

# Lösungen zu den Aufgaben

## 1. Aufgabe

Die Varianzanalyse (ANOVA) ist ein inferenzstatistisches Verfahren des Frequentismus. Welches Skalenniveau passt zu diesem Verfahren?

- a. UV: nominal (mehrstufig), AV: metrisch
- b. UV: metrisch, AV: nominal (zweistufig)
- c. UV: nominal (mehrstufig), AV: nominal (mehrstufig)
- d. UV: metrisch, AV: nominal (zweistufig)
- e. UV: nominal (zweistufig), AV: metrisch

## Lösung

- a. Wahr
- b. Falsch
- c. Falsch
- d. Falsch
- e. Falsch

## 2. Aufgabe

Welches der folgenden Beispiele ist *kein* Beispiel für eine Nullhypothese?

- a.  $\beta_1 \leq 0$
- b.  $\mu_1 = \mu_2$
- c.  $\mu_1 = \mu_2 = \dots = \mu_k$
- d.  $\rho = 0$
- e.  $\pi_1 = \pi_2$

## Lösung

- a. Wahr
- b. Falsch
- c. Falsch
- d. Falsch
- e. Falsch

## 3. Aufgabe

Der t-Test ist ein inferenzstatistisches Verfahren des Frequentismus. Welches Skalenniveau passt zu diesem Verfahren?

- a. UV: nominal (mehrstufig), AV: metrisch

- b. UV: metrisch, AV: nominal (zweistufig)
- c. UV: nominal (mehrstufig), AV: nominal (mehrstufig)
- d. UV: metrisch, AV: nominal (zweistufig)
- e. UV: nominal (zweistufig), AV: metrisch

## Lösung

- a. Falsch
- b. Falsch
- c. Falsch
- d. Falsch
- e. Wahr

## 4. Aufgabe

Für Statistiken (Stichprobe) verwendet man meist lateinische Buchstaben; für Parameter (Population) verwendet man entsprechend meist griechische Buchstaben.

Vervollständigen Sie folgende Tabelle entsprechend!

| Kennwert             | Statistik               | Parameter |
|----------------------|-------------------------|-----------|
| Mittelwert           | $\bar{X}$               | NA        |
| Mittelwertsdifferenz | $\bar{X}_1 - \bar{X}_2$ | NA        |
| Streuung             | sd                      | NA        |
| Anteil               | p                       | NA        |
| Korrelation          | r                       | NA        |
| Regressionsgewicht b |                         | NA        |

## Lösung

| Kennwert             | Statistik               | Parameter |
|----------------------|-------------------------|-----------|
| Mittelwert           | $\bar{X}$               | $\mu$     |
|                      |                         | $\mu_1$   |
| Mittelwertsdifferenz | $\bar{X}_1 - \bar{X}_2$ | $\mu_2$   |
| Streuung             | sd                      | $\sigma$  |
| Anteil               | p                       | $\pi$     |
| Korrelation          | r                       | $\rho$    |

| Kennwert | Statistik | Parameter |
|----------|-----------|-----------|
|----------|-----------|-----------|

|                      |  |         |
|----------------------|--|---------|
| Regressionsgewicht b |  | $\beta$ |
|----------------------|--|---------|

## 5. Aufgabe

Das Testen von Nullhypothesen wird u.a. deswegen kritisiert, weil die Nullhypothese zumeist apriori als falsch bekannt ist, weswegen es keinen Sinn macht, sie zu testen.

Nennen Sie ein Verfahren von John Kruschke, das einen Äquivalenzbereich testet und insofern eine Alternative zum Testen von Nullhypothesen anbietet.

Hinweise:

- Geben Sie nur Kleinbuchstaben ein.
- Geben Sie nur ein einziges Wort ein.

## Lösung

rope

## 6. Aufgabe

Diagramm Diagramm A

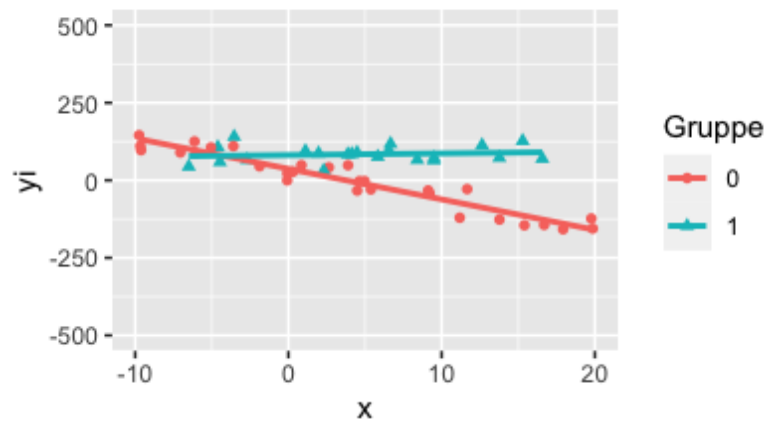


Diagramm Diagramm B

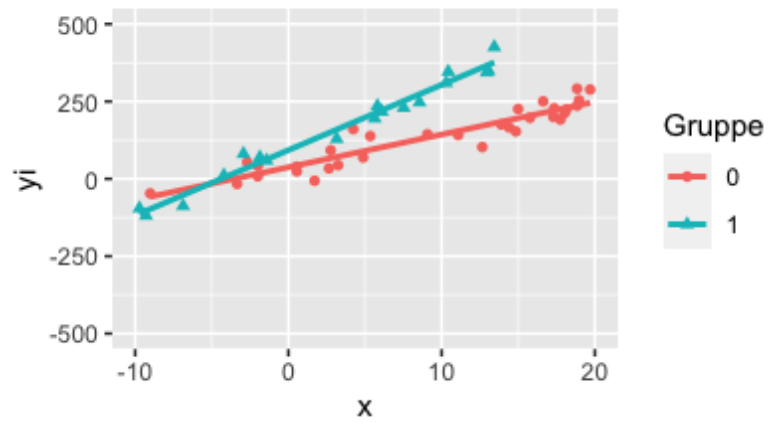
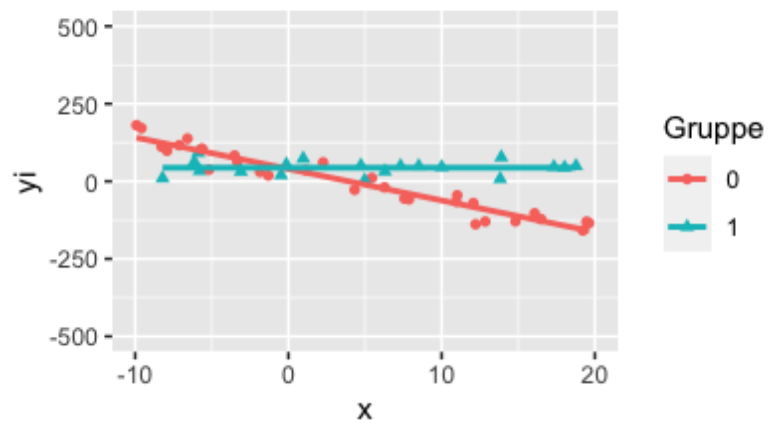
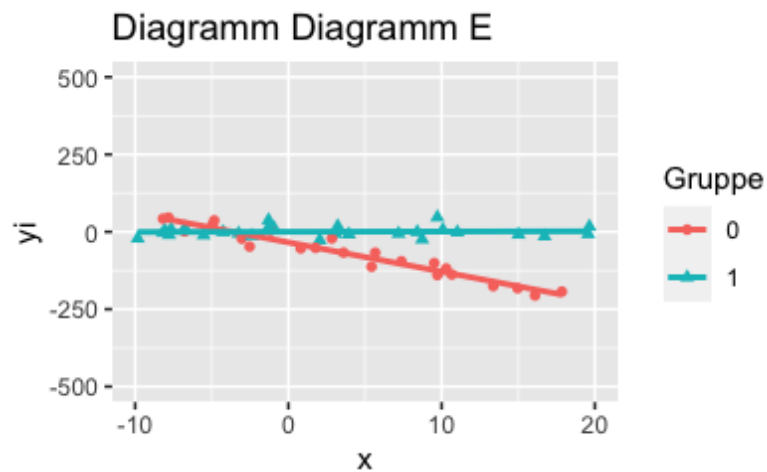
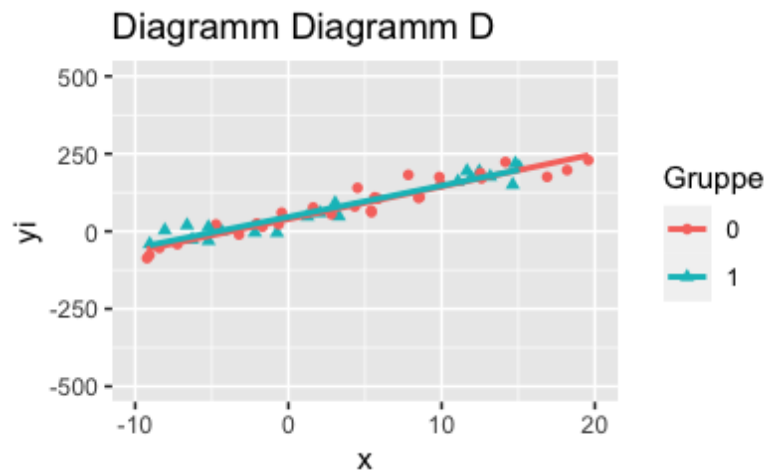


Diagramm Diagramm C





Wählen Sie das Diagramm, in dem *kein* Interaktionseffekt (in der Population) vorhanden ist (bzw. wählen Sie Diagramm, dass dies am ehesten darstellt).

- a. Diagramm A
- b. Diagramm B
- c. Diagramm C
- d. Diagramm D
- e. Diagramm E

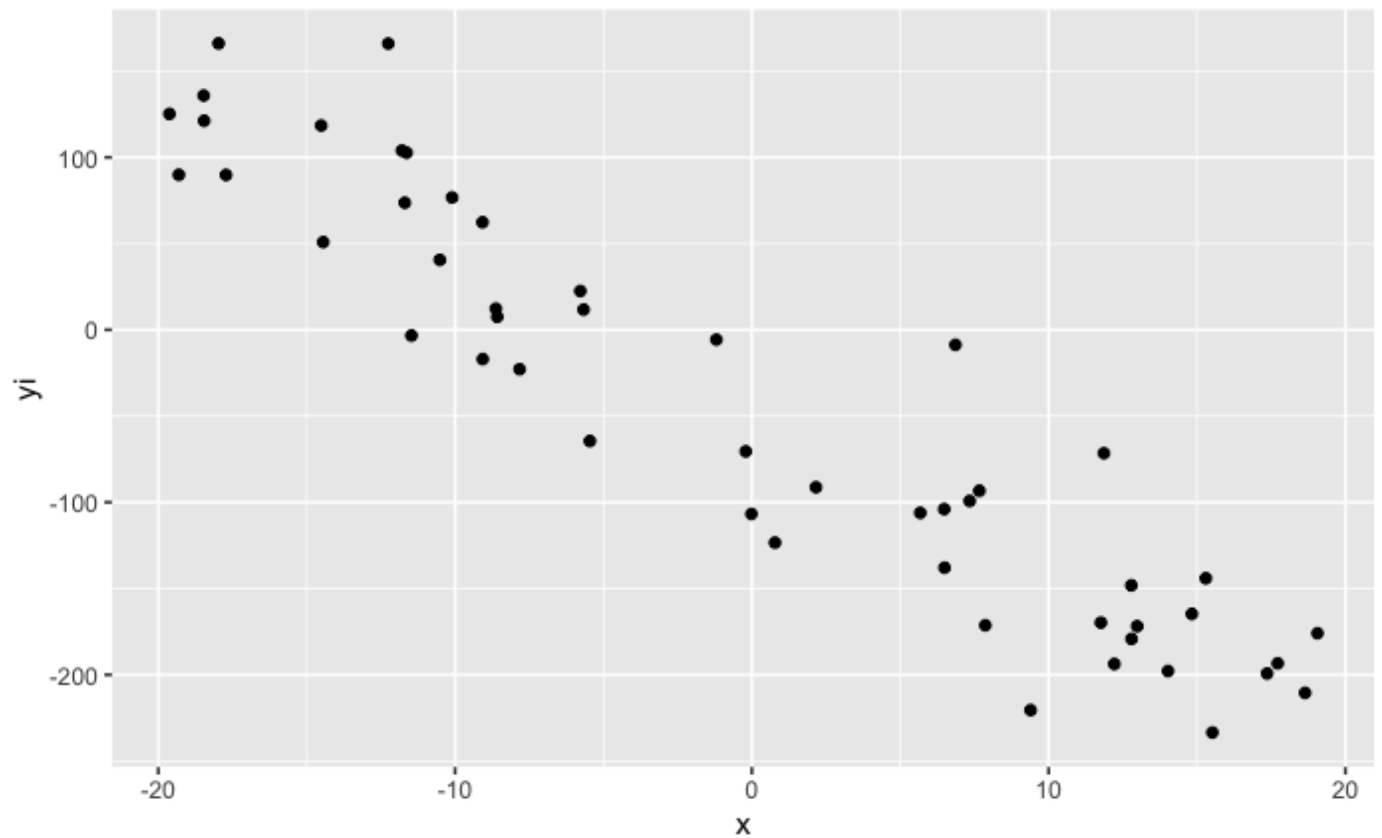
## Lösung

Das Streudiagramm Diagramm D zeigt *keinen* Interaktionseffekt.

- a. Falsch
- b. Falsch
- c. Falsch
- d. Wahr
- e. Falsch

## 7. Aufgabe

Ein Streudiagramm von  $x$  und  $y$  ergibt folgende Abbildung:

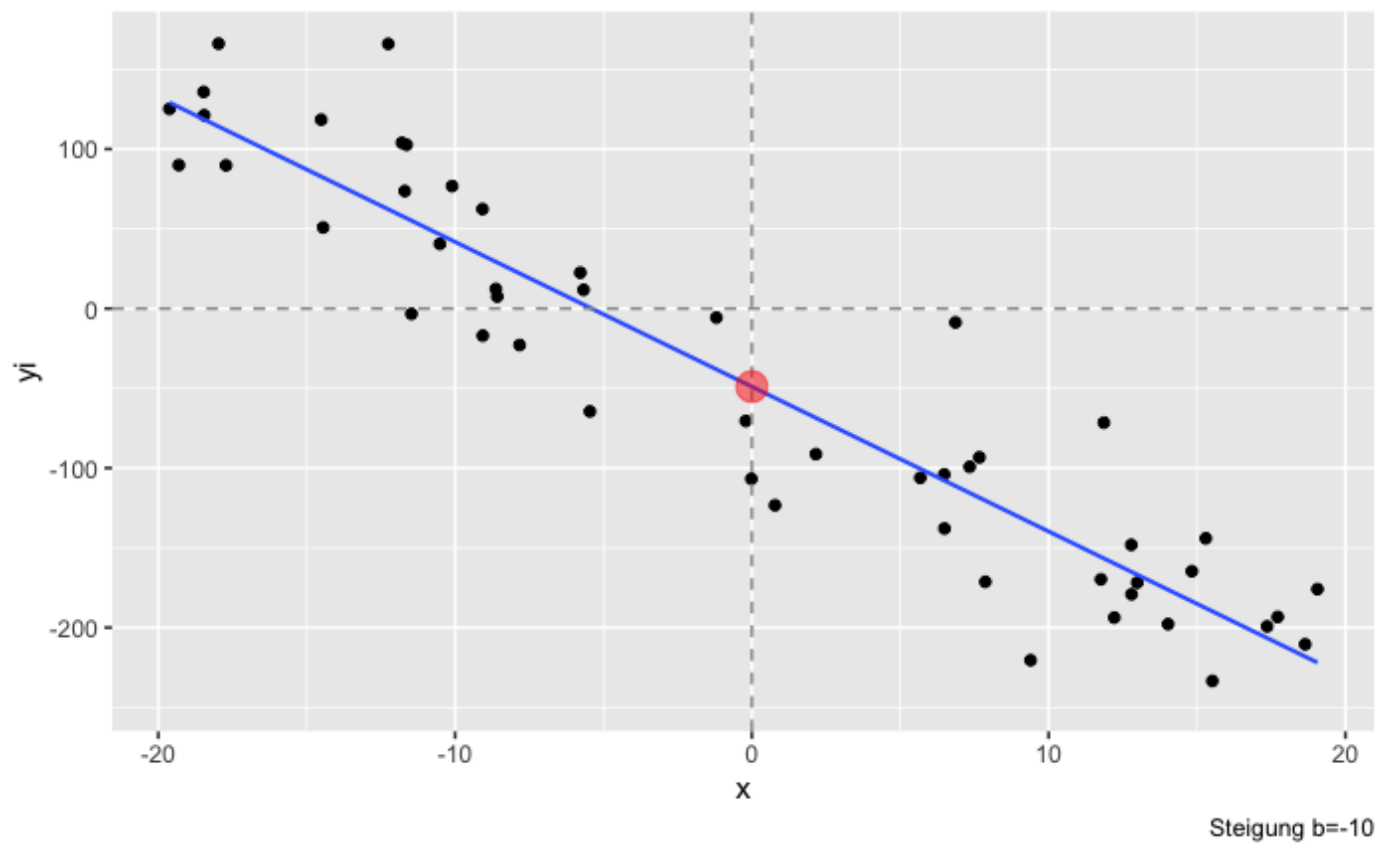


Wählen Sie das am besten passende Modell aus der Liste aus!

- a.  $y = 40 + -10 \cdot x + \epsilon$
- b.  $y = -40 + 10 \cdot x + \epsilon$
- c.  $y = 40 + 10 \cdot x + \epsilon$
- d.  $y = -40 + -10 \cdot x + \epsilon$
- e.  $y = 0 + -40 \cdot x + \epsilon$

### Lösung

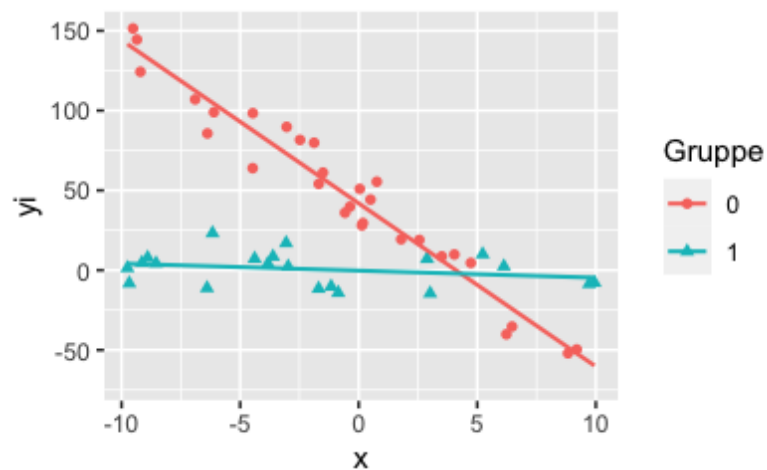
Das dargestellte Modell lautet  $y = -40 + -10 \cdot x + \epsilon$ .



- a. Falsch
- b. Falsch
- c. Falsch
- d. Richtig
- e. Falsch

## 8. Aufgabe

Ein Streudiagramm von  $x$  und  $y$  ergibt folgende Abbildung; dabei wird noch die Gruppierungsvariable  $g$  (mit den Stufen 0 und 1) berücksichtigt (vgl. Farbe und Form der Punkte). Zur besseren Orientierung ist die Regressionsgerade pro Gruppe eingezeichnet.



Wählen Sie das (für die Population) am besten passende Modell aus der Liste aus!

*Hinweis:* Ein Interaktionseffekt der Variablen  $x$  und  $g$  ist mit  $xg$  gekennzeichnet.

- a.  $y = 40 + -10 \cdot x + -40 \cdot g + 10 \cdot xg + \epsilon$
- b.  $y = -40 + -10 \cdot x + -40 \cdot g + -10 \cdot xg + \epsilon$
- c.  $y = 40 + -10 \cdot x + 0 \cdot g + -10 \cdot xg + \epsilon$
- d.  $y = -40 + 10 \cdot x + -40 \cdot g + 10 \cdot xg + \epsilon$

## Lösung

Das dargestellte Modell lautet  $y = 40 + -10 \cdot x + -40 \cdot g + 10 \cdot xg + \epsilon$ . Der Modellfehler  $\epsilon$  hat den Anteil 0.3 im Vergleich zur Streuung von  $y$ .

- a. Richtig
- b. Falsch
- c. Falsch
- d. Falsch

## 9. Aufgabe

Gegeben sei ein Datensatz mit folgenden Prädiktoren, wobei Studierende die Beobachtungseinheit darstellen:

- $X_1$ : Muttersprachler (0: nein, 1: ja)
- $X_2$ : Abitur-Durchschnitt (z-Wert)
- $X_3$ : Alter (z-Wert)
- $X_4$ : Interaktion von  $X_1$  und  $X_2$

Die vorherzusagende Variable ( $Y$ ; Kriterium) ist *Gehalt nach Studienabschluss*.

Folgende Modellparameter einer Regression (Least Squares) seien gegeben:

- $\beta_0$  : 30
- $\beta_1$  : 20
- $\beta_2$  : 1
- $\beta_3$  : 15
- $\beta_4$  : 5

Welche der Aussagen ist korrekt?

- a. Für einen bestimmten (festen) Wert von  $X_2 = \text{Abitur-Durchschnitt (z-Wert)}$  und  $X_3 = \text{Alter (z-Wert)}$  gilt, dass das Gehalt im Mittel höher ist bei  $X_1 = 1$  im Vergleich zu  $X_1 = 0$ , laut dem Modell.
- b. Für einen bestimmten (festen) Wert von  $X_2 = \text{Abitur-Durchschnitt (z-Wert)}$  und  $X_3 = \text{Alter (z-Wert)}$  gilt, dass das Gehalt im Mittel höher ist bei  $X_1 = 0$  im Vergleich zu  $X_1 = 1$ , laut dem Modell.
- c. Der mittlere Gehaltsunterschied  $Y$  zweier Personen  $a$  und  $b$ , wobei bei Person  $a$  gilt  $X_1 = 0$  und bei Person  $b$  gilt  $X_1 = 1$ , beträgt stets 30, laut dem Modell.
- d. Der mittlere Gehaltsunterschied  $Y$  zweier Personen  $a$  und  $b$ , wobei bei Person  $a$  gilt  $X_1 = 0$  und bei Person  $b$  gilt  $X_1 = 1$ , kann *nicht* ohne weitere Angaben auf eine Zahl fixiert



werden, laut dem Modell.

- e. Der mittlere Gehaltsunterschied von Menschen ist eine Wirkung von genau drei Ursachen: Muttersprachler (0: nein, 1: ja), Abitur-Durchschnitt (z-Wert), Alter (z-Wert), laut dem Modell.

## Lösung

- o Wahr
- o Falsch
- o Falsch
- o Falsch
- o Falsch

## 10. Aufgabe

Welches Ergebnis hat der R-Befehl `posterior_interval()` (R-Paket `rstanarm`)?

Wählen Sie die (am besten) passende Antwort aus.

Hinweis:

- o Soweit nicht anders benannt, ist immer die Voreinstellung der betreffenden Funktion gemeint.
- a. Er liefert einen Vorhersagewert aus der Posteriori-Verteilung.
- b. Er liefert ein Vorhersageintervall aus der Posteriori-Verteilung.
- c. Er liefert ein 90%-Vorhersageintervall aus der Posteriori-Verteilung.
- d. Er liefert ein 95%-Vorhersageintervall aus der Posteriori-Verteilung.
- e. Er liefert ein HDI-Vorhersageintervall aus der Posteriori-Verteilung.

## Lösung

So können Sie sich Hilfe zu diesem Befehl ausgeben lassen:

```
help(posterior_interval)
```

- a. Falsch
- b. Falsch
- c. Wahr
- d. Falsch
- e. Falsch

## 11. Aufgabe

Betrachten Sie folgende Ausgabe eines Bayesmodell, das mit `rstanarm` "gefittet" wurde:

```
## stan_glm
## family:      gaussian [identity]
## formula:     price ~ cut
## observations: 1000
```

```
## predictors: 5
## -----
##           Median MAD_SD
## (Intercept)  4571.7   675.1
## cutGood      -570.2   777.2
## cutIdeal     -1288.3   688.1
## cutPremium    362.5   709.8
## cutVery Good -807.4   706.3
##
## Auxiliary parameter(s):
##           Median MAD_SD
## sigma 3795.0    82.4
```

Welche Aussage passt (am besten)?

Hinweise:

- Mit "Nullhypothese" ist im Folgenden dieser Ausdruck gemeint:  $\mu_1 = \mu_2 = \dots = \mu_k$ .
- a. Die Nullhypothese muss verworfen werden.
- b. Die Nullhypothese muss beibehalten werden.
- c. Unter der Annahmen, dass die Posteriori-Verteilung für jeden Regressionsparameter normal verteilt ist, kann man schließen, dass bei allen Gruppenmittelwerte 95% der Posteriori-Verteilung ungleich Null sind.
- d. Unter der Annahmen, dass die Posteriori-Verteilung für jeden Regressionsparameter normal verteilt ist, kann man schließen, dass *nicht* bei allen Gruppenmittelwerte 95% der Posteriori-Verteilung ungleich Null sind.
- e. Unter der Annahmen, dass die Posteriori-Verteilung für jeden Regressionsparameter normal verteilt ist, kann man schließen, dass bei allen Gruppenmittelwerte 95% der Posteriori-Verteilung ungleich Null sind. Daher muss die Nullhypothese verworfen werden.

## Lösung

- a. Falsch
- b. Falsch
- c. Falsch
- d. Wahr
- e. Falsch

## 12. Aufgabe

Berechnet man eine Posteriori-Verteilung mit `stan_glm()`, so kann man entweder die schwach informativen Prioriwerte der Standardeinstellung verwenden, oder selber Prioriwerte definieren.

Betrachten Sie dazu dieses Modell:

```
stan_glm(price ~ cut, data = diamonds,
          prior = normal(location = c(100, 100, 100, 100),
                          scale = c(100, 100, 100, 100)),
          prior_intercept = normal(3000, 500))
```

Welche Aussage dazu passt (am besten)?

- a. Es wird für (genau) einen Parameter eine Priori-Verteilung definiert.

- b. Für keinen Parameter liegt apriori die Wahrscheinlichkeit für einen negativen Wert bei mehr als 5%.
- c. Mit `prior = normal()` werden Gruppenmittelwerte definiert.
- d. Alle Parameter des Modells sind normalverteilt.
- e. mit `prior_intercept = normal(3000, 500)` wird praktisch eine Gleichverteilung definiert (da die Streuung sehr hoch ist).

## Lösung

- a. Falsch
- b. Wahr
- c. Falsch
- d. Falsch
- e. Falsch

## 13. Aufgabe

Berechnet man eine Posteriori-Verteilung mit `stan_glm()`, so kann man entweder die schwach informativen Prioriwerte der Standardeinstellung verwenden, oder selber Prioriwerte definieren.

Betrachten Sie dazu dieses Modell:

```
stan_glm(price ~ cut, data = diamonds,
          prior = normal(location = c(100, 100, 100, 100),
                           scale = c(100, 100, 100, 100)),
          prior_intercept = normal(3000, 500))
```

Wie viele Parameter gibt es in diesem Modell?

Hinweise:

- Geben Sie nur eine (ganze) Zahl ein.

## Lösung

Die Anzahl der Parameter in diesem Modell ist: 11

- Achsenabschnitt: 2 Parameter (MW, SD)
- 4 Regressionsparameter: je 2 Parameter (MW, SD)
- Sigma (Streuung der y-Werte): 1 Parameter (Rate lambda)

## 14. Aufgabe

Sei  $X \sim \mathcal{N}(42, 7)$  und  $x_1 = 28$ .

Berechnen Sie den z-Wert für  $x_1$ !

Hinweis:

- Runden Sie ggf. auf die nächste ganze Zahl.

## Lösung

$$x1\_z = (x1 - x\_mw) / x\_sd$$

-2

## 15. Aufgabe

John Kruschke hat einen (Absolut-)Wert vorgeschlagen, als Grenze für Regressionskoeffizienten "vernachlässigbarer" Größe.

Nennen Sie diesen Wert!

Hinweise:

- Geben Sie nur Zahlen ein (und ggf. Dezimaltrennzeichen).
- Führende Nullen dürfen auch bei Zahlen kleiner als 1 nicht weggelassen werden.

## Lösung

0.05