **Vorzeichen**.

Interpretationshilfe für den Korrelationskoeffizienten:

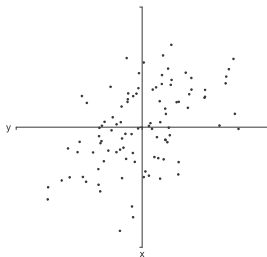
$\rho = 0$	kein Zusammenhang
$0 < \rho \leq 0.3$	schwacher Zusammenhang
$0.3 < \rho \leq 0.7$	mittlerer Zusammenhang
$0.7 < \rho < 1$	starker Zusammenhang
$ \rho = 1$	vollständiger Zusammenhang

Anmerkung: ρ^2 entspricht dem Bestimmtheitsmaß (R^2) bei einer linearen Einfachregression. Das bedeutet, dass ein Anteil ρ^2 der Variation eines Merkmals durch die Variation des zweiten Merkmals erklärt werden kann. Wir werden später genauer darauf zurückkommen.

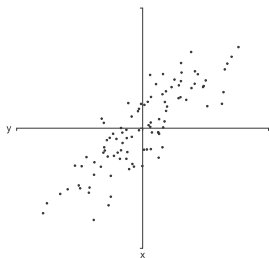
Korrelationskoeffizient: Beispiele



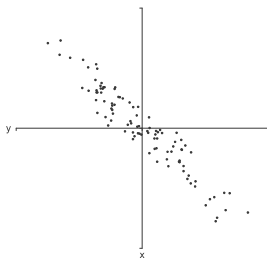
Korrelation $\rho = 0$



Korrelation $\rho = 0,5$



Korrelation $\rho = 0,85$



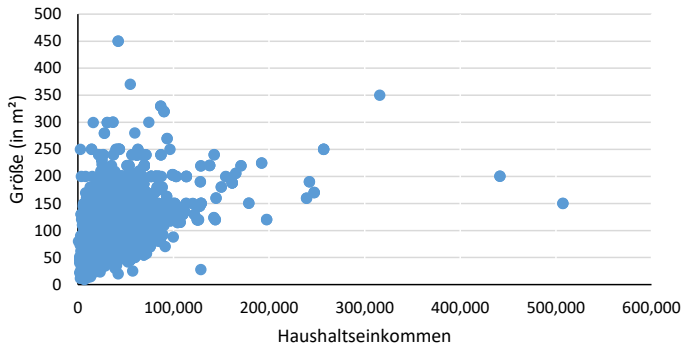
Korrelation $\rho = -0,95$

(siehe Duller, Abbildung 7.2)

Umsetzung in EXCEL (1)

Streudiagramm:

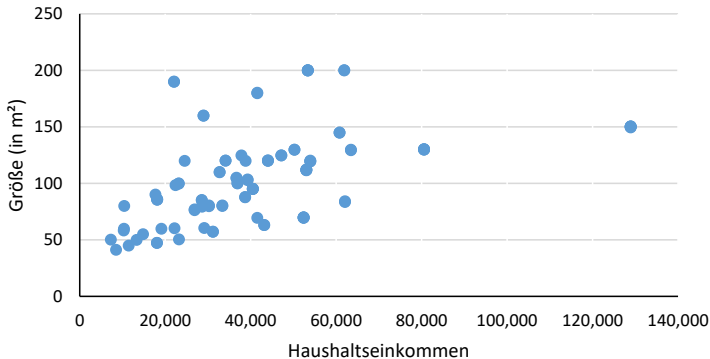
Mit **Einfügen** → **Diagramme** → **Punkt (XY)** kann ein **Streudiagramm** eingefügt werden. Streudiagramme sind oft unübersichtlich. Das kann damit zusammenhängen, (i) dass es **Ausreißer** gibt, und/oder (ii) dass der Stichprobenumfang n sehr groß ist (und es daher sehr viele Punkte im Streudiagramm gibt). Man kann auch nur eine (Zufalls-) Auswahl der Stichprobe für das Streudiagramm verwenden.



Umsetzung in EXCEL (2)

Streudiagramm:

Mit **Einfügen** → **Diagramme** → **Punkt (XY)** kann ein **Streudiagramm** eingefügt werden. Streudiagramm sind oft unübersichtlich. Das kann damit zusammenhängen, (i) dass es **Ausreißer** gibt, oder (ii) dass der Stichprobenumfang n sehr groß ist (und es daher sehr viele Punkte im Streudiagramm gibt). Man kann auch nur eine (Zufalls-) Auswahl der Stichprobe für das Streudiagramm verwenden.



Umsetzung in EXCEL (3)

- **Kovarianz:**

Mit den **Excel-Befehlen** KOVARIANZ.S und KOVARIANZ.P kann die **Kovarianz** zwischen zwei Merkmalen einer Stichprobe (".S") bzw. der Grundgesamtheit (".P") berechnet werden.

- **Korrelation:**

Mit den **Excel-Befehl** KORREL kann der **Korrelationskoeffizient** von zwei Merkmalen berechnet werden.

- ▶ Hier ist eine Unterscheidung in Stichprobe und Grundgesamtheit nicht notwendig, da sich $1/(n - 1)$ bzw. $1/N$ wegekürzt.

- Alternativ kann **Daten** → **Analyse** → **Datenanalyse** → **Kovarianz** bzw. **Daten** → **Analyse** → **Datenanalyse** → **Korrelation** verwendet werden.

- ▶ Vorteil: Gibt eine formatierte Tabelle zurück. Es können auch Kovarianzen bzw. Korrelationskoeffizienten von mehreren Merkmalen (bzw.: von mehreren Merkmals-Paaren) berechnet werden.
- ▶ Nachteil: "Eingabebereich muss ein zusammenhängender Bezug sein" (d.h. die Variablen müssen nebeneinander stehen).

Umsetzung in EXCEL (4)

Varianz-Kovarianz-Matrix:

	<i>persnr</i>	<i>ybirth</i>	<i>income</i>	<i>hhinc</i>	<i>size</i>	<i>Größe (in m2)</i>
<i>persnr</i>	9, 588, 648, 723, 916					
<i>ybirth</i>	-11, 977, 831	329				
<i>income</i>	206, 952, 278	52, 743	1, 400, 149, 494			
<i>hhinc</i>	-4, 608, 430, 303	74, 305	479, 177, 712	714, 252, 359		
<i>size</i>	-99, 474, 892	680	3, 073, 259	5, 895, 795	211, 692	
<i>Größe (in m2)</i>	-9, 241, 443	63	285, 513	547, 733	19, 667	1, 827

Korrelationsmatrix:

	<i>persnr</i>	<i>ybirth</i>	<i>income</i>	<i>hhinc</i>	<i>size</i>	<i>Größe (in m2)</i>
<i>persnr</i>	1					
<i>ybirth</i>	-0.213	1				
<i>income</i>	0.002	0.081	1			
<i>hhinc</i>	-0.056	0.153	0.483	1		
<i>size</i>	-0.070	0.081	0.178	0.479	1	
<i>Größe (in m2)</i>	-0.070	0.081	0.178	0.479	1	1