

TUTORIUM

Datenauswertung

ASSOZIATIONSMASSE

AGENDA

- Assoziationsmaße
- Pearson's r per Hand
- Probleme mit Pearson's r
- Pearson's r mit R
- Scatterplots in R

ASSOZIATIONSMAßE

Assoziationsmaße

- Messen die Stärke des Zusammenhangs einer Zusammenhangshypothese
- Erstrecken sich (bis auf V) von -1 bis 1
- Gibt es für verschiedene Skalenniveaus:
 - Nominal: Cramer's V
 - Ordinal: Spearman's ρ
 - Intervall: Kendall's τ
 - Rational: Pearson's r

KOVARIANZ UND VARIANZ

$$s^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n - 1}$$

| Formel Varianz

$$\leftrightarrow s^2 = \frac{\sum_{i=1}^n (x_i - \bar{x}) * (x_i - \bar{x})}{n - 1}$$

| Formel Varianz (umgeformt)

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

| Formel Kovarianz

KOVARIANZ – BEISPIEL


Fall i	1	2	3	4	5
Radiokonsum x_i	60	120	30	120	180
Zeitungskonsum y_i	15	0	10	20	15

$$\bar{x} = \frac{60 + 120 + 30 + 120 + 180}{5} = 102$$

$$\bar{y} = \frac{15 + 0 + 10 + 20 + 15}{5} = 12$$

Fall i	$x_i - \bar{x}$	$y_i - \bar{y}$	$(x_i - \bar{x}) * (y_i - \bar{y})$
1	60 - 102 = -42	15 - 12 = 3	-42 * 3 = -126
2	120 - 102 = 18	0 - 12 = -12	18 * (-12) = -216
3	30 - 102 = -72	10 - 12 = -2	-72 * (-2) = 144
4	120 - 102 = 18	20 - 12 = 8	18 * 8 = 144
5	180 - 102 = 78	15 - 12 = 3	78 * 3 = 234
			-126 + (-216) + 144 + 144 + 234 = 180

$$\text{cov}(x, y) = \frac{\sum_{i=1}^n (x_i - \bar{x}) * (y_i - \bar{y})}{n - 1} = \frac{180}{5 - 1} = 45$$

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


PEARSON'S R

Pearson's r

- Ist die durch das Produkt der Standardabweichungen standardisierte Kovarianz
- Ist ein sehr gängiges Zusammenhangsmaß in den Kommunikationswissenschaften
- Gibt an, wie stark zwei metrische Variablen zusammenhängen:
 - -1 → annähernd perfekter negativer linearer Zusammenhang
 - 0 → kein Zusammenhang
 - 1 → annähernd perfekter positiver linearer Zusammenhang

PEARSON'S R – FORMELN UND INTERPRETATION

$$r = \frac{\text{cov}(x, y)}{s_x * s_y} = \frac{\frac{\sum_{i=1}^n (x_i - \bar{x}) * (y_i - \bar{y})}{n - 1}}{\sqrt{\frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n - 1}} * \sqrt{\frac{\sum_{i=1}^n (y_i - \bar{y})^2}{n - 1}}}$$

r-Wert (negativ)	R-Wert (positiv)	Stärke des Einflusses
$-0,1 \leq r \leq 0$	$0 \leq r \leq 0,1$	nicht vorhanden
$-0,2 \leq r \leq -0,1$	$0,2 \leq r \leq 0,1$	schwach
$-0,4 \leq r \leq -0,2$	$0,2 \leq r \leq 0,4$	mittel
$-0,8 \leq r \leq -0,4$	$0,4 \leq r \leq 0,8$	stark
$-1,0 \leq r \leq -0,8$	$0,8 \leq r \leq 1,0$	(annähernd) linear

PEARSON'S R – BEISPIEL

$$\text{cov}(x, y) = 45, \quad \bar{x} = 102, \quad \bar{y} = 12$$

Fall i	$x_i - \bar{x}$	$(x_i - \bar{x})^2$
1	$60 - 102 = -42$	1764
2	$120 - 102 = 18$	324
3	$30 - 102 = -72$	5184
4	$120 - 102 = 18$	324
5	$180 - 102 = 78$	6084
		$\Sigma = 13680$

Fall i	$y_i - \bar{y}$	$(y_i - \bar{y})^2$
1	$15 - 12 = 3$	9
2	$0 - 12 = -12$	144
3	$10 - 12 = -2$	4
4	$20 - 12 = 8$	64
5	$15 - 12 = 3$	9
		$\Sigma = 230$

$$s = \sqrt{\frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n - 1}}$$

$$s_x = \sqrt{\frac{13680}{5 - 1}} \approx 58,48$$

$$s_y = \sqrt{\frac{230}{5 - 1}} \approx 7,58$$

$$r = \frac{\text{cov}(x, y)}{s_x * s_y} = \frac{45}{58,48 * 7,58} \approx 0,10$$

$r \leq 0,1 \rightarrow$ Kein Zusammenhang!

PEARSON'S R – AUFGABE

Aufgabe: Prüft, ob die Variablen Radiokonsum und Internetkonsum zusammenhängen, und gebt, falls das so ist, an, wie stark sie zusammenhängen.

$$\text{cov}(x, y) = \frac{\sum_{i=1}^n (x_i - \bar{x}) * (y_i - \bar{y})}{n - 1}, \quad \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i, \quad r = \frac{\text{cov}(x, y)}{s_x * s_y}, \quad s = \sqrt{\frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n - 1}}$$

$$\bar{y} = 192$$

$$s_y \approx 115,41$$

Fall i	1	2	3	4	5
Radiokonsum x_i	120	180	60	45	60
Internetkonsum y_i	180	240	120	60	360

PEARSON'S R – LÖSUNG I: MITTELWERT UND STANDARDABWEICHUNG (x)

Fall i	1	2	3	4	5
Radiokonsum x_i	120	180	60	45	60

$$\bar{x} = \frac{120 + 180 + 60 + 45 + 60}{5} = 93$$

$$s_x = \sqrt{\frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n - 1}} = \sqrt{\frac{12780}{5 - 1}} \approx 56,52$$

Fall i	$x_i - \bar{x}$	$(x_i - \bar{x})^2$
1	120 - 93 = 27	27 ² = 729
2	180 - 93 = 87	87 ² = 7569
3	60 - 93 = -33	-33 ² = 1089
4	45 - 93 = -48	-48 ² = 2304
5	60 - 93 = -33	-33 ² = 1089
		$\Sigma = 12780$

PEARSON'S R – LÖSUNG II: MITTELWERT UND STANDARDABWEICHUNG (Y)

Fall i	1	2	3	4	5
Internetkonsum y_i	180	240	120	60	360

$$\bar{y} = \frac{180 + 240 + 120 + 60 + 360}{5} = 192$$

$$s_y = \sqrt{\frac{\sum_{i=1}^n (y_i - \bar{y})^2}{n - 1}} = \sqrt{\frac{53280}{5 - 1}} \approx 115,41$$

Fall i	$y_i - \bar{y}$	$(y_i - \bar{y})^2$
1	$180 - 192 = -12$	$-12^2 = 144$
2	$240 - 192 = 48$	$48^2 = 2304$
3	$120 - 192 = -72$	$-72^2 = 5184$
4	$60 - 192 = -132$	$-132^2 = 17424$
5	$360 - 192 = 168$	$168^2 = 28224$
		$\sum = 53280$

PEARSON'S R – LÖSUNG III: KOVARIANZ (X,Y)

Fall i	1	2	3	4	5
Radiokonsum x_i	120	180	60	45	60
Internetkonsum y_i	180	240	120	60	360

Fall i	$x_i - \bar{x}$	$y_i - \bar{y}$	$(x_i - \bar{x}) * (y_i - \bar{y})$
1	$120 - 93 = 27$	$180 - 192 = -12$	$27 * (-12) = -324$
2	$180 - 93 = 87$	$240 - 192 = 48$	$87 * 48 = 4176$
3	$60 - 93 = -33$	$120 - 192 = -72$	$-33 * (-72) = 2376$
4	$45 - 93 = -48$	$60 - 192 = -132$	$-48 * (-132) = 6336$
5	$60 - 93 = -33$	$360 - 192 = 168$	$-33 * 168 = -5544$
			$\Sigma = 7020$

$$\text{cov}(x, y) = \frac{\sum_{i=1}^n (x_i - \bar{x}) * (y_i - \bar{y})}{n - 1} = \frac{7020}{5 - 1} = 1755$$

PEARSON'S R – LÖSUNG IV: STANDARDISIERUNG, INTERPRETATION

$$\text{cov}(x, y) = 1755, \quad s_x \approx 56,52, \quad s_y \approx 115,41$$

$$r = \frac{\text{cov}(x, y)}{s_x * s_y} = \frac{1755}{56,52 * 115,41} \approx 0,27$$

→ Es gibt einen mittelstarken positiven Zusammenhang zwischen den Variablen Radiokonsum und Internetkonsum.

PEARSON'S R – PROBLEME

Pearson's r

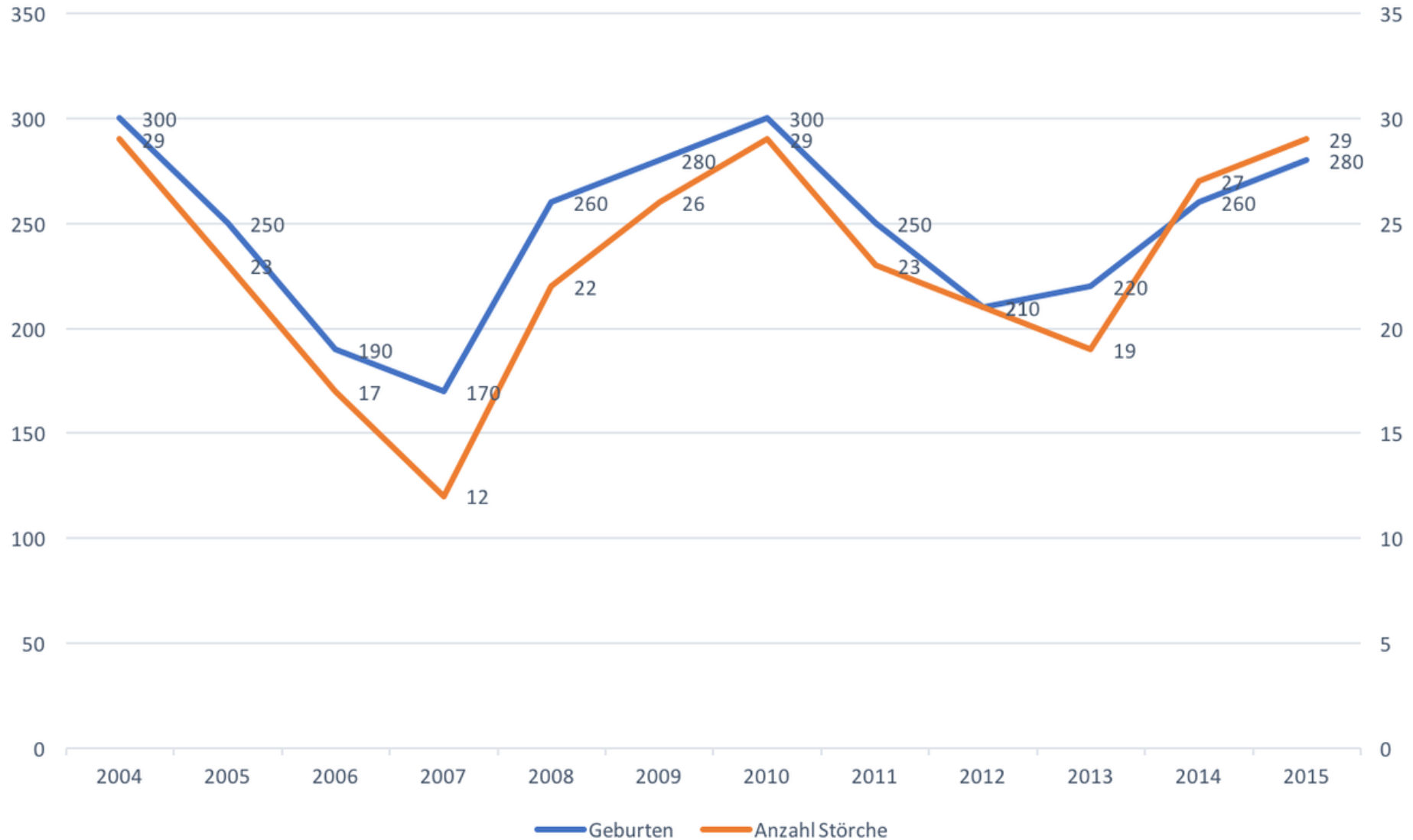
- Funktioniert nur für metrische und stetige Variablen
- Benötigt große Stichprobenzahlen ($n > 100$) oder kleinere Stichprobenzahlen, wobei beide Variablen normalverteilt sind
- Ist extrem ausreißerempfindlich
- Kann deshalb ersetzt werden durch
 - Kendall's τ
 - Spearman's ρ
- Sagt nicht, ob die Variablen *kausal* zusammenhängen

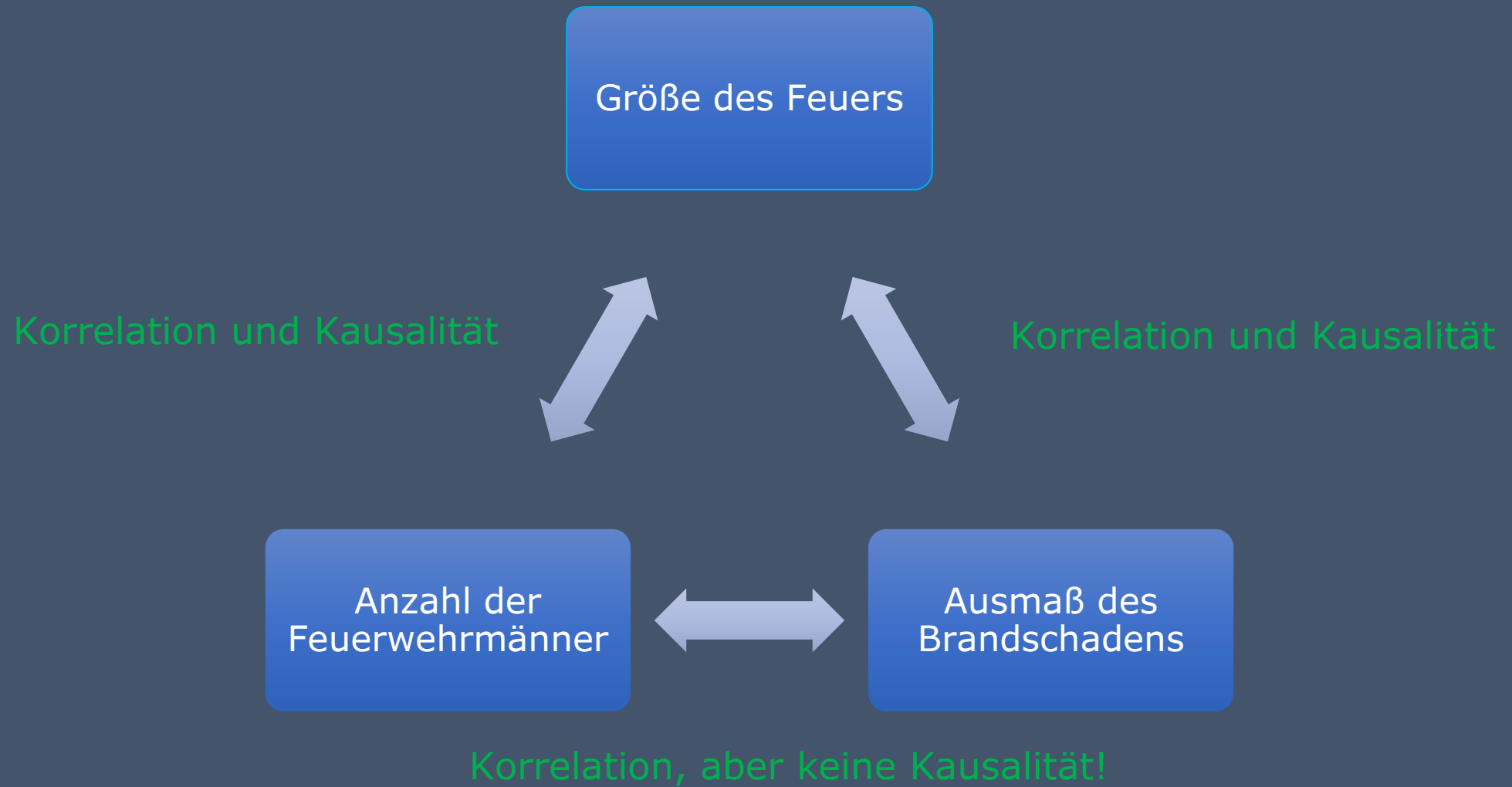
KORRELATIONEN UND KAUSALITÄT

Die Korrelation zwischen zwei Variablen

- ist notwendige Bedingung für einen Kausalzusammenhang dieser Variablen
- Ist dafür aber nicht hinreichend, denn
 - Die Korrelation könnte zufällig zustande gekommen sein
 - Die Variablen könnten aufgrund einer dritten Variable zusammenhängen
 - Die Korrelation sagt nichts über die Richtung aus (beeinflusst x y oder andersherum?)
- Korrelationen dürfen nie als Kausalzusammenhänge interpretiert werden!

Zahl der Geburten / Zahl der Störche





KORRELATIONEN IN R: R-MATRIX

```
#Laden und Auswahlen der Daten
```

```
library(knitr)
```

```
library(Hmisc)
```

```
load("daten_x.RData")
```

```
datensatz <- subset(daten_x, Bedingung)
```

```
#Berechnung
```

```
matrix= rcorr(cbind(datensatz$var1, datensatz$var2, datensatz$var3))
```

```
#Korrelationsmatrix zusammenstellen
```

```
rmatrix = data.frame(rmatrix$r, row.names = c("var1label", "var2label", "var3label"))
```

```
colnames(rmatrix) = c("var1label", "var2label", "var3label")
```

```
kable(round(rmatrix, 2), caption = "Name der Graphik")
```

KORRELATIONEN IN R: R MATRIX - BEISPIEL

```
#Laden und Auswählen der Daten
```

```
library(knitr)
```

```
library(Hmisc)
```

```
load("daten_2019.RData")
```

```
datensatz <- subset(daten_2019, jahr>2000)
```

```
#Berechnung
```

```
matrix = rcorr(cbind(datensatz$tv_minuten, datensatz$tz_minuten, datensatz$www_minuten))
```

```
#Korrelationsmatrix zusammenstellen
```

```
rmatrix = data.frame(matrix$r, row.names = c("TV-Konsum", "Zeitungskonsum", "Inernetkonsum"))
```

```
colnames(rmatrix) = c("TV-Konsum", "Zeitungskonsum", "Inernetkonsum")
```

```
kable(round(rmatrix, 2), caption = "Korrelationsmatrix")
```

Korrelationsmatrix

	TV-Konsum	Zeitungskonsum	Inernetkonsum
TV-Konsum	1.00	0.33	-0.20
Zeitungskonsum	0.33	1.00	-0.15
Inernetkonsum	-0.20	-0.15	1.00

KORRELATIONEN IN R: P-MATRIX

```
#Laden von Daten und Paketen
```

```
library(knitr)
```

```
library(Hmisc)
```

```
load("daten_x.RData")
```

```
datensatz <- subset(daten_x, Bedingung)
```

```
#p-Matrix erstellen
```

```
pmatrix = data.frame(matrix$P, row.names = c("var1label", "var2label", "var3label"))
```

```
colnames(pmatrix) = c("var1label", "var2label", "var1label")
```

```
kable(round(pmatrix, 2), caption = "Überschrift")
```

KORRELATIONEN IN R: P-MATRIX

```
#p-Matrix erstellen
```

```
matrix = rcorr(cbind(datensatz$tv_minuten,  
datensatz$tz_minuten, datensatz$www_minuten))
```

```
pmatrix = data.frame(matrix$P, row.names = c("TV-  
Konsum", "Zeitungskonsum", "Internetkonsum"))
```

```
colnames(pmatrix) = c("TV-Konsum", "Zeitungskonsum",  
"Internetkonsum")
```

```
kable(round(pmatrix, 2), caption = "p-Matrix")
```

p-Matrix

	TV-Konsum	Zeitungskonsum	Haushaltsgröße
TV-Konsum	NA	0	0
Zeitungskonsum	0	NA	0
Haushaltsgröße	0	0	NA

KORRELATIONEN IN R: N-MATRIX

#n-Matrix erstellen

```
nmatrix = data.frame(matrix$n, row.names = c("var1label", "var2label", "var3label"))
```

```
colnames(nmatrix) = c("var1label", "var2label", "var3label")
```

#n-Matrix ausgeben

```
kable(nmatrix, caption = "Überschrift")
```


KORRELATIONEN IN R: N-MATRIX

#n-Matrix erstellen

```
nmatrix = data.frame(matrix$n, row.names = c("TV-Konsum", "Zeitungskonsum", "Internetkonsum"))
```

```
colnames(nmatrix) = c("TV-Konsum", "Zeitungskonsum", "Internetkonsum")
```

#n-Matrix ausgeben

```
kable(nmatrix, caption = "n-Matrix")
```

n-Matrix

	TV-Konsum	Zeitungskonsum	Haushaltsgröße
TV-Konsum	483	473	481
Zeitungskonsum	473	476	474
Haushaltsgröße	481	474	486

DER DETERMINATIONSKOEFFIZIENT R^2

Der Determinationskoeffizient

- Gibt an, wie groß der Anteil der Varianz ist, den die beiden Variablen teilen → wie gut passt das Modell?
- Sagt dennoch nichts über Kausalität aus
- Kann näherungsweise auch für Spearman's ρ und Kendall's τ angegeben werden
- Formel: $R^2 = r^2 = r * r$

SCATTERPLOTS

Scatterplots

- Sind Diagramme, die den Zusammenhang zwischen zwei numerischen Variablen darstellen
- Stellen jeden Fall als einen Punkt dar, der aus den Ausprägungen der beiden Variablen besteht: $(x_i | y_i)$
- Können durch eine Regressionslinie ergänzt werden → Veranschaulichung des Determinationskoeffizienten

SCATTERPLOTS IN R

```
#Laden von Daten und Paketen
```

```
load("daten_x.RData")
```

```
library(ggplot2)
```

```
datensatz <- subset(daten_x, Bedingung)
```

```
#Scatterplot erstellen
```

```
ggplot(daten_x) + geom_point() +  
geom_smooth(method = "lm") + aes(var1, var2) + labs  
(x="var1label", y="var2label")
```

SCATTERPLOTS IN R – BEISPIEL

```
#Laden von Daten und Paketen
```

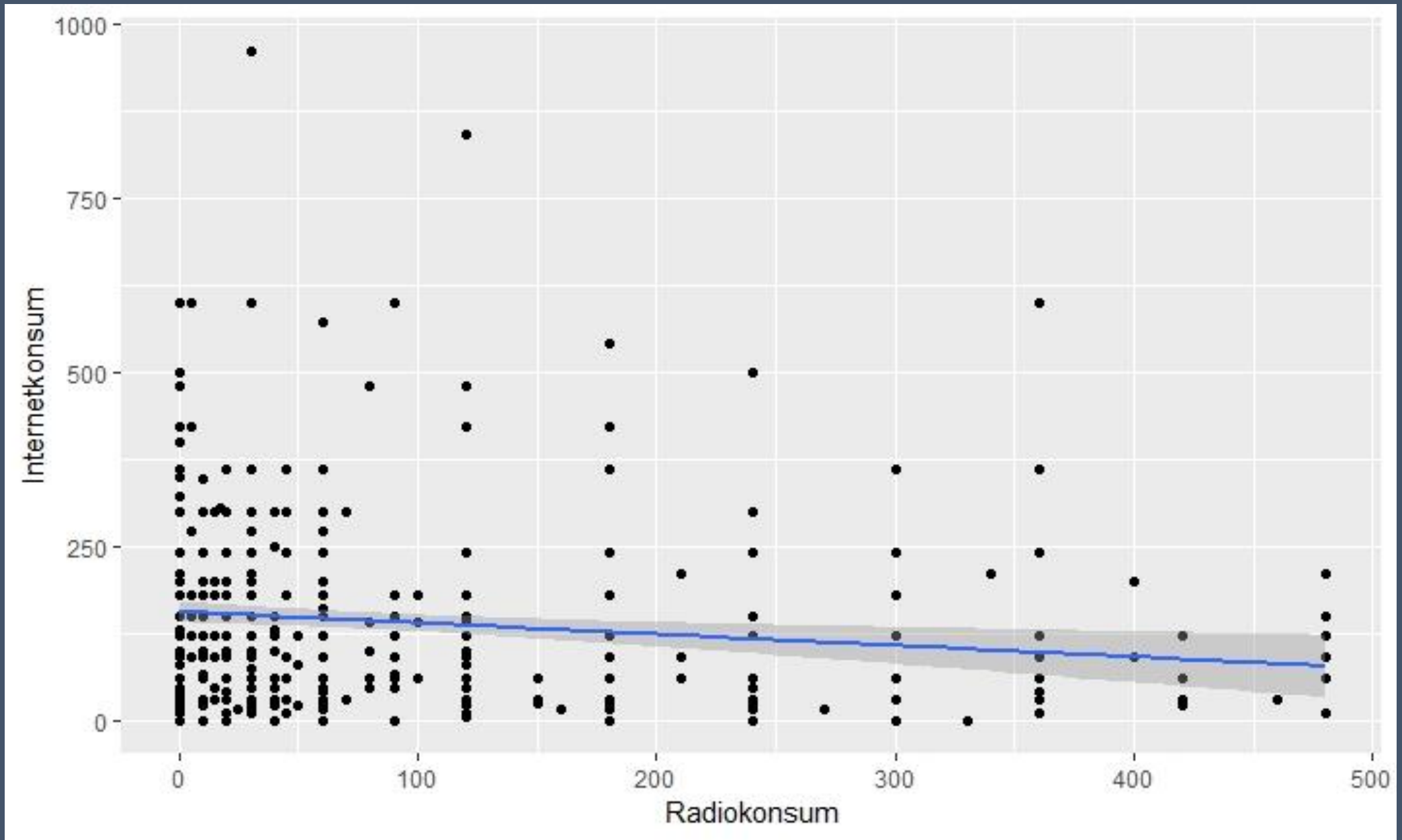
```
load("daten_2019.RData")
```

```
library(ggplot2)
```

```
datensatz <- subset(daten_2019, radio_minuten < 500 &  
www_minuten < 500)
```

```
#Scatterplot erstellen
```

```
ggplot(datensatz) + geom_point() + geom_smooth(method  
= "lm") + aes(radio_minuten, www_minuten) + labs  
(x="Radiokonsum", y="Internetkonsum")
```



SPEARMAN'S ρ

Spearman's ρ

- Ist ein Assoziationsmaß für ordinalskalierte Variablen
- Funktioniert genauso wie Pearson's r , nur arbeitet mit den Rangplätzen der Variablen
- Liefert kleinere Werte als Pearson's r
- benutzt man
 - Bei metrischen Variablen mit vielen Ausreißern
 - Bei natürlicherweise ordinalskalierten Variablen (bspw. Zustimmung: überhaupt nicht – voll und ganz)
- Formel: $\rho = \frac{\text{cov}(rp(x), rp(y))}{s(rp(x)) * s(rp(y))}$

SPEARMAN'S ρ – RANGPLÄTZE BESTIMMEN

Fall i	1	2	3	4	5
Radiokonsum x_i	60	120	30	120	180
Zeitungskonsum y_i	15	0	10	20	15

Fall i	1	2	3	4	5
Radiokonsum x_i	4	3	5	2	1
Zeitungskonsum y_i	2	5	4	1	3

FRAGEN?

SELBSTSTUDIUM

- Forschungsbericht (EFB und ZFB)
 - r-, n- und p-Matrix für `www_minuten`, `tz_minuten` und `soz_haushalt`
 - Inhaltliche Sätze zu den Ergebnissen
- Anfängen: Wiederholung
- Fragen zur Wiederholungssitzung sammeln und bis Freitag, 29.06.2019, 23:59 Uhr per Mail zuschicken

BIS NÄCHSTE WOCHEN!