

Statistik und Datenanalyse: Aufbau

Übung 1 – 4. Sitzung

Benjamin Fretwurst

▶ PDF-Version der Folien



Inhalt

- Übung 1 a+b
- Take Home – Ausblick – Vokabeln



Orga

Orga

- Es gibt neue Studienteilnahmestudien!



Lernziele

Grundprinzipien der Regression

- Sinn und Zweck von Quarto
- Einübung Regressionsanalyse
- Übung zur Berechnung von b_2 , um zu verstehen, wie die Zusammenhänge bivariat angelegt sind und die übrigen Dritteinflüsse “herausgerechnet” werden
- Vor- und Nachteile der visuellen Residualanalyse vs. Tests



Übung 1 a+b



Ü1.1 Erstellen Sie eine Quarto-Datei.qmd

1. Öffnen Sie R-Studio
2. In R-Studio ⇨ File ⇨ New File ⇨ Quarto Document...
3. Klicken Sie unten links auf «Create Empty Document»
4. (Wählen Sie als **title** «Erste Regression»)
5. Fügen Sie einen r-Chunk hinzu mit diesem Schalter: 
6. speichern Sie an einem günstigen Ort
(am besten in der Cloud + nicht auf Desktop)



Ü1.2 Installation und Setup

Die R-Befehle für die Installation von Paketen haben wir in die Datei “Installation.R” ausgelagert, weil man sie im Grunde jeweils nur einmal braucht. Ich habe Ihnen [hier eine Installationsdatei](#) gebaut, mit der Sie die Pakete mit höherer Erfolgchance installieren können. Mit folgendem Befehl wird diese Datei automatisch von unserer Homepage heruntergeladen und im Unterordner “files” des Projekts gespeichert. Sie können Sie dann dort öffnen und (am besten Zeilenweise) ausführen, wenn Sie die Pakete noch nicht installiert haben.

```
1 # Prüfe, ob es in dem Ordner in der die Uebung_1_ab.qmd gespeichert ist, ein  
2 if(dir.exists("files")){} else {dir.create("files")}
```

NULL

```
1 # Lade die Installations.R herunter und speichere sie im Unterordner des Pro  
2 download.file("https://stat.ikmz.uzh.ch/Aufbau/Folien/Sitzung_04/files/Insta
```



Für generelle Grundeinstellungen haben wir eine “_common.R” angelegt, in der wir den Aufruf der Basispakete des tidyverse geschrieben haben und andere Optionen und Einstellungen für Designs (wie Farben). Die Datei kann man dann immer am Anfang seiner Quarto-Dateien aufrufen und braucht diese Generaleinstellungen nicht immer wieder neu kopieren. Das ist doch praktisch.

```
1 # Prüfe, ob es in dem Ordner in der die Uebung_1_ab.qmd gespeichert ist, ein  
2 if(dir.exists("files")){} else {dir.create("files")}
```

NULL

```
1 download.file("https://stat.ikmz.uzh.ch/Aufbau/Folien/Sitzung_04/files/_common.R")  
2  
3 source("files/_common.R")
```



Ü1.3 Laden Sie die Daten

Laden Sie den Fragebogen [hier](#) runter und schauen ihn an.

Laden Sie die Daten und lassen Sie mal die Variablenlabel raus:

► R-Code anzeigen



Ü1.4 Rechnen Sie ein Regressionsmodell

```

1  DATEN <- DATEN |> haven::zap_formats()
2
3  Modell_1 <- lm(E201_10 ~ E102_02, data = DATEN)
4
5  summary(Modell_1)
6  ##
7  ## Call:
8  ## lm(formula = E201_10 ~ E102_02, data = DATEN)
9  ##
10 ## Residuals:
11 ##      Min       1Q   Median       3Q      Max
12 ## -1.2726 -0.7684 -0.2641  0.7274  2.7274
13 ##
14 ## Coefficients:
15 ##              Estimate Std. Error t value Pr(>|t|)
16 ## (Intercept)   1.28105    0.16912   7.575 2.60e-13
17 ## E102_02       0.49578    0.06668   7.435 5.73e-13

```

Was sehen Sie?

1. Wie gross ist R^2 ?
2. Wie gross ist die bivariate Korrelation r ? (selbst ausrechnen)
3. Ist der Zusammenhang positiv oder negativ?
4. Ist der Zusammenhang signifikant?

Ü1.5 Verändern Sie das Regressionsmodell

Kopieren Sie den r -Chunk der letzten Folie und setzen Sie andere Variablen ein: Nehmen Sie die Variablen für «Statistik Einführung hat mir viel Spass gemacht» und erklären Sie damit: «Ich freu mich auf Statistik Aufbau!».

Beantworten wieder die Fragen:

1. Wie gross ist R^2 ?
2. Wie gross ist die bivariate Korrelation r ? (selbst ausrechnen)
3. Ist der Zusammenhang positiv oder negativ?
4. Ist der Zusammenhang signifikant?

► Code



Ü1.6 b_2 aus Korrelationen und SDs berechnen

Note

Lassen Sie die Korrelationen durchlaufen, schauen Sie sich an, wo was steht und setzen Sie es in die Formel für $b_2 = \frac{r_{Y2} - r_{23}r_{Y3}}{(1 - R_{2.3}^2)} \frac{s_Y}{s_2}$, um es zu berechnen.

$$\text{Also: } b_2 = \frac{.50 - (-.32 \cdot -.23)}{(1 - .32^2)} \frac{1.05}{1.07}$$

► R-Code anzeigen



Man kann natürlich auch R nutzen

▶ R-Code anzeigen



Ü1.7 Berechnen Sie b_2 mit Hilfe einer Regressionsanalyse

```

1 Modell11 <- lm(E201_10 ~ E102_02 + E102_04, data = DATEN)
2
3 Modell_1_beta <- lm.beta::lm.beta(Modell11)
4
5 summary(Modell_1_beta, digits = digits, maxsum = maxsum)
6 ##
7 ## Call:
8 ## lm(formula = E201_10 ~ E102_02 + E102_04, data = DATEN)
9 ##
10 ## Residuals:
11 ##      Min       1Q   Median       3Q      Max
12 ## -1.3852 -0.7194 -0.2281  0.5642  2.7439
13 ##
14 ## Coefficients:
15 ##              Estimate Standardized Std. Error t value Pr(>|t|)
16 ## (Intercept)   1.57016             NA     0.32129   4.887 2.46e-06 ***
17 ## E102_02       0.47210             0.48031   0.07031   6.714 3.08e-10 ***

```



Ü1.8 Geben Sie mit folgendem Befehl die Tolerance und VIF-Werte raus

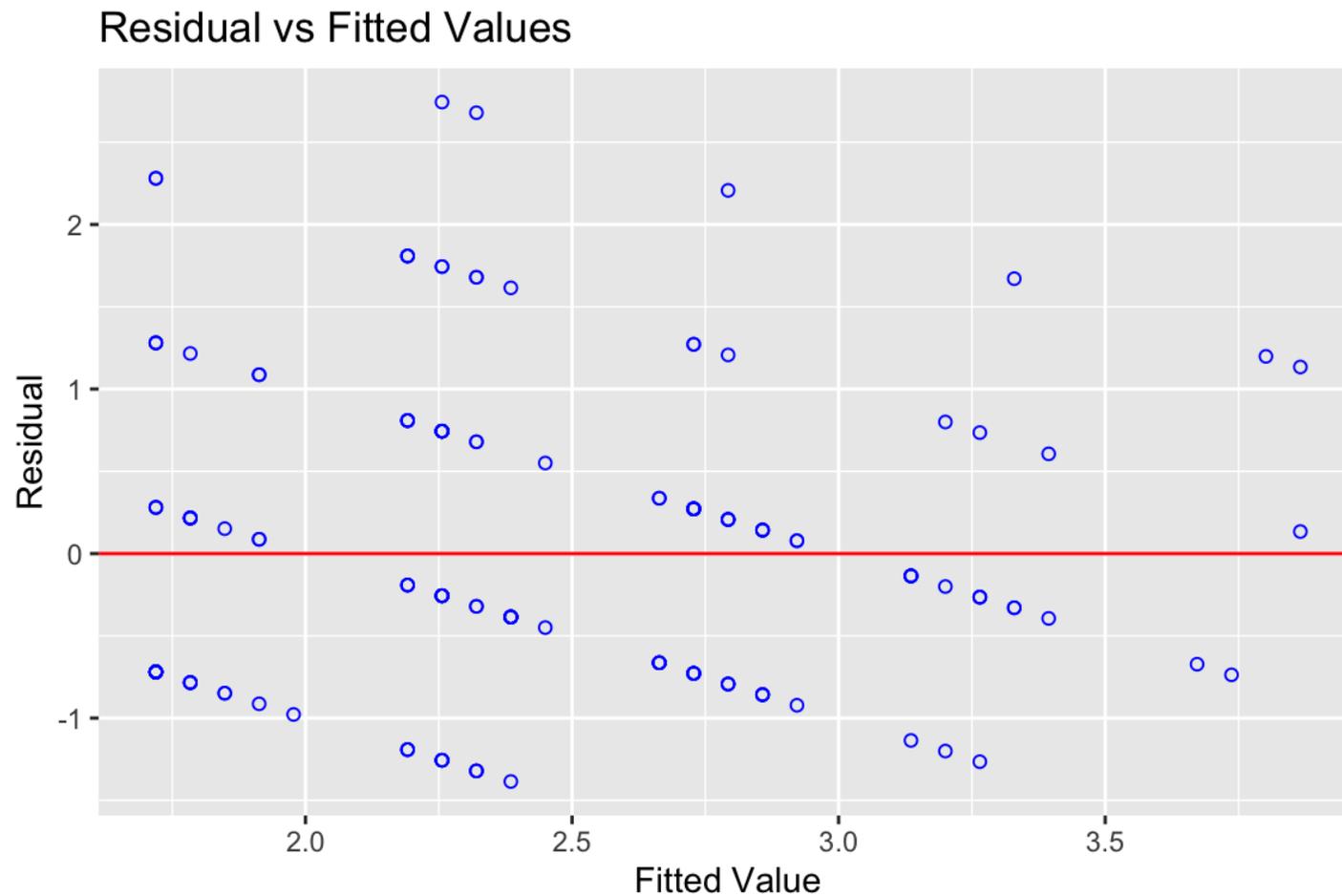
```

1  olsrr::ols_regress(Modell1)
2  ##                               Model Summary
3  ## -----
4  ## R                               0.509           RMSE           0.915
5  ## R-Squared                       0.260           Coef. Var       37.790
6  ## Adj. R-Squared                   0.250           MSE            0.837
7  ## Pred R-Squared                   0.235           MAE            0.736
8  ## -----
9  ## RMSE: Root Mean Square Error
10 ## MSE: Mean Square Error
11 ## MAE: Mean Absolute Error
12 ##
13 ##                               ANOVA
14 ## -----
15 ##                               Sum of
16 ##                               Squares           DF           Mean Square           F           Sig.
17 ## -----

```

Ü1.9 Schauen Sie sich die Residualplotts an

► R-Code anzeigen



Plot der Residuen



Ü1.10 Testen Sie auf Homoskedastizität

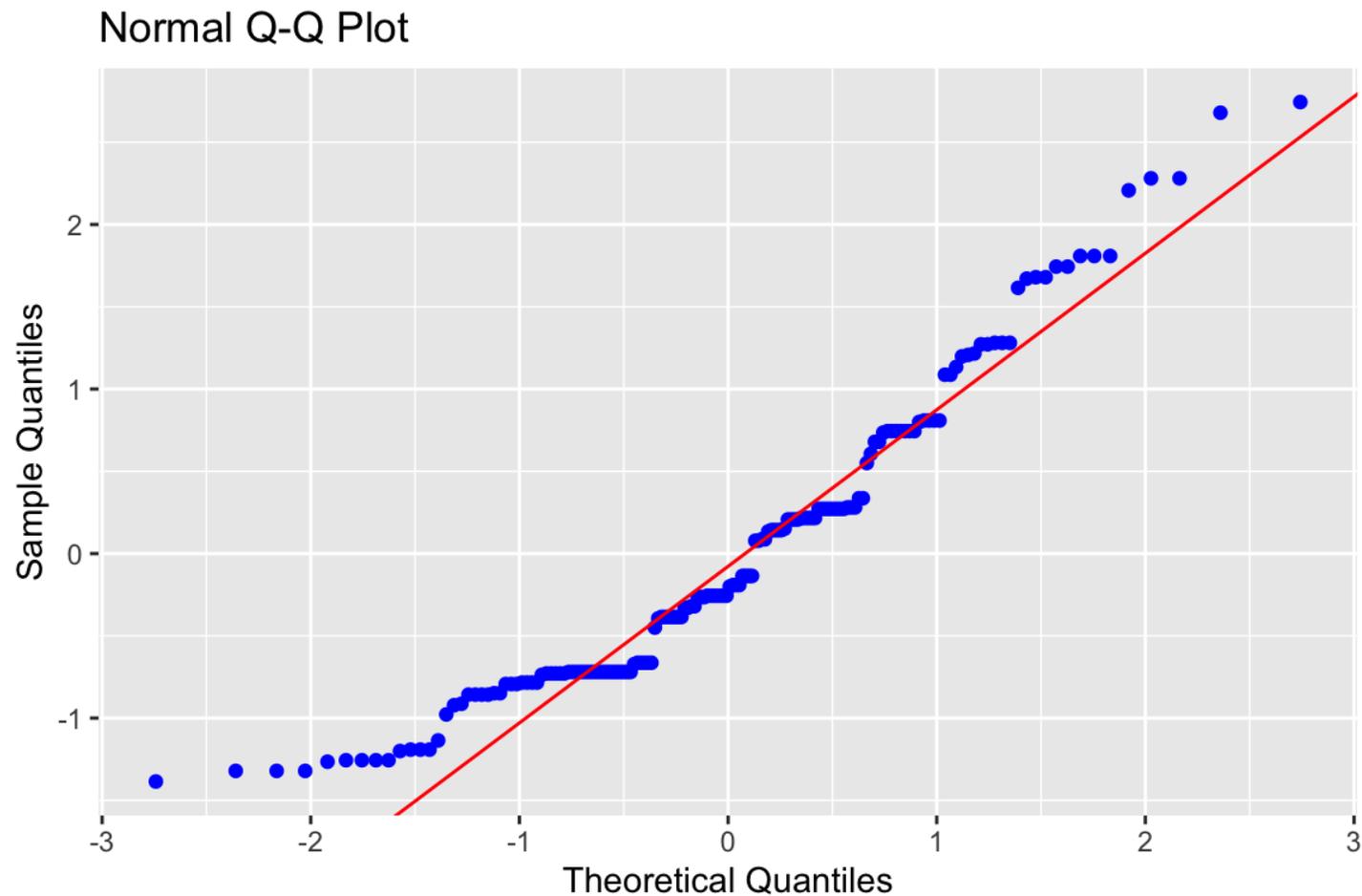
► R-Code anzeigen

Was sagt Ihnen das?



Ü1.11 Gucken Sie sich den N-Q-Q-Plot an

► R-Code anzeigen

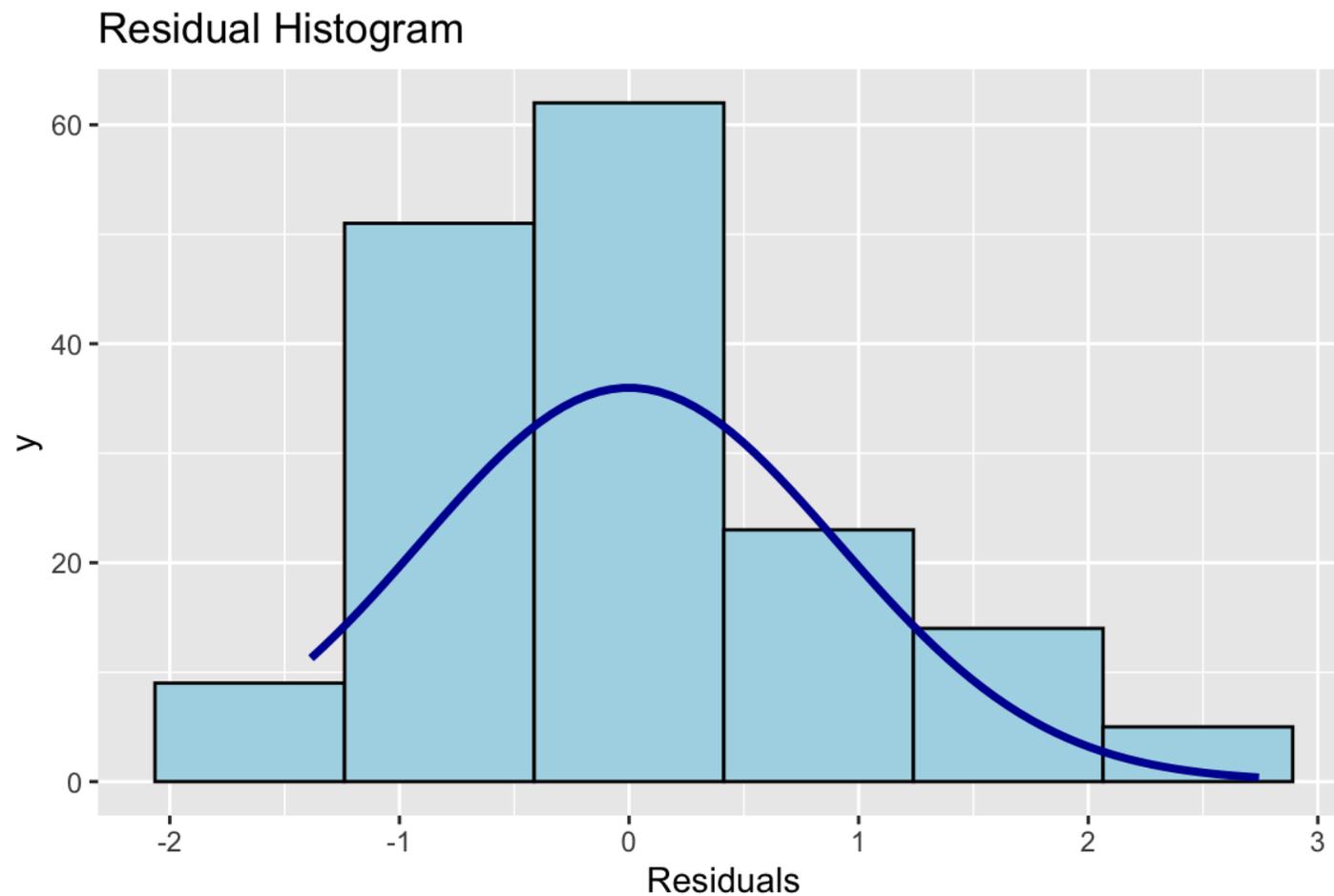


Normal-Q-Q-Plot



Ü1.12 Und das Histogramm

► R-Code anzeigen



Histogramm der Residuen



Ü1.13 Jetzt auf Normalverteilung testen

```
1 # Führe Tests auf signifikante Verletzungen
2 # der Normalverteilungsannahme aus.
3
4 olsrr::ols_test_normality(Modell1)
5 ## -----
6 ##           Test                Statistic         pvalue
7 ## -----
8 ## Shapiro-Wilk                0.9349          0.0000
9 ## Kolmogorov-Smirnov          0.127           0.0101
10 ## Cramer-von Mises           15.381          0.0000
11 ## Anderson-Darling          3.2969          0.0000
12 ## -----
```



Ü1.14 Fazit

Was ist Ihr Fazit aus der Regressionsrechnung?



Weiterführung

Predictors	B	BETA	std.err	t	p
(Intercept)	1.57		0.32	4.89	<.001
E102_02	0.47	.480	0.07	6.71	<.001
E102_04	-0.06	-.076	0.06	-1.06	.292

^a $R^2 = 0.25$

($F = 28$, $df =$

161 , $p =$

161),

$R^2_{adj.} = 0.25$



Take Home – Ausblick – Vokabeln



Take Home

Note

- Sie können eine Regressionsanalyse in Quarto berechnen und die Ergebnisse interpretieren
- Ihnen ist klar, was es bedeutet, dass Dritteinflüsse herausgerechnet werden.
- Sie können Residualanalysen anschauen und erkennen, wann es Probleme gibt (mögliche Lösungen kommen später)
- Sie können den notwendigen Code für Regressionsanalysen in Ihre Projekte kopieren und an den richtigen Stellen anpassen.



Ausblick

Wie beschäftigen uns mit kategorialen Variablen in den UVs.



Vokabeln

Search:

Nr	Sitzung	Inhalt	Deutsch	Englisch	Erläuterung
<input type="text" value="All"/>					
46	3	Voraussetzungen	Allgemeine kleinste Quadrate	Generalized Least Squares (GLS)	Schätzmethode OLS ersetzen wenn es Heteroskedastizität gibt, also die Residuen nicht gleichmäßig streuen.
43	3	Voraussetzungen	BLUE	BLUE	Akronym für Linear Unbiased Estimator
42	3	Voraussetzungen	Bias	Bias	Grad der Verzerrung



Nr	Sitzung	Inhalt	Deutsch	Englisch	Erläuterung
44	3	Voraussetzungen	Effizienz	efficiency	Die Genauigkeit eines Schätzers, also wie stark er streut.
45	3	Voraussetzungen	Fehlervarianz	error variance	Streuung eines Kennwertes (b 's).
47	3	Voraussetzungen	Heteroskedastizität	heteroscedasticity	Die Residuen streuen nicht gleichmässig nach grösser UV. Also hängt die Streubreite der Residuen mit der Grösse einer UV zusammen.
48	3	Voraussetzungen	Homoskedastizität	homoscedasticity	Die Residuen streuen



Nr	Sitzung	Inhalt	Deutsch	Englisch	Erläuterung
					gleichmässig
49	3	Voraussetzungen	Modellspezifikation	model specification	Formulierung Modells, also welche UVs f AV wichtig si und wie deren Beziehung gestaltet ist.
50	3	Voraussetzungen	Multikollinearität	multicollinearity	Eine UV hängt einer oder mehreren der übrigen UVs zusammen.
51	3	Voraussetzungen	Toleranz (TOL)	tolerance	Die übrige Va die eine Variä noch hat, wei gemeinsame Varianz mit a



Nr	Sitzung	Inhalt	Deutsch	Englisch	Erläuterung
					anderen UVs rausgerechnet wurde.
53	3	Voraussetzungen	Unterspezifikation	under estimation	Zu wenige UV Modell (hat verzerrt b's zu Folge).
54	3	Voraussetzungen	Unverzerrtheit	unbiasedness	Eigenschaft einer Methode valid Messungen c Kennwerte für Parameter zu messen.
55	3	Voraussetzungen	Varianzinflationsfaktor (VIF)	variance inflation factor	Der Faktor, um die Ungenauigkeit (Fehlervarianz) einer UV steigt wenn

