



LUDWIG-
MAXIMILIANS-
UNIVERSITÄT
MÜNCHEN

Institut für Soziologie
Dipl. Soz. Maximilian Sonnauer

Methoden II

Zusammenhangsmaße für kategoriale
und metrische Variablen



Prüfungsanmeldung Methoden II

- Die Anmeldung zur Prüfung läuft über die Tafelübung.
- Aus Einfachheitsgründen ist Maximilian Sonnauer als einziger Prüfer eingetragen.
- Unabhängig von Ihrem Übungstermin melden ALLE Teilnehmer*innen von Methoden II ihre Prüfung bei Maximilian Sonnauer an.

1. Wiederholung χ^2 Unabhängigkeitstest
2. Zusammenhang bei kategorialen Variablen
 - Cramers V
3. Metrische Zusammenhänge grafisch darstellen
 - Streudiagramme
 - Nichtlineare Verfahren
4. Zusammenhangsmaße metrischer Variablen
 - Pearsons Korrelationskoeffizient r

1. Wiederholung χ^2 Unabhängigkeitstest
2. Zusammenhang bei kategorialen Variablen
 - Cramers V
3. Metrische Zusammenhänge grafisch darstellen
 - Streudiagramme
 - Nichtlineare Verfahren
4. Zusammenhangsmaße metrischer Variablen
 - Pearsons Korrelationskoeffizient r

Welche Eigenschaften hat der χ^2 -Unabhängigkeitstest?

- Anwendung bei kategorialen Variablen
- Vergleicht beobachtete Häufigkeiten mit erwarteten Häufigkeiten unter Unabhängigkeit

$$\chi^2 = \sum_k \sum_m \frac{(O_{km} - E_{km})^2}{E_{km}}$$

O_{km} : beobachtete Häufigkeiten
 E_{km} : erwartete Häufigkeiten



		Geschlecht			
		weiblich	männlich		
Politische Ausrichtung	Rechts von der Mitte	O ₁ .	125	125	250
		E ₁ .	150	100	
	Mitte	O ₂ .	250	200	450
		E ₂ .	270	180	
	Links von der Mitte	O ₃ .	225	75	300
		E ₃ .	180	120	
			600	400	1000

$$\chi^2 = \frac{(125 - 150)^2}{150} + \frac{(125 - 100)^2}{100} + \frac{(250 - 270)^2}{270} + \frac{(200 - 180)^2}{180} + \frac{(225 - 180)^2}{180} + \frac{(75 - 120)^2}{120} = 42,25$$



Signifikanztest

Vergleich der berechneten Teststatistik (χ^2) mit dem kritischen Wert aus der Vergleichstabelle. Unabhängigkeit wird abgelehnt, wenn gilt:

$$\chi^2 = 42,25 > \chi_{1-\alpha}^2(\underbrace{df}_{\text{Freiheitsgrade}})$$

(Zeilen-1)(Spalten-1)

α entspricht dem festgelegten Signifikanzniveau (Fehler 1. Art). In der Tabelle muss für $1-\alpha$ nachgeschlagen werden.

df \ P	0.005	0.01	0.025	0.05	0.1	0.9	0.95
1	0.00004	0.0002	0.001	0.004	0.016	2.706	3.842
2	0.01	0.02	0.05	0.10	0.21	4.61	5.99
3	0.07	0.11	0.22	0.35	0.58	6.25	7.81
4	0.21	0.30	0.48	0.71	1.06	7.78	9.49
5	0.41	0.55	0.83	1.15	1.61	9.24	11.07
6	0.68	0.87	1.24	1.64	2.20	10.64	12.59
7	0.99	1.24	1.69	2.17	2.83	12.02	14.07
8	1.34	1.65	2.18	2.73	3.49	13.36	15.51
9	1.73	2.09	2.70	3.33	4.17	14.68	16.92
10	2.16	2.56	3.25	3.94	4.87	15.99	18.31

df = (3-1)(2-1) = 2

1- α = 0.95



Nachteile von χ^2

Keine Interpretation der *Stärke* eines Zusammenhangs, denn...

- die Höhe des Koeffizienten hängt von der Fallzahl ab.
- die Höhe des Koeffizienten hängt von der Anzahl der Zellen einer Tabelle ab, welche die Freiheitsgrade bestimmen $(k - 1) * (m - 1)$.

1. Wiederholung χ^2 Unabhängigkeitstest
2. Zusammenhang bei kategorialen Variablen
 - Cramers V
3. Metrische Zusammenhänge grafisch darstellen
 - Streudiagramme
 - Nichtlineare Verfahren
4. Zusammenhangsmaße metrischer Variablen
 - Pearsons Korrelationskoeffizient r

Cramers V

- Ausgangspunkt ist der χ^2 Test
- Normierung des Wertebereichs zwischen 0 und 1.
- Aussagen über die *Stärke* eines Zusammenhangs möglich
- *Keine* Aussage über die Richtung des Zusammenhangs möglich (Vorzeichen ignorieren)



Formel:	$V = \sqrt{\frac{\chi^2}{N * \min(k-1, m-1)}}$ <ul style="list-style-type: none"> ▪ k steht für „Zeile“ und m steht für „Spalte“ ▪ Man wählt entweder Anzahl der Spalten oder der Zeilen, je nachdem, was kleiner ist 	
Eigenschaften:	Wertebereich von [0,1]	
	Interpretation	<p>$V = 0 \rightarrow$ kein Zusammenhang zwischen X und Y</p> <p>$V = 1 \rightarrow$ perfekter Zusammenhang zwischen X und Y</p>
Achtung:	Keine Interpretation des Vorzeichens möglich	
	Kann für k x m-Tabellen berechnet werden (wobei k=m sein darf)	

Stata Syntax für Cramers V:

```
tabulate abhvar unabhvar, chi2 v
```

„chi2“ ruft den Chi-Quadrat Test auf [siehe letzte Woche]

Durch „V“ wird Cramers V angezeigt

Beispiel von letzter Woche:

H1: Amerikanische und ausländische Autos unterscheiden sich in ihrer Reparaturanfälligkeit.

Verwendete Variablen:

Dummy zu Herkunft: `foreign`

Anzahl der bisherigen Reparaturen: `rep78`



```
tab rep78 foreign, nofreq column
```

Repair Record 1978	Car type		Total
	Domestic	Foreign	
1	4.17	> 0.00	2.90
2	16.67	> 0.00	11.59
3	56.25	> 14.29	43.48
4	18.75	< 42.86	26.09
5	4.17	< 42.86	15.94
Total	100.00	100.00	100.00

Unterdrückt die
Anzeige der
absoluten
Häufigkeiten



```
tab rep78 foreign, nofreq column chi2 V
```

Repair Record 1978	Car type		Total
	Domestic	Foreign	
1	4.17	0.00	2.90
2	16.67	0.00	11.59
3	56.25	14.29	43.48
4	18.75	42.86	26.09
5	4.17	42.86	15.94
Total	100.00	100.00	100.00

Pearson chi2(4) = 27.2640 Pr = 0.000
Cramér's V = 0.6286

Der Chi-Quadrat Wert ist mit einem p-Wert = 0.000 höchst signifikant.

Es muss davon ausgegangen werden, dass ein Unterschied besteht.

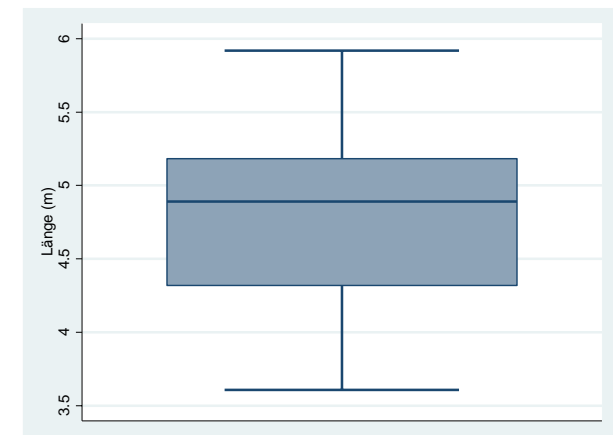
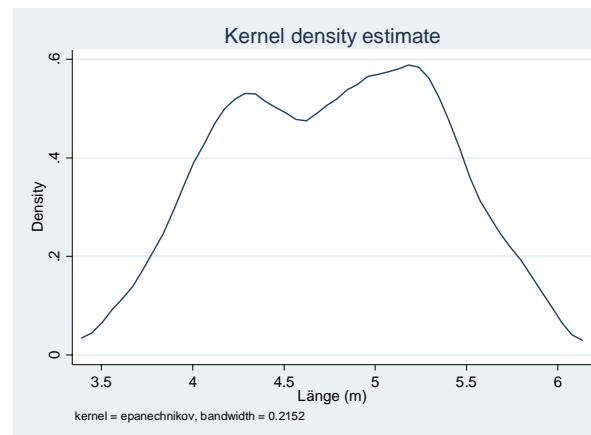
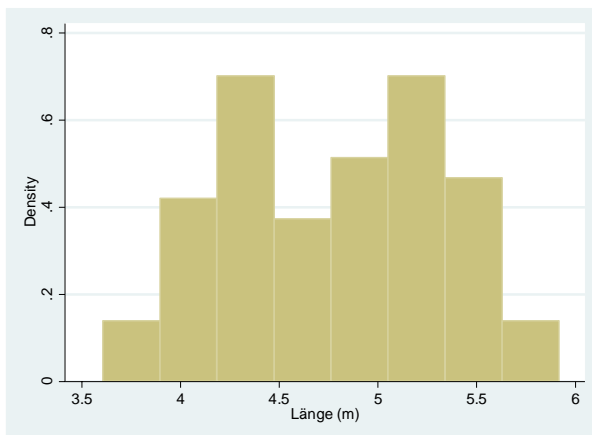
Cramers V ist mit 0.6286 mittel bis stark.

1. Wiederholung χ^2 Unabhängigkeitstest
2. Zusammenhang bei kategorialen Variablen
 - Cramers V
3. Metrische Zusammenhänge grafisch darstellen
 - Streudiagramme
 - Nichtlineare Verfahren
4. Zusammenhangsmaße metrischer Variablen
 - Pearsons Korrelationskoeffizient r



- Bisher wurden *univariate* grafische Verfahren betrachtet

Beispiel: Variable Länge



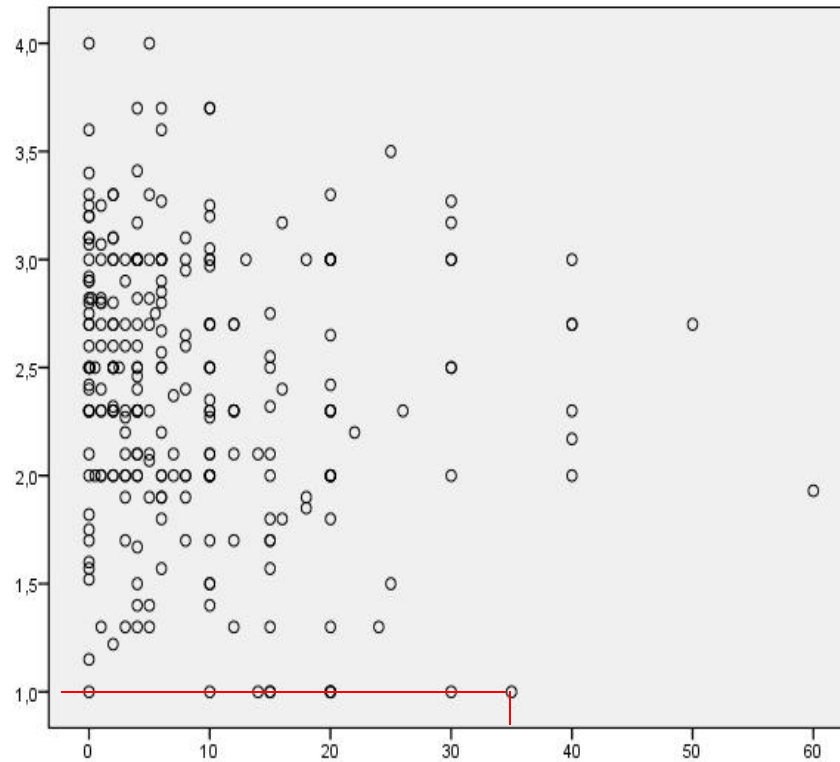
- Diese Grafen zeigen immer nur Verteilung *einer* Variable an



- Eine Kernkompetenz empirischer Sozialforscher*innen ist die Fähigkeit relevante Ergebnisse grafisch darzustellen.
- „Gut gestaltete Graphiken sind nämlich mitunter die einfachste und zugleich wirkungsvollste Möglichkeit zur Analyse und Kommunikation statistischer Information.“ (Bauer 2010)



Abhängige Variable auf der Y-Achse



Unabhängige Variable auf der X-Achse

- Wertepaare der zwei Variablen werden als Punktwolke (Scatterplot) im Koordinatensystem abgebildet.
- Mögliches Muster der Punkte kann Aufschluss über Beziehung zwischen den Variablen geben
- Achtung: Grafische Verfahren sollten immer mit statistischen Tests kombiniert werden

Stata Syntax für Streudiagramme:

`graph twoway scatter abhvar unabhvar`

zweidimensionale Grafik

scatterplot

Oder (Kurzform)

`scatter abhvar unabhvar`

Beispiel mit dem Auto-Datensatz:

H2: Je länger ein Auto ist, desto größer ist der Wendekreis.

Variablen:

UV: Länge der Autos in Metern (laenge, bereits erstellt)

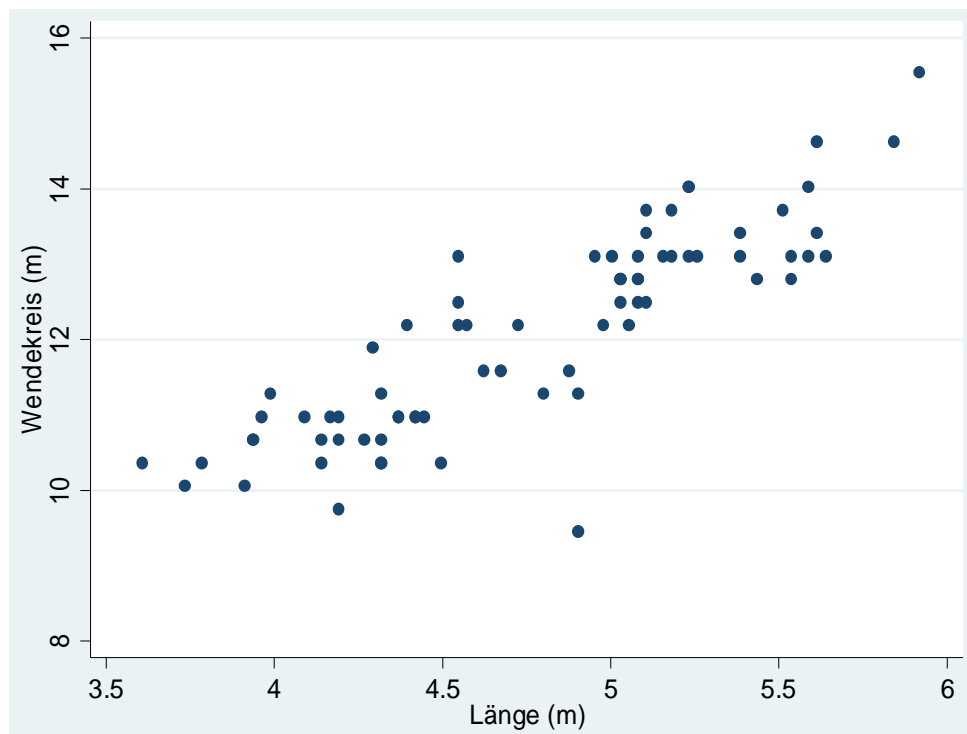
AV: Wendekreis in Metern (wend)

```
gen wend = turn*0.3048
```

```
lab var wend "Wendekreis (m) "
```



```
graph twoway scatter wend laenge
```



- Tendenziell linearer Zusammenhang erkennbar
- Längere Autos haben einen größeren Wendekreis

- Auf der Basis von Streudiagrammen kann ein erster Eindruck über Zusammenhänge gewonnen werden
- Nichtlineare (oder nonparametrische) Verfahren unterteilen Punktwolke in einzelne Abschnitte auf der X-Achse
- In den jeweiligen Abschnitten werden lokale Werte errechnet
 - Gibt Einblick in Verlauf eines Zusammenhangs, ohne weitere Überlegungen hinsichtlich Verteilung etc. zu benötigen

Median Bands

- Unterteilung der Daten in k Abschnitte
- Für jeden Abschnitt wird jeweils ein lokaler Median gebildet

Stata Syntax:

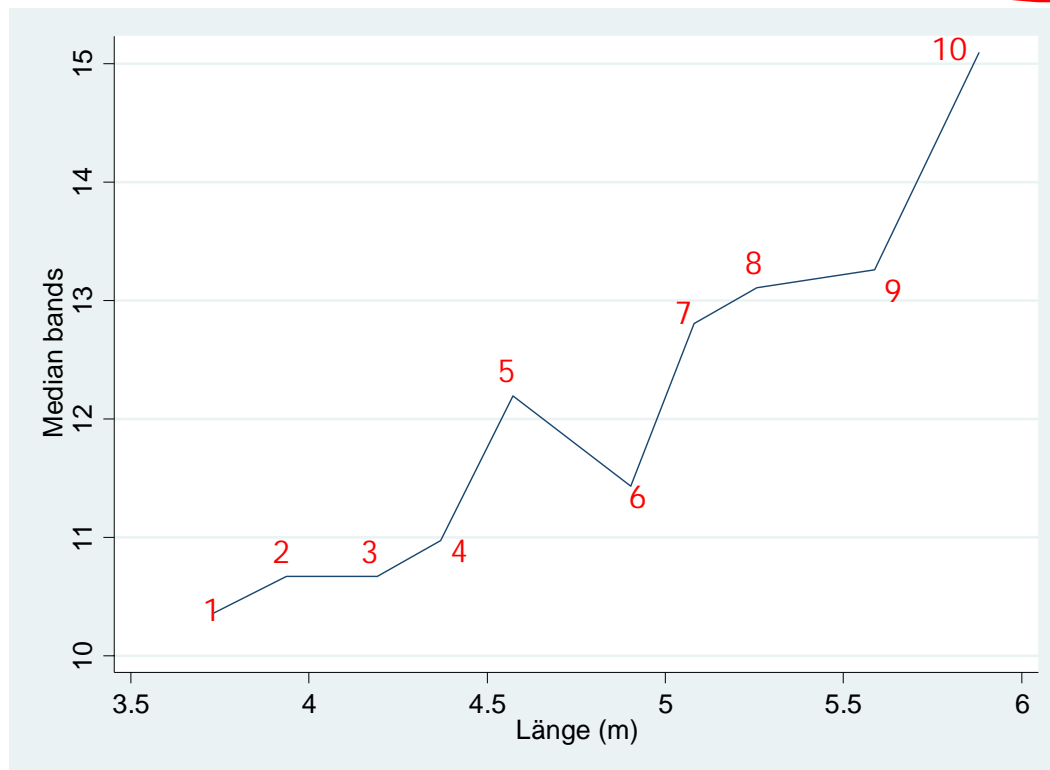
```
graph twoway mband AV UV , bands(k)
```

Anzahl der Abschnitte



Zurück zum Beispiel

graph twoway mband wend laenge, bands(10)



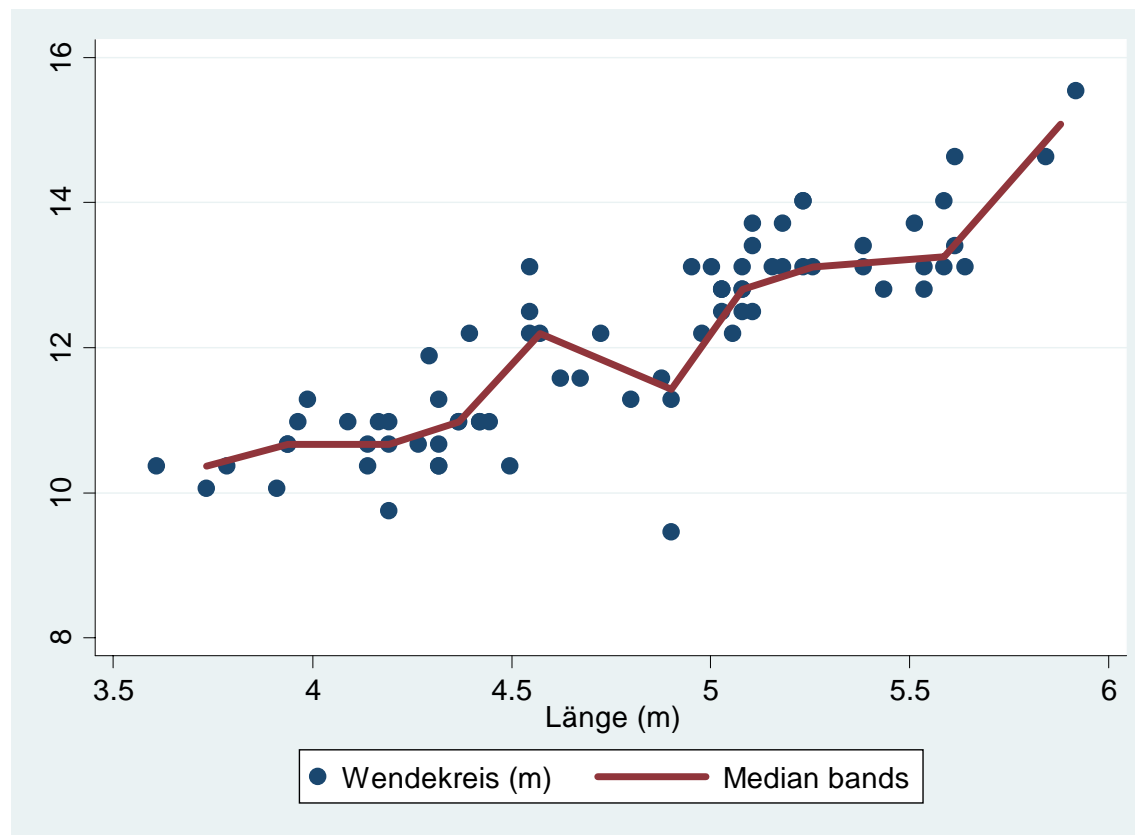
- Linear-positiver Zusammenhang scheint sich zu bestätigen
- In Abschnitt 6 jedoch „Delle“ nach unten



- In Stata können mehrere twoway-Grafiken miteinander kombiniert werden
- Zur Kombination verwendet man
`||` (also zweimal den Operator für „oder“)



```
graph twoway scatter wend laenge || mband wend laenge ,bands(10)
```



1. Wiederholung Chi² Unabhängigkeitstest
2. Zusammenhang bei kategorialen Variablen
 - Cramers V
3. Metrische Zusammenhänge grafisch darstellen
 - Streudiagramme
 - Nichtlineare Verfahren
4. Zusammenhangsmaße metrischer Variablen
 - Pearsons Korrelationskoeffizient r



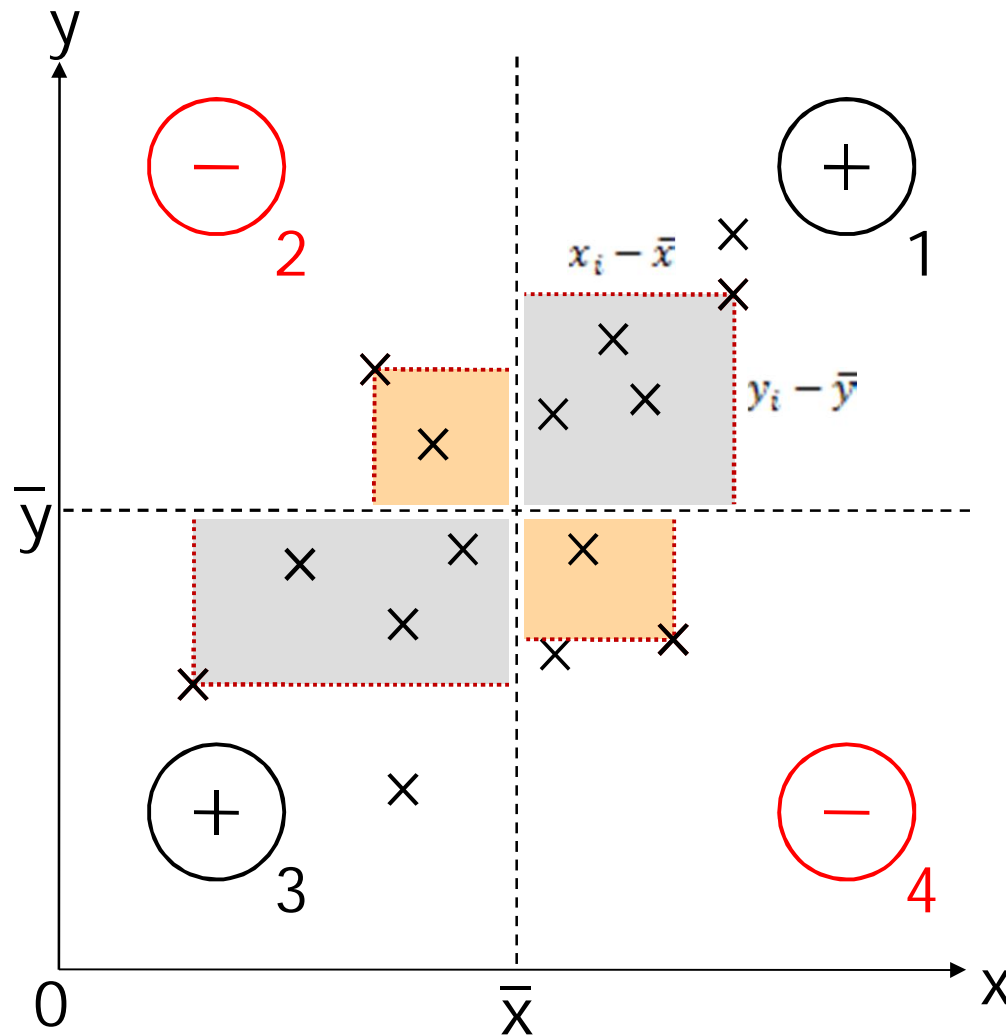
Fragestellung:

Besteht ein linearer Zusammenhang zwischen zwei metrischen Variablen?

Wie stark ist dieser Zusammenhang?

Anforderungen von Kovarianz und Korrelationskoeffizient nach Pearson:

- Beide Variablen müssen metrisch skaliert (Intervall-, Verhältnis- oder Absolutskala) sein
- Ein linearer Zusammenhang muss unterstellt werden können



Berechnung der Kovarianz:

$$\text{cov}(x, y) = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x}) * (y_i - \bar{y})$$

$$(x_i - \bar{x}) * (y_i - \bar{y})$$

- ist *positiv* im 1. und 3. Quadranten
- ist *negativ* im 2. und 4. Quadranten



Formel:	$r = \frac{\text{cov}(x, y)}{\sqrt{\text{var}(y)} \cdot \sqrt{\text{var}(x)}}$	
Eigenschaften:	Wertebereich von $[-1, +1]$	
	Interpretation	<p> $r = -1$ Perfekter, negativ-linearer Zusammenhang \uparrow $\left. \vphantom{\uparrow} \right\}$ <i>Je größer X, desto kleiner Y</i> </p> <p> $r = 0$ Kein linearer Zusammenhang. </p> <p> \downarrow $\left. \vphantom{\downarrow} \right\}$ <i>Je größer X, desto größer Y</i> </p> <p> $r = 1$ Perfekter, positiv-linearer Zusammenhang. </p>

Vorteile

- Symmetrisches Maß, d.h. abhängige und unabhängig Variablen können vertauscht werden.
- Das Korrelationsmaß informiert darüber, wie nah die Datenpunkte um eine lineare Gerade streuen, welche an die Beobachtungen angepasst ist.
- Die Richtung des Zusammenhangs lässt sich interpretieren.

Nachteile

- Das Korrelationsmaß ist ausreißeranfällig, d.h. sehr große Werte haben einen starken Einfluss auf die Maßzahl
- Das Maß informiert nicht darüber, wie stark bei Veränderung der einen Variable die Veränderung der anderen Variable ist (gibt keinen Aufschluss über Kausalität).

Stata Syntax für Korrelationen:

```
correlate abhvar unabvar
```

Zurück zum Beispiel:

```
correlate wend laenge
```

	wend	laenge
wend	1.0000	
laenge	0.8643	1.0000

- Es gibt eine stark positive Korrelation zwischen der Länge eines Autos und dem Wendekreis

Aber: Ist dieser Zusammenhang auch überzufällig (signifikant)?

Für detailliertere Korrelationstabellen verwenden wir den Befehle
pairwise-correlation:

```
pwcorr abhvar unabvar , sig
```

```
pwcorr wend laenge , sig
```

	wend	laenge
wend	1.0000	
laenge	0.8643	1.0000
	0.0000	

p < 0,0000:

→ Zusammenhang ist höchst
signifikant

→ Längere Autos haben einen
signifikant größeren
Wendekreis

„sig“ zeigt das Signifikanzniveau jeder Variable an



	Griechischer Buchstabe	Skalenniveau	Tabellengröße	Wertebereich	Symmetrisch	PRE-Maß	Stärke/Richtung
Chi ²	χ^2	nominal					-
Phi	Φ	nominal	2*2 k*m	[-1, 1] [0, 1]	Ja	Nein	Stärke
Cramers V	V	nominal	k*m	[0, 1]	Ja	Nein	Stärke
Lambda	λ	nominal	k*m	[0, 1]	Nein	Ja	Stärke
Gamma	γ	ordinal	k*m	[-1, 1]	Ja	Ja	Stärke Richtung
Kendall's Tau b	τ_b	ordinal	k*m	[-1, 1]	Ja	Nein	Stärke Richtung
Kendall's Tau c	τ_c	ordinal	k*m	[-1, 1]	Ja	Nein	Stärke Richtung
eta und eta ²	η η^2	UV: nominal/ordinal AV: metrisch		[0, 1]	Nein	eta ² : ja eta: nein	Stärke Richtung
Korrelation	r	metrisch		r: [-1;1]	Ja	r: nein	Stärke Richtung



Cramers V: `tabulate AV UV , chi2 v`

Scatterplot: `graph twoway scatter AV UV if X==`

Median Bands: `graph twoway mband AV UV, bands(k)`

komb. Grafen: `graph twoway plot1 || plot2`

Korrelation: `correlate AV UV`

Korr. detailliert: `pwcorr AV UV , sig`

Sie vermuten, dass der Zusammenhang zwischen der Länge von Autos und dem Wendekreis sich je nach der Herkunft der Autos unterscheidet.

Überprüfen Sie diese Vermutung.

- 1) Erstellen Sie vier Variablen um folgende Unterscheidung treffen zu können:

	Herkunft	
	Domestic	Foreign
Wendekreis	<i>wend_dom</i>	<i>wend_for</i>
Länge	<i>laenge_dom</i>	<i>laenge_for</i>

- 2) Unterziehen Sie beide Herkunftsgruppen einer grafischen Analyse
- 3) Stützen Sie die grafische Analyse mittels einer geeigneten statistischen Maßzahl
- 4) Interpretieren Sie die Ergebnisse inhaltlich