

Checkbox für die Version vom 27. Dezember 2022

Die gesamte Klausur beinhaltet aktuell in Summe **109** Fragen. Davon sind **50** Multiple Choice Fragen sowie **59** Rechen- und Textaufgaben.

Frequently asked questions (FAQ)

Was ist das hier? Im Folgenden findet sich die Sammlung *aller* Klausurfragen der Bio Data Science über *alle* Veranstaltungen, die ich an der Fakultät für Agrarwissenschaften und Landschaftsarchitektur anbiete.

Sind aber ein bisschen viele Fragen... Ja, das stimmt. Die Sortierung und Überlegung welche Fragen zur Veranstaltung passen obliegt dem Studierenden. Gerne stehe ich für Rückfragen bereit. Teilweise sind Fragen auch ähnlich.

Sind die Fragen fix? Ein klares Jein. Die Zahlen und die *Reihenfolge* der Aufgaben - auch im Multiple Choice Teil - werden sich ändern, da die Klausurfragen zufällig erstellt werden. Die *Aufgabenfragen* hindoch werden die gleichen Fragen bleiben.

Okay, aber woher weiß ich jetzt welche Fragen zu meiner Veranstaltung gehören? Das ist der Trick. Durch das Durchlesen und das selbstständige Sortieren der Fragen zu Themen und Inhalten merkt man ziemlich schnell, welche Inhalte zu der Veranstaltung gehören und welche nicht. Ist also alles Teil des Lernprozesses. *Und* wenn Unsicherheiten da sind, gerne in der Wiederholungsveranstaltung - letzte Vorlesung - einfach mich fragen.

Wie sieht denn die finale Klausur aus? Die Klausur hat am Ende 10 Multiple Choice Fragen mit jeweils 2 Punkten sowie Rechen- und Textaufgaben mit in Summe ca. 60 Punkten. Ich peile daher eine Klausur mit 80 Punkten an, wobei 40 Punkte zum Bestehen der Klausur notwendig sind. **Bei geteilten Veranstaltungen mit mehr als einem Dozenten ändert sich die Zusammensetzung der endgültigen Punkteanzahl!**

Sind aber mehr als zehn Multiple Choice Fragen... Ja, aber es werden in der finalen Klausur nur zehn Multiple Choice Fragen sein. Ich wähle die Fragen dann zufällig aus. Ich berücksichtige natürlich die Veranstaltung und das Lernniveau.

Solange kann ich nicht warten... Dann einfach eine Mail an mich schreiben:
j.kruppa@hs-onsabrueck.de

Ich versuche dann die Frage kurzfristig zu beantworten oder aber in der Vorlesung nochmal (anonym) aufzugreifen.

Name: _____

Nicht bestanden: ☐

Vorname: _____

Matrikelnummer: _____

Endnote: _____

Klausurfragen der Bio Data Science

Hochschule Osnabrück

Prüfer: Prof. Dr. Jochen Kruppa
Fakultät für Agrarwissenschaften und Landschaftsarchitektur
j.kruppa@hs-osnabrueck.de

Version vom 27. Dezember 2022

Erlaubte Hilfsmittel für die Klausur

- Normaler Taschenrechner ohne Möglichkeit der Kommunikation mit anderen Geräten - also ausdrücklich kein Handy!
- Eine DIN A4-Seite als beidseitig, selbstgeschriebene, handschriftliche Formelsammlung - keine digitalen Ausdrucke.

Ergebnis der Klausur

_____ von 20 Punkten sind aus dem Multiple Choice Teil erreicht.
_____ von 60 Punkten sind aus dem Rechen- und Textteil erreicht.
_____ von 80 Punkten in Summe.

Es wird folgender Notenschlüssel angewendet.

Punkte	Note
78 - 80	1,0
75 - 77	1,3
70 - 74	1,7
65 - 69	2,0
59 - 64	2,3
54 - 58	2,7
49 - 53	3,0
44 - 48	3,3
41 - 43	3,7
40	4,0

Es ergibt sich eine Endnote von _____.

Multiple Choice Aufgaben

- Pro Multiple Choice Frage ist *genau* eine Antwort richtig.
- **Übertragen Sie Ihre Kreuze in die Tabelle auf dieser Seite.**
- Es werden nur Antworten berücksichtigt, die in dieser Tabelle angekreuzt sind!

	A	B	C	D	E	✓
1 Aufgabe						
2 Aufgabe						
3 Aufgabe						
4 Aufgabe						
5 Aufgabe						
6 Aufgabe						
7 Aufgabe						
8 Aufgabe						
9 Aufgabe						
10 Aufgabe						

- Es sind ____ von 20 Punkten erreicht worden.

Rechen- und Textaufgaben

- Die Tabelle wird vom Dozenten ausgefüllt.

Aufgabe	11	12	13	14	15	16	17
Punkte							

- Es sind ____ von 60 Punkten erreicht worden.

Multiple Choice Aufgaben

- Es wird nie mehr als fünfzig Multiple Choice Fragen geben.
- Im Laufe der Zeit werden einzelne Fragen durch andere Fragen *ersetzt*, bitte beachten Sie diesen Sachstand, wenn Sie eine *Wiederholungsklausur* im nächsten Semester schreiben.

1 Aufgabe

(2 Punkte)

Welche Aussage über die α Adjustierung ist richtig?

- A** ☐ Die α Adjustierung wird meist ignoriert. Wenn die Annahmen an den statistischen Test richtig sind, kann auf eine Adjustierung verzichtet werden.
- B** ☐ Die α Adjustierung wird durchgeführt um den Effekt von Interesse, meist die Behandlung, von anderen Effekten zu trennen. Daher eine Adjustierung auf den β -Werten einer Regression.
- C** ☐ Die α ist notwendig um Effekte gegeneinander aufzurechnen. Ohne diese Adjustierung würde der eigentliche Effekt nicht richtig geschätzt. Daher handelt es sich um eine Adjustierung der Fehlerwahrscheinlichkeiten.
- D** ☐ Die α Adjustierung wird durchgeführt um bei multiplen Vergleichen den Fehler 1. Art zu kontrollieren. Es wird die Irrtumswahrscheinlichkeit adjustiert, daher das α -Niveau.
- E** ☐ Die α Adjustierung wird durchgeführt um den Fehler 2. Art zu kontrollieren. Ohne diese Adjustierung würde der Fehler 2. Art nicht bei 80% liegen sondern sehr schnell gegen 0 laufen.

2 Aufgabe

(2 Punkte)

Sie haben folgende unadjustierten p-Werte gegeben: 0.03, 0.001, 0.01 und 0.89. Sie adjustieren die p-Werte nach Bonferroni. Welche Aussage ist richtig?

- A** ☐ Nach der Bonferroni-Adjustierung ergeben sich die adjustierten p-Werte von 0.0075, 0.0003, 0.0025 und 0.2225. Die adjustierten p-Werte werden zu einem α -Niveau von 5% verglichen.
- B** ☐ Nach der Bonferroni-Adjustierung ergeben sich die adjustierten p-Werte von 0.0075, 0.0003, 0.0025 und 0.2225. Die adjustierten p-Werte werden zu einem α -Niveau von 1.25% verglichen.
- C** ☐ Nach der Bonferroni-Adjustierung ergeben sich die adjustierten p-Werte von 0.12, 0.004, 0.04 und 1. Die adjustierten p-Werte werden zu einem α -Niveau von 5% verglichen.
- D** ☐ Nach der Bonferroni-Adjustierung ergeben sich die adjustierten p-Werte von 0.12, 0.004, 0.04 und 3.56. Die adjustierten p-Werte werden zu einem α -Niveau von 5% verglichen.
- E** ☐ Nach der Bonferroni-Adjustierung ergeben sich die adjustierten p-Werte von 0.12, 0.004, 0.04 und 1. Die adjustierten p-Werte werden zu einem α -Niveau von 1.25% verglichen.

3 Aufgabe

(2 Punkte)

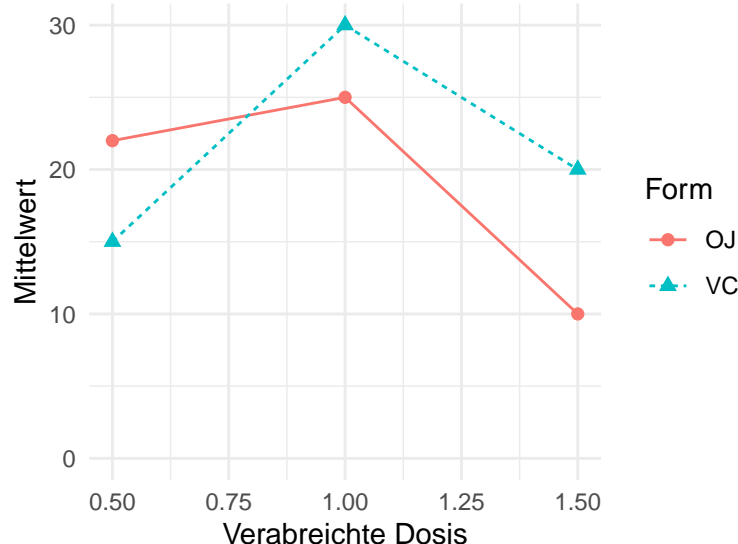
Der Datensatz PlantGrowth enthält das Gewicht von Pflanzen, die unter einer Kontrolle und zwei verschiedenen Behandlungsbedingungen erzielt wurden. Nach der Berechnung einer einfaktoriellen ANOVA ergibt sich ein $\eta^2 = 0.27$. Welche Aussage ist richtig?

- A** ☐ Das η^2 ist ein Wert für die Güte der ANOVA. Je kleiner desto besser. Ein η^2 von 0 bedeutet ein perfektes Modell mit keiner Abweichung. Die Varianz ist null.
- B** ☐ Die Berechnung von η^2 ist ein Wert für die Interaktion.
- C** ☐ Das η^2 beschreibt den Anteil der Varianz, der von den Behandlungsbedingungen nicht erklärt wird. Somit der Rest an nicht erklärbarer Varianz.
- D** ☐ Das η^2 ist die Korrelation der ANOVA. Mit der Ausnahme, dass 0 der beste Wert ist.
- E** ☐ Das η^2 beschreibt den Anteil der Varianz, der von den Behandlungsbedingungen erklärt wird. Das η^2 ist damit mit dem R^2 aus der linearen Regression zu vergleichen.

4 Aufgabe

(2 Punkte)

Die folgende Abbildung enthält die Daten aus einer Studie zur Bewertung der Wirkung von Vitamin C auf das Zahnwachstum bei Meerschweinchen. Der Versuch wurde an 60 Schweinen durchgeführt, wobei jedes Tier eine von drei Vitamin-C-Dosen (0.5, 1 und 1.5 mg/Tag) über eine von zwei Verabreichungsmethoden mit Orangensaft (OJ) oder Ascorbinsäure (VC) erhielt.



Welche Aussage ist richtig im Bezug auf eine zweifaktorielle ANOVA?

- A** ☐ Keine Interaktion liegt vor. Die Geraden scheiden sich und laufen nicht parallel.
- B** ☐ Eine starke Interaktion ist zu erwarten. Die Geraden schneiden sich und die Abstände sind nicht gleichbleibend.
- C** ☐ Eine starke Interaktion liegt vor. Die Geraden laufen parallel und schneiden sich nicht.
- D** ☐ Eine leichte Interaktion ist zu erwarten. Die Geraden schneiden sich noch nicht, aber die Abstände unterscheiden sich stark.
- E** ☐ Keine Interaktion ist zu erwarten. Die Geraden der Verabreichungsmethode laufen parallel und mit ähnlichen Abständen.

5 Aufgabe

(2 Punkte)

Eine einfaktorielle ANOVA berechnet eine Teststatistik um zu die Nullhypothese abzulehnen. Welche Aussage über die Teststatistik der ANOVA ist richtig?

- A** ☐ Die ANOVA berechnet die F-Statistik indem die MS des Fehlers durch die MS der Behandlung geteilt werden. Wenn die F-Statistik sich der 1 annähert kann die Nullhypothese nicht abgelehnt werden.
- B** ☐ Die ANOVA berechnet die F-Statistik indem die MS der Behandlung durch die MS des Fehlers geteilt werden. Wenn die F-Statistik sich der 0 annähert kann die Nullhypothese nicht abgelehnt werden.
- C** ☐ Die ANOVA berechnet die F-Statistik aus den SS Behandlung geteilt durch die SS Fehler.
- D** ☐ Die ANOVA berechnet die T-Statistik indem den Mittelwertsunterschied der Gruppen simultan durch die Standardabweichung der Gruppen teilt. Wenn die T-Statistik höher als 1.96 ist, kann die Nullhypothese abgelehnt werden.
- E** ☐ Die ANOVA berechnet die T-Statistik aus der Multiplikation der MS Behandlung mit der MS der Fehler. Wenn die F-Statistik 0 ist, kann die Nullhypothese abgelehnt werden.

6 Aufgabe

(2 Punkte)

Sie haben das abstrakte Modell $Y \sim X$ mit X als Faktor mit zwei Leveln vorliegen. Welche Aussage über $s_1^2 \neq s_2^2$ ist richtig?

- A** ☐ Es handelt sich um ein balanciertes Design.
- B** ☐ Es handelt sich um ein unbalanciertes Design
- C** ☐ Es liegt Varianzheterogenität vor.
- D** ☐ Es liegt Varianzhomogenität vor.
- E** ☐ Es handelt sich um abhängige Beobachtungen.

7 Aufgabe

(2 Punkte)

Die Mindestanzahl an Beobachtungen für eine Zelle der Vierfeldertafel bei der Nutzung eines Chi-Quadrat-Testes ist...

- A** ☐ 2 Beobachtungen
- B** ☐ 1 Beobachtung
- C** ☐ 10 Beobachtungen
- D** ☐ 5 Beobachtungen
- E** ☐ 0 Beobachtungen

8 Aufgabe

(2 Punkte)

Welche Aussage über den Korrelationskoeffizienten nach Spearman ist richtig?

- A** ☐ Der Korrelationskoeffizienten nach Spearman wird genutzt, wenn das Outcome Y nicht normalverteilt ist. Der Korrelationskoeffizienten liegt zwischen 0 und 1.
- B** ☐ Der Korrelationskoeffizienten nach Spearman wird genutzt, wenn der Korrelationskoeffizienten zwischen -1 und 1 liegt. Dann sind die Residuen normalverteilt.
- C** ☐ Der Korrelationskoeffizienten nach Spearman wird genutzt, wenn das Outcome Y nicht normalverteilt ist. Der Korrelationskoeffizienten liegt zwischen -1 und 1.

- D** ☐ Der Korrelationskoeffizienten nach Spearman wird genutzt, wenn das Outcome Y normalverteilt ist. Der Korrelationskoeffizienten liegt zwischen 0 und 1.
- E** ☐ Der Korrelationskoeffizienten nach Spearman wird genutzt, wenn das Outcome Y normalverteilt ist. Der Korrelationskoeffizienten liegt zwischen -1 und 1.

9 Aufgabe

(2 Punkte)

Berechnen Sie den Mittelwert und Standardabweichung von y mit 8, 4, 16, 14 und 14.

- A** ☐ Es ergibt sich 10.2 +/- 12.6
- B** ☐ Es ergibt sich 11.2 +/- 5.02
- C** ☐ Es ergibt sich 11.2 +/- 2.51
- D** ☐ Es ergibt sich 11.2 +/- 25.2
- E** ☐ Es ergibt sich 12.2 +/- 2.51

10 Aufgabe

(2 Punkte)

Berechnen Sie den Median und das IQR von x mit 9, 16, 15, 32, 7, 21, 23, 19, 14, 9 und 63.

- A** ☐ Es ergibt sich 16 [9, 23]
- B** ☐ Es ergibt sich 21 [9, 23]
- C** ☐ Es ergibt sich 16 +/- 9
- D** ☐ Es ergibt sich 21 +/- 9
- E** ☐ Es ergibt sich 16 +/- 23

11 Aufgabe

(2 Punkte)

Eine der gängigsten Methode der Statistik um einen Fehler zu bestimmen ist...

- A** ☐ ... das Produkt der kleinsten Quadrate.
- B** ☐ ... die Methode des absoluten Abstands.
- C** ☐ ... die Methode der aufaddierten, absoluten Abstände.
- D** ☐ ... die kleinste Quadrate Methode oder auch least square method genannt.
- E** ☐ ... die Methode des absoluten, quadrierten Abstands.

12 Aufgabe

(2 Punkte)

Welche Aussage über Cook's d und Cohen's d ist richtig?

- A** ☐ Wir nutzen Cook's d um Outlier in den Daten zu finden und Cohen's d um einen standardisierten Effektschätzer für Gruppenvergeliche zu erhalten.
- B** ☐ Wir nutzen Cook's d um Outlier in den Daten zu finden. Cohen's d findet auch Outlier, ist aber ein veraltetes Konzept in der Statistik.

- C** ☐ Wir nutzen Cook's d um Outlier in den Daten zu finden und Cohen's d um einen nicht standardisierten Effektschätzer für Gruppenvergeliche zu erhalten.
- D** ☐ Wir nutzen Cohen's d um Outlier in den Daten zu finden und Cook's d um einen standardisierten Effektschätzer für Gruppenvergeliche zu erhalten.
- E** ☐ Wir nutzen Cook's d um Outlier in den Daten zu finden und Cohen's d um standardisierte Outlier für Gruppenvergeliche zu erhalten.

13 Aufgabe

(2 Punkte)

Die empfohlene Mindestanzahl an Beobachtungen für ein Histogramm sind...

- A** ☐ 5 und mehr Beobachtungen.
- B** ☐ 2-5 Beobachtungen.
- C** ☐ mindestens 20 Beobachtungen.
- D** ☐ 10 Beobachtungen.
- E** ☐ 1 Beobachtung.

14 Aufgabe

(2 Punkte)

Mit einem Dotplot wird Folgendes in der Statistik und der explorativen Datenanalyse hauptsächlich dargestellt.

- A** ☐ Die Verteilung von Daten. Auch 2 bis 5 Beobachtungen können noch dargestellt werden.
- B** ☐ Die Verteilung von Daten. Meistens dem Outcome oder auch y genannt. Hierbei ist eine hohe Fallzahl notwendig mit mehr als 20 Beobachtungen.
- C** ☐ Die Eigenschaften von Daten anhand der fehlenden Werte und den Leveln eines Faktors x .
- D** ☐ Die Verteilung von Daten. Meistens dem Outcome oder auch y genannt. Die Darstellung ist auch mit einer kleinen Fallzahl von Minimum 5 Beobachtungen möglich.
- E** ☐ Die Eigenschaften von Daten aufgeteilt nach zwei Gruppen eines Faktors x .

15 Aufgabe

(2 Punkte)

Nachdem Sie in einem Experiment die Daten D erhoben haben, berechnen Sie den Mittelwert und den Median. Der Mittelwert \bar{y} und der Median \tilde{y} unterscheiden sich nicht. Welche Aussage ist richtig?

- A** ☐ Da sich der Mittelwert und der Median unterscheiden, liegen vermutlich Outlier in den Daten vor. Wir untersuchen den Datensatz nach auffälligen Beobachtungen.
- B** ☐ Da sich der Mittelwert und der Median unterscheiden, liegen vermutlich keine Outlier in den Daten vor.
- C** ☐ Da sich der Mittelwert und der Median nicht unterscheiden, liegen vermutlich Outlier in den Daten vor.
- D** ☐ Da sich der Mittelwert und der Median unterscheiden, ist der Datensatz nicht zu verwenden. Mittelwert und Median müssen gleich sein.
- E** ☐ Da sich der Mittelwert und der Median nicht unterscheiden, liegen vermutlich keine Outlier in den Daten vor. Wir verwenden den Datensatz so wie er ist.

16 Aufgabe

(2 Punkte)

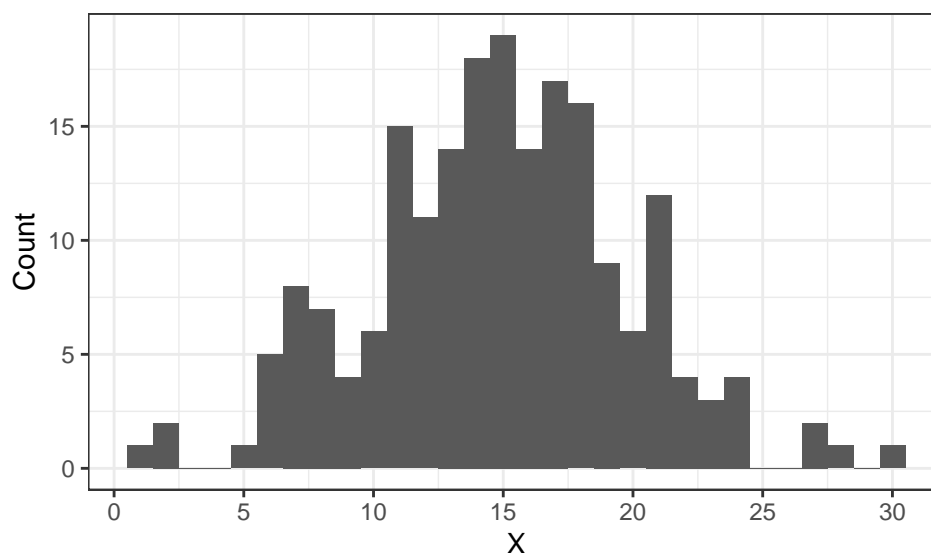
Price et al. (2016) untersuchte die Auswirkungen des Bergbaus und der Talauffüllung auf den Bestand und die Häufigkeit von Bachsalamandern. Um den Effekt zu Berechnen nutze Price et al. (2016) eine Poisson-Regression auf die Anzahl an aufgefundenen Bachsalamandern an den jeweiligen Suchorten. Welche Aussage zur Poisson-Regression auf Zählraten ist richtig?

- A** ☐ Die Poisson-Regression schätzt nur einen Verteilungsparameter. Deshalb muss überprüft werden, ob Overdispersion vorliegt. Mit einer geschätzten Overdispersion von 3.03 liegt Overdispersion vor. Die Lösung ist die Nutzung einer anderen Verteilungsfamilie wie die Quasipoisson Verteilung.
- B** ☐ Die Poisson-Regression schätzt drei Verteilungsparameter. Deshalb muss überprüft werden, ob Overdispersion vorliegt. Mit einer geschätzten Overdispersion von 3.03 liegt Overdispersion vor. Die Lösung ist die Nutzung von nur zwei der drei Verteilungsparameter: γ_1 und γ_3 .
- C** ☐ Die Poisson-Regression schätzt nur einen Verteilungsparameter. Deshalb muss überprüft werden, ob Overdispersion vorliegt. Mit einer geschätzten Overdispersion von 3.03 liegt Overdispersion vor. Damit kann keine Poisson-Regression gerechnet werden. Die Lösung ist eine Gaussian Regression mit Nullanpassung.
- D** ☐ Die Poisson-Regression schätzt zwei Verteilungsparameter. Deshalb muss überprüft werden, ob Overdispersion vorliegt. Mit einer geschätzten Overdispersion von 3.03 liegt keine Overdispersion vor.
- E** ☐ Die Poisson-Regression schätzt nur einen Verteilungsparameter. Deshalb muss überprüft werden, ob Overdispersion vorliegt. Mit einer geschätzten Overdispersion von 3.03 liegt keine Overdispersion vor. Overdispersion liegt vor, wenn die geschätzte Overdispersion unter 1 liegt.

17 Aufgabe

(2 Punkte)

In dem folgenden Histogramm von $n = 200$ Pflanzen ist welche Verteilung mit welchen korrekten Verteilungsparametern dargestellt?



- A** ☐ Eine Standardnormalverteilung mit $N(0,1)$.
- B** ☐ Es handelt sich um eine Binomial-Verteilung mit $\text{Binom}(10)$.

- C** ☐ Es handelt sich um eine Normalverteilung mit $N(15, 5)$.
- D** ☐ Es handelt sich um eine Poisson-Verteilung mit $Pois(15)$.
- E** ☐ Eine rechtsschiefe, multivariate Normalverteilung.

18 Aufgabe

(2 Punkte)

Price et al. (2016) untersuchte die Auswirkungen des Bergbaus und der Talauffüllung auf den Bestand und die Häufigkeit von Bachsalamandern. Um den Effekt zu Berechnen nutze Price et al. (2016) eine Poisson-Regression auf die Anzahl an aufgefundenen Bachsalamandern an den jeweiligen Suchorten. Nach einer statistischen Beratung wurde ihm nahegelegt auf Overdispersion zu achten, wenn er statistische Aussagen zur Signifikanz treffen will. Price et al. (2016) schätzt zwei Modelle. Modell 1 mit einer Poisson Verteilung und Modell 2 mit einer Quasi-Poisson Verteilung. Welche Aussage zu einer geschätzten Overdispersion von 2.74 ist richtig?

- A** ☐ Das Vergleichen von verschiedenen Modellen muss erst über ein AIC Kriterium erfolgen. Die Abschätzung über die Overdispersion ist nicht notwendig. Die Varianzen werden später in einer ANOVA adjustiert. Die Confounder Adjustierung.
- B** ☐ Bei einer geschätzten Overdispersion höher als 1.5 ist von keiner Overdispersion in den Daten auszugehen. Dennoch sind die p-Werte zu klein, dass diese p-Werte natürlich entstanden sein könnten. Die p-Werte müssen adjustiert werden.
- C** ☐ Bei einer geschätzten Overdispersion höher als 1.5 ist von Overdispersion in den Daten auszugehen. Daher wird die Varianz systematisch unterschätzt, was zu kleineren p-Werten führt. Daher gibt es mehr signifikante Ergebnisse als es in Wirklichkeit gibt. Daher ist das Modell 2 die bessere Wahl.
- D** ☐ Bei einer geschätzten Overdispersion höher als 1.5 ist von Overdispersion in den Daten auszugehen. Daher wird die Varianz systematisch unterschätzt, was zu höheren p-Werten führt. Daher gibt es weniger signifikante Ergebnisse als es in Wirklichkeit gibt. Daher ist das Modell 1 die bessere Wahl.
- E** ☐ Bei einer geschätzten Overdispersion höher als 1.5 ist von Overdispersion in den Daten auszugehen. Daher wird die Varianz systematisch überschätzt, was zu höheren p-Werten führt. Daher gibt es mehr signifikante Ergebnisse als es in Wirklichkeit gibt. Daher ist das Modell 1 die bessere Wahl.

19 Aufgabe

(2 Punkte)

In einem Zuchtexperiment messen wir die Ferkel verschiedener Sauen. Die Ferkel einer Muttersau sind daher im statistischen Sinne...

- A** ☐ Untereinander abhängig. Die Ferkel stammen von einem Muttertier und haben vermutlich eine ähnliche Varianzstruktur.
- B** ☐ Untereinander unabhängig. Sollten die Mütter verwandt sein, so ist die Varianzstruktur ähnlich und muss modelliert werden.
- C** ☐ Untereinander stark korreliert. Die Ferkel sind von einer Mutter und somit miteinander korreliert. Dies wird in der Statistik jedoch meist nicht modelliert.
- D** ☐ Untereinander unabhängig. Die Ferkel sind eigenständig und benötigen keine zusätzliche Behandlung.
- E** ☐ Untereinander abhängig, wenn die Mütter ebenfalls miteinander verwandt sind. Erst die Abhängigkeit 2. Grades wird in der Statistik modelliert.

20 Aufgabe

(2 Punkte)

Sie haben das abstrakte Modell $y \sim x_1 + x_2$ vorliegen. Welche Aussage über X ist richtig?

- A ☐ X beinhaltet die Zeilen. Die Zeilen geben die Verteilungsfamilie vor.
- B ☐ X beinhaltet mehrere Spalten. Die Spalten geben die Verteilungsfamilie vor.
- C ☐ X beinhaltet eine Spalte. Die Spalte gibt nicht die Verteilungsfamilie vor.
- D ☐ X beinhaltet mehrere Spalten. Die Spalten enthalten die Behandlung und weitere potenzielle Einflussvariablen
- E ☐ X beinhaltet eine Spalte. Die Spalte gibt die Verteilungsfamilie vor.

21 Aufgabe

(2 Punkte)

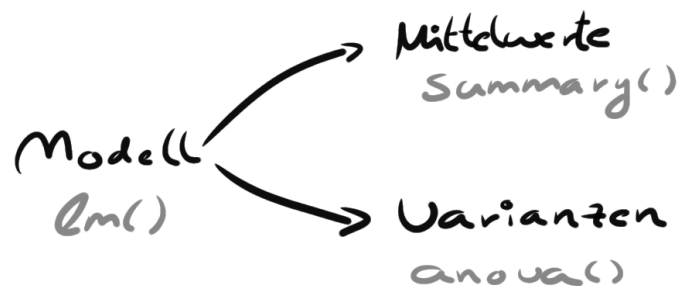
Sie haben das Modell $Y \sim X$ vorliegen und wollen nun ein prädiktives Modell rechnen. Welche Aussage ist richtig?

- A ☐ Ein prädiktives Modell benötigt mindestens eine Fallzahl von über 100 Beobachtungen und darf keine fehlenden Werte beinhalten. Die Varianzkomponenten müssen homogen sein.
- B ☐ Ein prädiktives Modell möchte die Zusammenhänge von X auf Y modellieren. Hierbei geht es um die Effekte von X auf Y . Man sagt, wenn X um 1 ansteigt ändert sich Y um einen Betrag β .
- C ☐ Ein prädiktives Modell basiert auf einem Trainingsdatensatz und einem Testdatensatz. Auf dem Trainingsdatensatz wird das Modell trainiert und auf dem Testdatensatz validiert.
- D ☐ Ein prädiktives Modell schliesst grundsätzlich lineare Modell aus. Es muss ein Graph gefunden werden, der alle Punkte beinhaltet. Erst dann kann das R^2 berechnet werden.
- E ☐ Ein prädiktives Modell wird auf einem Trainingsdatensatz trainiert und anschliessend über eine explorative Datenanalyse validiert. Signifikanzen über β_i können hier nicht festgestellt werden.

22 Aufgabe

(2 Punkte)

In der folgenden Abbildung ist der Zusammenhang vom Modell zu der linearen Regression und der ANOVA skizziert.



Welche der folgenden Aussagen ist richtig?

- A ☐ Die Effektschätzer aus einem Modell, in diesem Fall ein lineares Modell, erlauben es sowohl eine ANOVA zurechnen sowie auch eine Zusammenfassung der Mittelwerte zu betrachten.
- B ☐ Die Effektschätzer aus einem Modell, erlauben es nur eine ANOVA zu rechnen.

- C** ☐ Die Effektschätzer aus einem Modell, in diesem Fall ein lineares Modell, erlauben es nur eine ANOVA zurechnen oder eine Zusammenfassung der Mittelwerte zu betrachten. Beides ist nicht möglich.
- D** ☐ Die Effektschätzer aus einem Modell, in diesem Fall ein polynomes Modell, erlauben es sowohl eine ANOVA zurechnen sowie auch eine Zusammenfassung der Mittelwerte zu betrachten.
- E** ☐ Die Effektschätzer aus einem Modell, erlauben es nur einen Mittelwertsvergleich zu rechnen.

23 Aufgabe

(2 Punkte)

In der Statistik werden die Daten D modelliert in dem ein Modell der Form $Y \sim X$ aufgestellt wird. Welche statistische Kenngrösse wird modelliert?

- A** ☐ Die Y werden modelliert.
- B** ☐ Die Mittelwerte werden modelliert.
- C** ☐ Die X werden modelliert.
- D** ☐ Die Varianz der X unabhängig vom Y wird modelliert.
- E** ☐ Die Varianzstruktur wird modelliert.

24 Aufgabe

(2 Punkte)

Sie führen ein Experiment zur Behandlung von Klaueninfektionen bei Kühen durch. Bei 3 Tieren finden Sie eine Erkrankung der Klauen vor und 12 Tiere sind gesund. Welche Aussage über den Risk ratio Effektschätzer ist richtig?

- A** ☐ Es ergibt sich ein Risk ratio von 0.25, da es sich um ein Anteil handelt.
- B** ☐ Es ergibt sich ein Risk ratio von 0.25, da es sich um eine Chancenverhältnis handelt.
- C** ☐ Es ergibt sich ein Risk ratio von 0.2, da es sich um ein Anteil handelt.
- D** ☐ Es ergibt sich ein Risk ratio von 0.2, da es sich um eine Chancenverhältnis handelt.
- E** ☐ Es ergibt sich ein Risk ratio von 4, da es sich um ein Anteil handelt.

25 Aufgabe

(2 Punkte)

Welche Aussage über die nicht-parametrische Statistik ist richtig?

- A** ☐ Die nicht-parametrische Statistik ist ein Vorgänger der parametrischen Statistik und wurde wegen dem Mangel an Effektschätzern nicht mehr ab 1960 genutzt.
- B** ☐ Die nicht-parametrische Statistik basiert auf dem Schätzen von Parametern aus einer a priori festgelegten Verteilung. Daher gibt es auch direkt zu interpretierenden Effektschätzer.
- C** ☐ Die nicht-parametrische Statistik basiert auf Rängen. Daher wird jeder Zahl ein Rang zugeteilt. Nur auf den Rängen wird die Auswertung gerechnet. Daher gibt es auch keinen direkt zu interpretierenden Effektschätzer.
- D** ☐ Die nicht-parametrische Statistik basiert auf dem Schätzen von Parametern aus einer festgelegten Verteilung. Daher gibt es auch direkt zu interpretierenden Effektschätzer.
- E** ☐ Die nicht-parametrische Statistik basiert auf Rängen. Daher gibt es auch direkt zu interpretierenden Effektschätzer.

26 Aufgabe

(2 Punkte)



Die Randomisierung von Beobachtungen bzw. Samples zu den Versuchseinheiten ist bedeutend in der Versuchsplanung. Welche der folgenden Aussagen ist richtig?

- A ☐ Randomisierung bringt starke Unstrukturiertheit in das Experiment und erlaubt erst von der Stichprobe auf die Grundgesamtheit zurückzuschliessen.
- B ☐ Randomisierung erlaubt erst die Varianzen zu schätzen. Ohne eine Randomisierung ist die Berechnung von Mittelwerten und Varianzen nicht möglich.
- C ☐ Randomisierung erlaubt erst die Mittelwerte zu schätzen. Ohne Randomisierung keine Mittelwerte.
- D ☐ Randomisierung sorgt für Strukturgleichheit und erlaubt erst von der Stichprobe auf die Grundgesamtheit zurückzuschliessen.
- E ☐ Randomisierung war bis 1952 bedeutend, wurde dann aber in Folge besserer Rechnerleistung nicht mehr verwendet. Aktuelle Statistik nutzt keine Randomisierung mehr.

27 Aufgabe

(2 Punkte)

Wenn Sie einen Datensatz erstellen, dann ist es ratsam die Spalten und die Einträge in englischer Sprache zu verfassen, wenn Sie später die Daten in  auswerten wollen. Welcher folgende Grund ist richtig?

- A ☐ Die Spracherkennung von  ist nicht in der Lage Deutsch zu verstehen.
- B ☐ Programmiersprachen können nur englische Begriffe verarbeiten. Zusätzliche Pakete können zwar geladen werden, aber meist funktionieren diese Pakete nicht richtig. Deutsch ist International nicht bedeutend genug.
- C ☐ Alle Funktionen und auch Anwendungen sind in  in englischer Sprache. Die Nutzung von deutschen Wörtern ist nicht schick und das ist zu vermeiden.
- D ☐ Es gibt keinen Grund nicht auch deutsche Wörter zu verwenden. Es ist ein Stilmittel.
- E ☐ Im Allgemeinen haben Programmiersprachen Probleme mit Umlauten und Sonderzeichen, die in der deutschen Sprache vorkommen. Eine Nutzung der englischen Sprache umgeht dieses Problem auf einfache Art.

28 Aufgabe

(2 Punkte)


Bei der explorativen Datenanalyse (EDA) in  gibt es eine richtige Abfolge von Prozessschritten, auch *Circle of life* genannt. Wie lautet die richtige Reihenfolge für die Erstellung einer EDA?

- A ☐ Wir transformieren die Spalten über `mutate()` in ein `tibble` und können dann über `ggplot()` uns die Abbildungen erstellen lassen. Dabei beachten wir das wir keine Faktoren in den Daten haben.
- B ☐ Wir lesen die Daten über eine generische Funktion `read()` ein und müssen dann die Funktion `ggplot()` nur noch installieren. Dann haben wir die Abbildungen als `*.png` vorliegen.
- C ☐ Wir lesen als erstes die Daten über `read_excel()` ein, transformieren die Spalten über `mutate()` in die richtige Form und können dann über `ggplot()` uns die Abbildungen erstellen lassen. Wichtig ist, dass wir keine Faktoren sondern nur numerische Variablen vorliegen haben.

- D** ☐ Wir lesen als erstes die Daten über `read_excel()` ein, transformieren die Spalten über `mutate()` in die richtige Form und können dann über `ggplot()` uns die Abbildungen erstellen lassen.
- E** ☐ Wir lesen die Daten ein und mutieren die Daten. Dabei ist wichtig, dass wir nicht das Paket `tidyverse` nutzen, da dieses Paket veraltet ist. Über die Funktion `library(tidyverse)` entfernen wir das Paket von der Analyse.

29 Aufgabe

(2 Punkte)

In einem Stallexperiment mit $n = 109$ Ferkeln wurden verschiedene Outcomes gemessen: der Gewichtszuwachs, Überleben nach 21 Tagen sowie Anzahl Verletzungen pro 7 Tagen. Zwei Lichtregime wurden als Einflussfaktor gemessen. Sie erhalten den  Output der Funktion `tidy()` einer simplen *logistischen* linearen Regression. Welche Aussage über den **Effekt** ist richtig?

term	estimate	std.error
(Intercept)	-0.11	0.27
light_binhigh	-2.72	0.65

- A** ☐ In einer logistischen Regression kann kein Effekt roh interpretiert werden. Es muss erst eine Confounderadjustierung durchgeführt werden.
- B** ☐ Eine logistischen Regression basiert auf dem maximum Likelihood Prinzip. Hierbei kann kein Effekt beschrieben werden. Im Zweifel hilft aber eine Quadrierung der Fehlerquadrate ϵ .
- C** ☐ In einer logistischen Regression wird die Mittelwertsdifferenz betrachtet. Daher ist der Effekt zwischen den beiden Lichtregimen eine Gewichtsänderung von -2.72
- D** ☐ In einer logistischen Regression berechnet man das RR. Daher muss der Schätzer des Effektes β_1 noch quadriert werden. Somit liegt das RR bei 7.42
- E** ☐ In einer logistischen Regression muss für die Interpretation des Effektes das β_1 quadriert werden. Somit liegt das OR bei 7.42

30 Aufgabe

(2 Punkte)

In einer linearen Regression werden die ϵ oder Residuen geschätzt. Welcher Verteilung folgen die Residuen bei einer optimalen Modellierung?

- A** ☐ Die Residuen folgen einer Poissonverteilung mit $\text{Pois}(0)$.
- B** ☐ Die Residuen sind normalverteilt mit $\mathcal{N}(0, 1)$.
- C** ☐ Die Residuen sind normalverteilt mit $\mathcal{N}(0, s^2)$.
- D** ☐ Die Residuen sind normalverteilt mit $\mathcal{N}(\bar{y}, s^2)$.
- E** ☐ Die Residuen sind binomialverteilt.

31 Aufgabe

(2 Punkte)

Welche Aussage über das *generalisierte lineare Modell (GLM)* ist richtig?

- A** ☐ Das GLM ist eine Vereinfachung des LM in R. Mit dem GLM lassen polygonale Regressionen rechnen.

- B** ☐ Das GLM ist ein faktisch maschineller Lernalgorithmus, der selbstständig die Verteilungsfamilie für Y wählt.
- C** ☐ Das GLM ist eine allgemeine Erweiterung der linearen Regression auf die Normalverteilung.
- D** ☐ Das GLM erlaubt auch weitere Verteilungsfamilien für das Y bzw. das Outcome in einer linearen Regression zu wählen.
- E** ☐ Das GLM erlaubt auch nicht normalverteilte Residuen in der Schätzung der Regressionsgrade.

32 Aufgabe

(2 Punkte)

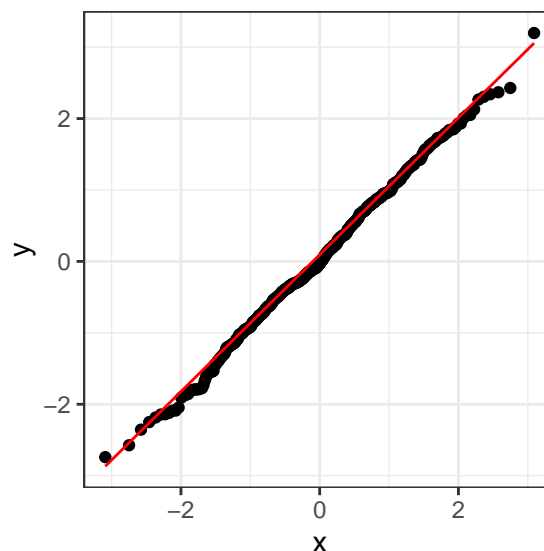
Sie rechnen in einer linearen Regression das Modell A und das Modell B. Nun stellt sich die Frage, welches der beiden Modelle das bessere Modell ist. Um die Modelle bewerten zu können berechnen Sie dafür das AIC_A für Modell A mit 125 und für das Modell B das AIC_B von 435. Welche Aussage über die beiden Modelle ist richtig?

- A** ☐ Da $AIC_A > 0$ ist das Modell A das bessere Modell. Der AIC Wert für B wird verworfen.
- B** ☐ Da $AIC_A < AIC_B$ ist das Modell B das bessere Modell.
- C** ☐ Da $AIC_A > AIC_B$ ist das Modell B das bessere Modell.
- D** ☐ Da $AIC_A > AIC_B$ ist das Modell A das bessere Modell.
- E** ☐ Da $AIC_A < AIC_B$ ist das Modell A das bessere Modell.

33 Aufgabe

(2 Punkte)

Sie rechnen in eine linearen Regression und erhalten folgenden QQ Plot. Welche Aussage ist richtig?



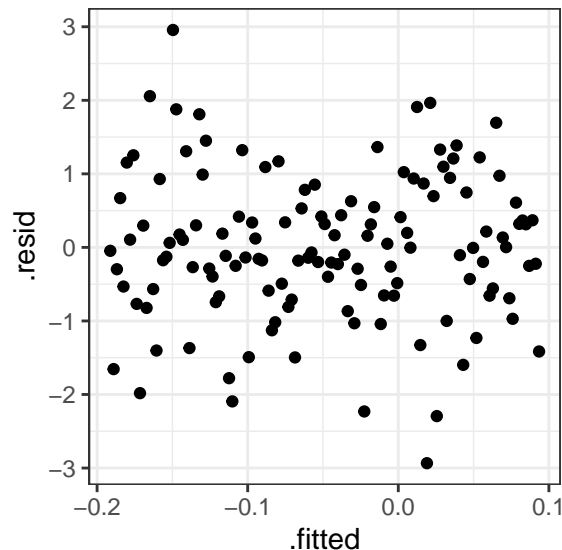
- A** ☐ Die Annahme der normalverteilten Residuen ist nicht erfüllt. Die Punkte liegen zum überwiegenden Teil nicht auf der Geraden.
- B** ☐ Die Annahme der normalverteilten Residuen ist nicht erfüllt. Die Punkte liegen zum überwiegenden Teil auf der Geraden.
- C** ☐ Die Annahme der normalverteilten Residuen ist erfüllt. Die Punkte liegen zum überwiegenden Teil nicht auf der Geraden und Korrelation ist negativ.

- D** ☐ Die Annahme der normalverteilten Residuen ist erfüllt. Die Punkte liegen zum überwiegenden Teil nicht auf der Geraden.
- E** ☐ Die Annahme der normalverteilten Residuen ist erfüllt. Die Punkte liegen zum überwiegenden Teil auf der Geraden.

34 Aufgabe

(2 Punkte)

Sie rechnen eine lineare Regression und erhalten folgenden Residual Plot. Welche Aussage ist richtig?



- A** ☐ Die Annahme der normalverteilten Residuen ist nicht erfüllt. Vereinzelt Punkte liegen oberhalb bzw. unterhalb der Geraden um die 0 Linie weiter entfernt. Ein klares Muster ist zu erkennen.
- B** ☐ Die Annahme der normalverteilten Residuen ist nicht erfüllt. Es ist kein Muster zu erkennen.
- C** ☐ Die Annahme der normalverteilten Residuen ist erfüllt. Kein Muster ist zu erkennen und keine Outlier zu beobachten.
- D** ☐ Die Annahme der normalverteilten Residuen ist erfüllt. Die Punkte liegen zum überwiegenden Teil auf der Diagonalen.
- E** ☐ Die Annahme der normalverteilten Residuen ist erfüllt. Es ist ein Muster zu erkennen.

35 Aufgabe

(2 Punkte)

Sie wollen ein Feldexperiment mit zwei Düngestufen durchführen. Berechnen Sie die benötigte Fallzahl mit $t_{1-\alpha/2} = 1.645$ und $t_{1-\beta} = 1.282$ sowie $s = 1$ und $\delta_0 = 1$. Es ergibt sich somit folgende Fallzahl.

- A** ☐ Es ergibt sich eine Fallzahl von 9
- B** ☐ Es ergibt sich eine Fallzahl von 36
- C** ☐ Es ergibt sich eine Fallzahl von 18
- D** ☐ Es ergibt sich eine Fallzahl von 34
- E** ☐ Es ergibt sich eine Fallzahl von 17

36 Aufgabe

(2 Punkte)

Welche Aussage zum mathematische Ausdruck $Pr(D|H_0)$ ist richtig?

- A ☐ Die Inverse der Wahrscheinlichkeit unter der die Nullhypothese nicht mehr die Alternativhypothese überdeckt.
- B ☐ Die Wahrscheinlichkeit für die Nullhypothese, wenn die Daten wahr sind.
- C ☐ Die Wahrscheinlichkeit der Daten unter der Nullhypothese in der Grundgesamtheit.
- D ☐ $Pr(D|H_0)$ ist die Wahrscheinlichkeit der Alternativhypothese und somit $1 - Pr(H_A)$
- E ☐ $Pr(D|H_0)$ ist die Wahrscheinlichkeit die Daten D zu beobachten wenn die Nullhypothese wahr ist.

37 Aufgabe

(2 Punkte)

Das Falsifikationsprinzip besagt...

- A ☐ ... dass ein schlechtes Modell durch ein weniger schlechtes Modell ersetzt wird. Die Wissenschaft lehnt ab und verifiziert nicht.
- B ☐ ... dass Annahmen an statistische Modelle meist falsch sind.
- C ☐ ... dass Fehlerterme in statistischen Modellen nicht verifiziert werden können.
- D ☐ ... dass Modelle meist falsch sind und selten richtig.
- E ☐ ... dass in der Wissenschaft immer etwas falsch sein muss. Sonst gebe es keinen Fortschritt.

38 Aufgabe

(2 Punkte)

Der Fehler 1. Art oder auch Signifikanzniveau α genannt, liegt bei 5%. Welcher der folgenden Gründe für diese Festlegung auf 5% ist richtig?

- A ☐ Die Festlegung von $\alpha = 5\%$ ist eine Kulturkonstante. Wissenschaftler benötigt eine Schwelle für eine statistische Testentscheidung, der Wert von α wurde aber historisch mehr zufällig gewählt.
- B ☐ Auf einer Statistikkonferenz in Genf im Jahre 1942 wurde dieser Cut-Off nach langen Diskussionen festgelegt. Bis heute ist der Cut Off aber umstritten, da wegen dem 2. Weltkrieg viele Wissenschaftler nicht teilnehmen konnten.
- C ☐ Im Rahmen eines langen Disputs zwischen Neyman und Fischer wurde $\alpha = 5\%$ festgelegt. Leider werden die Randbedingungen und Voraussetzungen an statistische Modelle heute immer wieder ignoriert.
- D ☐ Der Wert ergab sich aus einer Auswertung von 1042 wissenschaftlichen Veröffentlichungen zwischen 1914 und 1948. Der Wert 5% wurde in 28% der Veröffentlichungen genutzt. Daher legte man sich auf diese Zahl fest.
- E ☐ Der Begründer der modernen Statistik, R. Fischer, hat die Grenze simuliert und berechnet. Dadurch ergibt sich dieser optimale Cut-Off.

39 Aufgabe

(2 Punkte)

Welche Aussage über die Power ist richtig?

- A** ☐ Die Power beschreibt die Wahrscheinlichkeit die H_A abzulehnen. Wir testen die Power jedoch nicht.
- B** ☐ Die Power $1 - \beta$ wird auf 80% gesetzt. Damit liegt die Wahrscheinlichkeit für die H_0 bei 20%.
- C** ☐ Die Power ist nicht in der aktuellen Testtheorie mehr vertreten. Wir rechnen nur noch mit dem Fehler 1. Art.
- D** ☐ Es gilt $\alpha + \beta = 1$ und somit liegt β meist bei 95%.
- E** ☐ Die Power $1 - \beta$ wird auf 80% gesetzt. Alle statistischen Tests sind so konstruiert, dass die H_A mit 80% "bewiesen wird".

40 Aufgabe

(2 Punkte)

Beim statistischen Testen wird `signal` mit `noise` zur Teststatistik `T` verrechnet. Welche der Formel berechnet korrekt die Teststatistik `T`?

- A** ☐ Es gilt $T = (\text{signal} \cdot \text{noise})^2$
- B** ☐ Es gilt $T = \text{signal} \cdot \text{noise}$
- C** ☐ Es gilt $T = \frac{\text{noise}}{\text{signal}}$
- D** ☐ Es gilt $T = \frac{\text{signal}}{\text{noise}^2}$
- E** ☐ Es gilt $T = \frac{\text{signal}}{\text{noise}}$

41 Aufgabe

(2 Punkte)

In der Theorie zur statistischen Testentscheidung kann „ H_0 ablehnen obwohl die H_0 gilt“ in welche richtige Analogie gesetzt werden?

- A** ☐ In die Analogie eines Rauchmelders: *Alarm with fire*.
- B** ☐ In die Analogie eines Rauchmelders: *Alarm without fire*, dem α -Fehler.
- C** ☐ In die Analogie eines brennenden Hauses ohne Rauchmelder: *House without noise*.
- D** ☐ In die Analogie eines Feuerwehrautos: *Car without noise*.
- E** ☐ In die Analogie eines Rauchmelders: *Fire without alarm*, dem β -Fehler.

42 Aufgabe

(2 Punkte)

Sie rechnen eine simple Gaussian Regression. Welche Aussage betreffend der Konfidenzintervalle ist für die Gaussian Regression richtig?

- A ☐ Wenn die Relevanzschwelle mit enthalten ist, kann die Nullhypothese abgelehnt werden.
- B ☐ Wenn die Konfidenzintervalle den p-Wert der Regression enthalten, kann die Nullhypothese abgelehnt werden.
- C ☐ Wenn die 0 im Konfidenzintervall enthalten ist, kann die Nullhypothese abgelehnt werden.
- D ☐ Wenn die 0 im Konfidenzintervall enthalten ist, kann die Nullhypothese nicht abgelehnt werden.
- E ☐ Wenn die 1 im Konfidenzintervall enthalten ist, kann die Nullhypothese nicht abgelehnt werden.

43 Aufgabe

(2 Punkte)

In der Bio Data Science wird häufig mit sehr großen Datensätzen gerechnet. Historisch ergibt sich nun ein Problem bei der Auswertung der Daten und deren Bewertung hinsichtlich der Signifikanz. Welche Aussage ist richtig?

- A ☐ Aktuell werden immer grössere Datensätze erhoben. Eine erhöhte Fallzahl führt automatisch auch zu mehr signifikanten Ergebnissen, selbst wenn die eigentlichen Effekte nicht relevant sind.
- B ☐ Big Data ist ein Problem der parametrischen Statistik. Parameter lassen sich nur auf kleinen Datensätzen berechnen, da es sich sonst nicht mehr um eine Stichprobe im engen Sinne der Statistik handelt.
- C ☐ Relevanz und Signifikanz haben nichts miteinander zu tun. Daher gibt es auch keinen Zusammenhang zwischen hoher Fallzahl ($n > 10000$) und einem signifikanten Test. Ein Effekt ist immer relevant und somit signifikant.
- D ☐ Aktuell werden immer grössere Datensätze erhoben. Dadurch wird auch die Varianz immer höher was automatisch zu mehr signifikanten Ergebnissen führt.
- E ☐ Aktuell werden zu grosse Datensätze für die gängige Statistik gemessen. Daher wendet man maschinelle Lernverfahren für kausale Modelle an. Hier ist die Relevanz gleich Signifikanz.

44 Aufgabe

(2 Punkte)

Welche statistische Masszahl erlaubt es *Relevanz* mit *Signifikanz* zuverbinden? Welche Aussage ist richtig?

- A ☐ Das OR. Als Chancenverhältnis gibt es das Verhältnis von Relevanz und Signifikanz wieder.
- B ☐ Das Konfidenzintervall. Durch die Visualisierung des Konfidenzintervalls kann eine Relevanzschwelle vom Anwender definiert werden. Zusätzlich erlaubt das Konfidenzintervall auch eine Entscheidung über die Signifikanz.
- C ☐ Die Teststatistik. Durch den Vergleich von T_c zu T_k ist es möglich die H_0 abzulehnen. Die Relevanz ergibt sich aus der Fläche rechts vom dem T_c -Wert.
- D ☐ Der p-Wert. Durch den Vergleich mit α lässt sich über die Signifikanz entscheiden und der β -Fehler erlaubt über die Power eine Einschätzung der Relevanz.
- E ☐ Das Δ . Durch die Effektstärke haben wir einen Wert für die Relevanz, die vom Anwender bewertet werden muss. Da Δ antiproportional zum p-Wert ist, bedeutet auch ein hohes Δ ein sehr kleinen p-Wert.

45 Aufgabe

(2 Punkte)

Welche Aussage über die frequentistischen Testtheorie ist richtig?

- A ☐ Wir machen keine Aussagen über Wahrscheinlichkeiten!
- B ☐ Wir machen Aussagen über Wahrscheinlichkeiten!
- C ☐ Wir machen Aussagen über die individuelle Wahrscheinlichkeit eines Effektes!
- D ☐ Wir machen Aussagen über Individuen!
- E ☐ Wir machen Aussagen über den Effekt!

46 Aufgabe

(2 Punkte)

Welche Aussage über den p -Wert und dem Signifikanzniveau α gleich 5% ist richtig?

- A ☐ Wir vergleichen mit dem p -Wert und dem Signifikanzniveau α Wahrscheinlichkeiten und damit die absoluten Werte auf einem Zahlenstrahl, wenn die H_0 gilt.
- B ☐ Wir machen eine Aussage über die individuelle Wahrscheinlichkeit des Eintretens der Nullhypothese H_0 .
- C ☐ Wir vergleichen mit dem p -Wert und dem Signifikanzniveau α Wahrscheinlichkeiten und damit die Flächen unter der Kurve der Teststatistik, wenn die H_0 gilt.
- D ☐ Wir vergleichen die Effekte des p -Wertes mit den Effekten der Signifikanzschwelle unter der Annahme der Nullhypothese.
- E ☐ Wir vergleichen mit dem p -Wert und dem Signifikanzniveau α absolute Werte auf einem Zahlenstrahl und damit den Unterschied der Teststatistiken, wenn die H_0 gilt.

47 Aufgabe

(2 Punkte)

Welche Aussage über den t-Test ist richtig?

- A ☐ Der t-Test ist ein Vortest der ANOVA und basiert daher auf dem Vergleich von Streuungsparametern
- B ☐ Der t-Test testet generell zu einem erhöhten α -Niveau von 20%.
- C ☐ Der t-Test vergleicht die Varianzen von mindestens zwei oder mehr Gruppen
- D ☐ Der t-Test vergleicht die Mittelwerte von zwei Gruppen unter der strikten Annahme von Varianzhomogenität. Sollte keine Varianzhomogenität vorliegen, so gibt es keine Möglichkeit den t-Test in einer Variante anzuwenden.
- E ☐ Der t-Test vergleicht die Mittelwerte von zwei Gruppen.

48 Aufgabe

(2 Punkte)

Welche Aussage über den Welch t-Test ist richtig?

- A ☐ Der Welch t-Test vergleicht die Varianz von zwei Gruppen.
- B ☐ Der Welch t-Test ist die veraltete Form des Student t-Test und wird somit nicht mehr verwendet.

- C** ☐ Der Welch t-Test ist ein Post-hoc Test der ANOVA und basiert daher auf dem Vergleich der Varianz.
- D** ☐ Der Welch t-Test wird angewendet, wenn Varianzheterogenität zwischen den beiden zu vergleichenden Gruppen vorliegt.
- E** ☐ Der Welch t-Test vergleicht die Mittelwerte von zwei Gruppen unter der strikten Annahme von Varianzhomogenität.

49 Aufgabe

(2 Punkte)

Nach einem Experiment mit fünf Weizensorten ergibt eine ANOVA ($p = 0.041$) einen signifikanten Unterschied für den Ertrag. Sie führen anschließend die paarweisen t-Tests für alle Vergleiche der verschiedenen Weizensorten durch. Nach der Adjustierung für multiples Testen ist kein p-Wert unter der α -Schwelle. Sie schauen sich auch die rohen, unadjustierten p-Werte an und finden hier als niedrigsten p-Wert $p_{3-2} = 0.053$. Welche Aussage ist richtig?

- A** ☐ Die adjustierten p-Werte deuten in die richtige Richtung. Zusammen mit den nicht signifikanten rohen p-Werten ist von einem Fehler in der ANOVA auszugehen.
- B** ☐ Es gibt einen Fehler in der Varianzstruktur. Daher kann die ANOVA nicht richtig sein und paarweise t-Tests liefern das richtige Ergebnis.
- C** ☐ Die ANOVA testet auf der gesamten Fallzahl. Es wäre besser die ANOVA auf der gleichen Fallzahl wie die einzelnen t-Tests zu rechnen.
- D** ☐ Die ANOVA testet auf der gesamten Fallzahl. Die einzelnen t-Tests immer nur auf einer kleineren Subgruppe. Da mit weniger Fallzahl weniger signifikante Ergebnisse zu erwarten sind, kann eine Diskrepanz zwischen der ANOVA und den paarweisen t-Tests auftreten.
- E** ☐ Der Fehler liegt in den t-Tests. Wenn eine ANOVA signifikant ist, dann muss zwangsweise auch ein t-Test signifikant sein.

50 Aufgabe

(2 Punkte)

Welche Aussage über den gepaarten t-Test für verbundene Stichproben ist richtig?

- A** ☐ Beim gepaarten t-Test kombinieren wir die Vorteile des Student t-Test für Varianzhomogenität mit den Vorteilen des Welch t-Test für Varianzheterogenität. Wir bilden dafür die Differenz der Einzelbeobachtungen.
- B** ☐ Der gepaarte t-Test wird genutzt, wenn die Differenzen der Beobachtungen verbunden sind und wir dadurch die Unabhängigkeit nicht mehr vorliegen haben.
- C** ☐ Der gepaarte t-Test wird gerechnet, wenn die Beobachtungen abhängig voneinander sind. Wir messen jede Beobachtung nur einmal und berechnen dann die Differenz zu dem Mittel der anderen Beobachtungen.
- D** ☐ Der gepaarte t-Test wird gerechnet, wenn die Beobachtungen nicht unabhängig voneinander sind. Wir messen wiederholt an dem gleichen Probanden oder Tier oder Pflanze. Wir bilden die Differenzen um den gepaarten t-Test rechnen zu können.
- E** ☐ Der gepaarte t-Test nutzt die Varianz der Beobachtungen jeweils paarweise und bildet dafür eine verbundene Stichprobe. Dieser Datensatz d dient dann zur Differenzbildung.

Deskriptive Statistik & Explorative Datenanalyse

Mehr Informationen zu den Aufgaben in den folgenden Kapiteln aus dem Skript Bio Data Science.

- Kapitel 15 - Deskriptive Statistik
- Kapitel 16 - Visualisierung von Daten
- Kapitel 18 - Verteilung von Daten

51 Aufgabe

(10 Punkte)

Sie haben folgende Zahlenreihe y vorliegen $y = \{18, 21, 18, 20, 17, 19, 21\}$. Berechnen Sie folgende deskriptive Maßzahlen. Geben Sie Formeln und Rechenwege mit an!

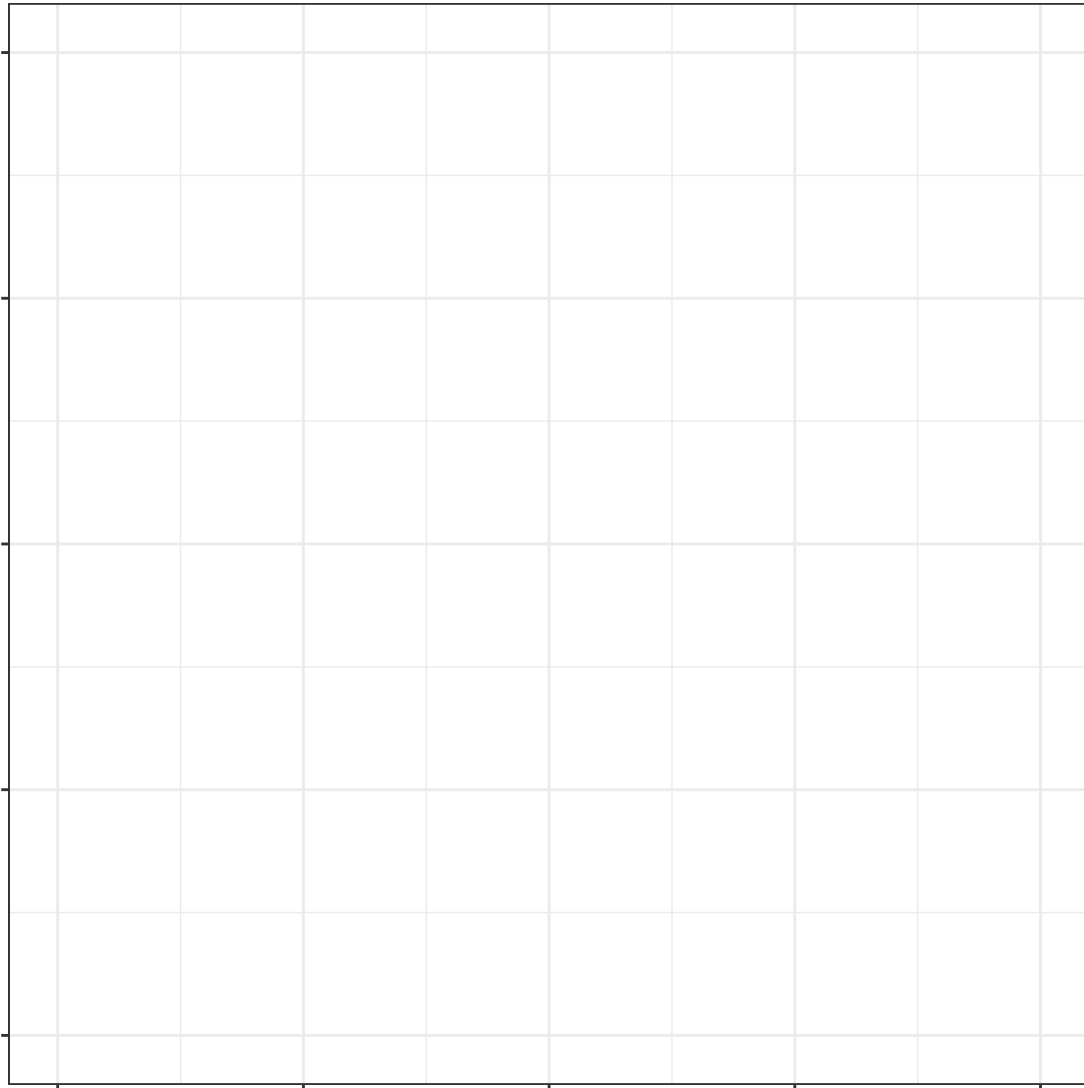
1. Das 3rd Quartile **(2 Punkte)**
2. Die Range oder Spannweite **(2 Punkte)**
3. Den Mittelwert **(2 Punkte)**
4. Die Standardabweichung **(2 Punkte)**
5. Die Varianz **(2 Punkte)**

52 Aufgabe

(8 Punkte)

Sie haben folgende Zahlenreihe y vorliegen $y = \{27, 25, 21, 23, 20\}$.

1. Visualisieren Sie den Mittelwert von y in der untenstehenden Abbildung! **(4 Punkte)**
2. Beschriften Sie die Y und X -Achse entsprechend! **(2 Punkte)**
3. Für die Berechnung der Varianz wird der Abstand der einzelnen Werte x_i zum Mittelwert \bar{x} quadriert. Warum muss der Abstand, $x_i - \bar{x}$, in der Varianzformel quadriert werden? Erklären Sie den Zusammenhang unter Berücksichtigung der Abbildung! **(2 Punkte)**



53 Aufgabe

(7 Punkte)

Nach einem Gewächshausexperiment mit drei Bewässerungstypen (*low*, *mid* und *high*) ergibt sich die folgende Datentabelle mit dem gemessenen Frischgewicht (*freshmatter*).

water_type	freshmatter
low	22
low	21
high	18
high	20
low	22
high	23
low	21
mid	24
mid	30
high	18
mid	25
low	19

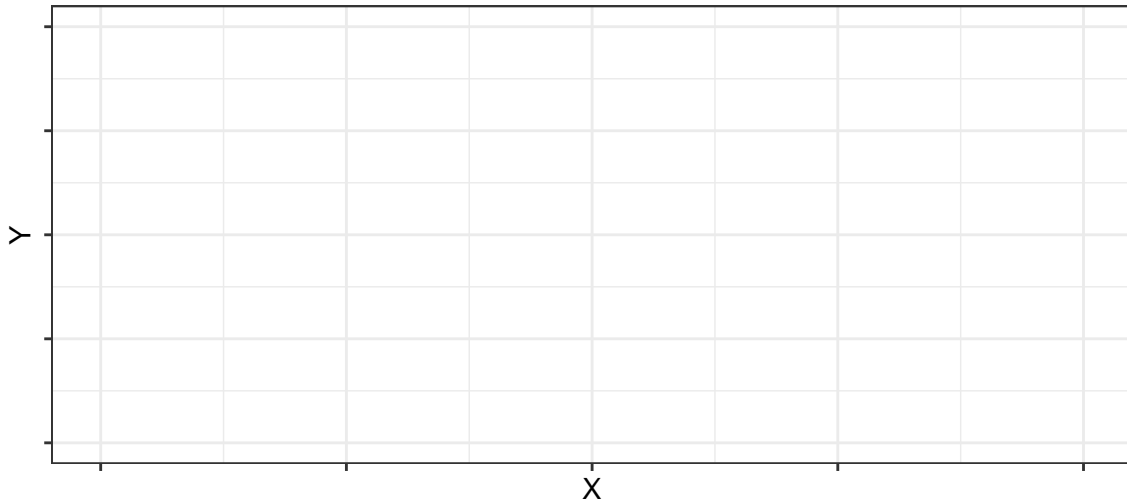
1. Zeichnen Sie in *einer* Abbildung die Barplots für die Bewässerungstypen! Beschriften Sie die Achsen entsprechend! **(4 Punkte)**
2. Beschriften Sie *einen* Barplot mit den gängigen statistischen Maßzahlen! **(2 Punkte)**
3. Wenn Sie *keinen Effekt* zwischen der Bewässerungstypen erwarten würden, wie sehen dann die Barplots aus? **(1 Punkt)**

54 Aufgabe

(6 Punkte)

1. Skizzieren Sie in die unten stehenden, freien Abbildungen die Verteilungen, die sich nach der Abbildungsüberschrift ergeben! **(4 Punkte)**
2. Achten Sie auf die entsprechende Skalierung der beiden Verteilungen in der ersten Abbildung! **(2 Punkte)**

$N(4, 1)$ und $N(1, 3)$



$\text{Pois}(5)$



55 Aufgabe

(6 Punkte)

1. Skizzieren Sie 2 Normalverteilungen *in einer Abbildung* mit $\mu_1 \neq \mu_2$ und $s_1 = s_2$! **(2 Punkte)**
2. Beschriften Sie die Normalverteilungen mit den entsprechenden Parametern! **(2 Punkte)**
3. Liegt Varianzhomogenität oder Varianzheterogenität vor? Begründen Sie Ihre Antwort! **(2 Punkte)**

56 Aufgabe

(6 Punkte)

Nach einem Experiment zählen Sie folgende Anzahl an Läsionen auf den Blättern von Sonnenblumen nach einer durchgestandenen Infektion.

5, 10, 9, 3, 7, 3, 8, 4, 3, 4, 1, 7, 7, 8

1. Zeichnen Sie ein Histogramm um die Verteilung der Daten zu visualisieren. **(3 Punkte)**
2. Beschriften Sie die Achsen der Abbildung! **(2 Punkte)**
3. Ergänzen Sie die relativen Häufigkeiten in der Abbildung! **1 Punkt**

57 Aufgabe

(9 Punkte)

Nach einem Feldexperiment mit zwei Düngestufen (A und B) ergibt sich die folgende Datentabelle mit dem gemessenen Trockengewicht (*drymatter*).

trt	drymatter
A	18.3
A	19.3
B	11.7
B	15.8
A	12.4
A	14.8
B	12.5
B	17.0
A	12.1
A	18.1
A	15.2
B	11.6
B	9.4
B	19.2

1. Zeichnen Sie in *einer* Abbildung die beiden Boxplots für die zwei Düngestufen A und B! Beschriften Sie die Achsen entsprechend! **(6 Punkte)**
2. Beschriften Sie *einen* der beiden Boxplots mit den gängigen statistischen Maßzahlen! **(2 Punkte)**
3. Wenn Sie *keinen Effekt* zwischen den Düngestufen erwarten würden, wie sehen dann die beiden Boxplots aus? **(1 Punkt)**

58 Aufgabe

(9 Punkte)

Nach einem Feldexperiment mit mehreren Düngestufen stellt sich die Frage, ob die Düngestufe *low* im Bezug auf das Trockengewicht normalverteilt sei. Sie erhalten folgende Datentabelle.

fertilizer	drymatter
low	26
low	27
low	29
low	20
low	26
low	22
low	27
low	29
low	19
low	26
low	28

1. Zeichnen Sie eine passende Abbildung in der Sie visuell überprüfen können, ob eine Normalverteilung des Trockengewichts vorliegt! **(4 Punkte)**
2. Beschriften Sie die Achsen und ergänzen Sie die statistischen Maßzahlen. **(3 Punkte)**
3. Entscheiden Sie, ob eine Normalverteilung vorliegt. Begründen Sie Ihre Antwort. **(2 Punkte)**

59 Aufgabe

(4 Punkte)

1. Zeichnen Sie über den untenstehenden Boxplot die entsprechende zugehörige Verteilung! (2 Punkte)
2. Zeichnen Sie unter den untenstehenden Boxplot die entsprechende zugehörige Beobachtungen! (2 Punkte)



60 Aufgabe

(6 Punkte)

Nach einer Bonitur von Schnittlauch mit einer Kontrolle und drei Pestiziden (ctrl, pestKill, roundUp, zeroX) ergibt sich die folgende Datentabelle mit den Boniturnoten (*grade*).

pesticide	grade
zeroX	2
ctrl	4
pestKill	3
roundUp	2
roundUp	3
zeroX	0
ctrl	8
pestKill	4
zeroX	1
zeroX	1
pestKill	4
roundUp	2
ctrl	7
roundUp	4
pestKill	3

1. Zeichnen Sie in *einer* Abbildung die Dotplots für die vier Pestizidlevel! Beschriften Sie die Achsen entsprechend! **(3 Punkte)**
2. Ergänzen Sie die Dotplots mit den gängigen statistischen Maßzahlen. **(2 Punkte)**
3. Wenn Sie *keinen Effekt* zwischen den Pestizidlevel erwarten würden, wie sehen dann die Dotplots aus? **(1 Punkt)**

61 Aufgabe

(6 Punkte)

Nach einem Feldexperiment mit zwei Pestiziden (*RoundUp* und *OutEx*) ergibt sich die folgende Datentabelle mit dem jeweiligen beobachteten Infektionsstatus.

pesticide	infected
RoundUp	no
OutEx	yes
OutEx	no
RoundUp	no
OutEx	yes
OutEx	yes
OutEx	yes
RoundUp	yes
RoundUp	yes
RoundUp	no
RoundUp	yes
OutEx	no
OutEx	yes
OutEx	yes
RoundUp	no
OutEx	yes
RoundUp	no
OutEx	no
RoundUp	no
OutEx	no

1. Stellen Sie in einer 2x2 Tafel den Zusammenhang zwischen dem Pestizid und dem Infektionsstatus dar! **(2 Punkte)**
2. Zeichnen Sie den zugehörigen Mosaic-Plot. Berechnen Sie das Verhältnis pro Spalte! **(2 Punkte)**
3. Wenn das Pestizid keine Auswirkung auf den Infektionsstatus hätte, wie sehe dann der Mosaic-Plot aus? **(2 Punkte)**

62 Aufgabe

(10 Punkte)

In einem Feldexperiment für die Bodendurchlässigkeit wurde der Niederschlag pro Parzelle sowie der durchschnittliche Ertrag gemessen. Es ergibt sich folgende Datentabelle.

water	drymatter
19	16
20	20
18	21
25	22
21	20

1. Erstellen Sie den Scatter-Plot für die Datentabelle. Beschriften Sie die Achsen entsprechend! **(4 Punkte)**
2. Zeichnen Sie eine Gerade durch die Punkte! **(1 Punkt)**
3. Beschriften Sie die Gerade mit den gängigen statistischen Maßzahlen! **(3 Punkte)**
4. Wenn kein Effekt von dem Niederschlag auf das Trockengewicht vorhanden wäre, wie würde die Gerade verlaufen und welche Werte würden die statistischen Maßzahlen annehmen? **(2 Punkt)**

Statistisches Testen

Mehr Informationen zu den Aufgaben in den folgenden Kapiteln aus dem Skript Bio Data Science.

- Kapitel 3 - Falsifikationsprinzip
- Kapitel 19 - Die Testentscheidung
- Kapitel 20 - Die Testtheorie

63 Aufgabe

(6 Punkte)

1. Erklären Sie den Zusammenhang zwischen Stichprobe und Grundgesamtheit an einem Schaubild! **(3 Punkte)**
2. Was ist der Unterschied zwischen μ und σ und \bar{x} und s im Kontext der Stichprobe und Grundgesamtheit? **(2 Punkte)**
3. Warum müssen wir überhaupt zwischen einer Stichprobe und einer Grundgesamtheit unterscheiden? **(1 Punkt)**

64 Aufgabe

(6 Punkte)

Geben ist folgende 2x2 Kreuztabelle.

1. Tragen Sie folgende Fachbegriffe korrekt in die 2x2 Kreuztabelle ein! (6 Punkte)

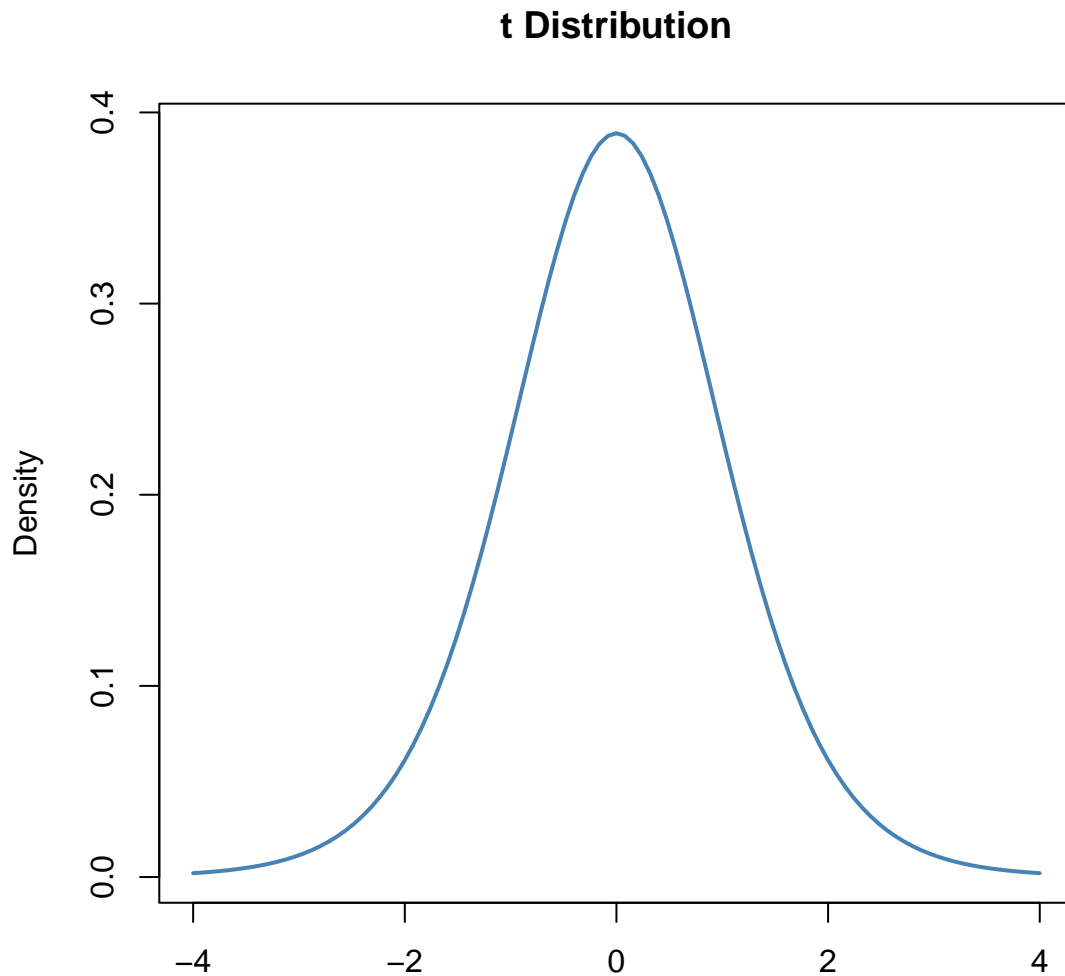
- (Unbekannte) Wahrheit
- H_0 wahr
- H_0 falsch
- H_0 abgelehnt
- H_0 beibehalten
- Testentscheidung
- α -Fehler
- β -Fehler
- Richtige Entscheidung
- 5%
- 20%

65 Aufgabe

(4 Punkte)

Im folgenden ist eine t-Verteilung abgebildet. Ergänzen Sie die Abbildung wie folgt.

1. Zeichnen Sie das Signifikanzniveau α in die Abbildung! **(1 Punkte)**
2. Zeichnen Sie einen signifikant p-Wert in die Abbildung! **(1 Punkt)**
3. Ergänzen Sie „ $\bar{y}_1 = \bar{y}_2$ “! **(1 Punkt)**
4. Ergänzen Sie „ $A = 0.95$ “! **(1 Punkt)**



66 Aufgabe

(8 Punkte)

Sie rechnen einen t-Test. Sie schätzen einen Mittelwertsunterschied.

1. Beschriften Sie die untenstehende Abbildung mit der Signifikanzschwelle! **(1 Punkt)**
2. Ergänzen Sie eine Relevanzschwelle! **(1 Punkt)**
3. Skizzieren Sie in die untenstehende Abbildung fünf einzelne Konfidenzintervalle (a-e) mit den jeweiligen Eigenschaften! **(6 Punkte)**
 - (a) Ein signifikantes, nicht relevantes Konfidenzintervall
 - (b) Ein Konfidenzintervall mit niedriger Varianz s_p in der Stichprobe als der Rest der Konfidenzintervalle
 - (c) Ein nicht signifikantes, relevantes Konfidenzintervall
 - (d) Ein signifikantes, relevantes Konfidenzintervall
 - (e) Ein Konfidenzintervall mit höherer Varianz s_p in der Stichprobe als der Rest der Konfidenzintervalle



67 Aufgabe

(6 Punkte)

Gegeben ist die vereinfachte Formel für den Zweistichproben t-Test mit der gepoolten Standardabweichung s_p und gleicher Gruppengrösse n_g der beiden Sample.

$$T = \frac{\bar{y}_1 - \bar{y}_2}{s_p \cdot \sqrt{\frac{2}{n_g}}}$$

Welche Auswirkung hat die Änderungen der jeweiligen statistischen Masszahl auf den T-Wert und damit auf die *vermutliche* Signifikanz? Füllen Sie hierzu die untenstehende Tabelle aus! **(6 Punkte)**

	T Statistik	$Pr(D H_0)$	$KI_{1-\alpha}$		T Statistik	$Pr(D H_0)$	$KI_{1-\alpha}$
$\Delta \uparrow$				$\Delta \downarrow$			
$s \uparrow$				$s \downarrow$			
$n \uparrow$				$n \downarrow$			

68 Aufgabe

(8 Punkte)

Sie haben folgende Aussage gegeben.

Bin ich im Winter?

1. Erklären Sie den Gedankengang der Testtheorie sowie des Falsifikationsprinzips an der Aussage! **(4 Punkte)**
2. Erklären Sie Ihre Entscheidung zu der Aussage! **(3 Punkte)**
3. Schätzen Sie den p-Wert zu der Aussage ab! **(1 Punkt)**

69 Aufgabe

(8 Punkte)

»Ohne ausreichende Zufuhr von Vitamin C stellen sich nach 60 Tagen die ersten Symptome ein; die ersten Toten sind nach 101 Tagen zu beklagen; nach 116 Tagen rafft die Skorbut eine ganze Schiffsbesatzung dahin.«

1. Stellen Sie den Verlauf der Skorbuterkrankung auf einem Schiff mit einer typischen Anzahl an Matrosen dar. Welche Annahmen haben Sie getroffen? Beschriften Sie die Achsen entsprechend! **(4 Punkte)**
2. Schifffarzt James Lind (1716-1794) ist dem Phänomen Skorbut systematisch nachgegangen. Wie würden Sie solch ein Experiment planen und welche (ethischen) Probleme sehen Sie? **(4 Punkte)**

70 Aufgabe

(11 Punkte)

Der Wissenschaftler Pettenkofer arbeitete streng naturwissenschaftlich-experimentell und gilt als Begründer der experimentellen Hygiene ("Konditionalhygiene"). Auch seine Untersuchungen zu Kleidung, Heizung, Lüftung, Kanalisation und Wasserversorgung trugen experimentelle Züge. Pettenkofer war ein Positivist, das heißt, er erkannte ausschließlich sichtbare, zum Beispiel in Experimenten gewonnene Tatsachen als Erkenntnisquelle an.

Pettenkofer unterlief ein Irrtum, der bis heute nachwirkt, indem viele Menschen glauben, es gebe eine "Atmende Wand": Er stellte bei frühen Luftwechsel-Messungen in einem Zimmer fest, dass sich nach dem vermeintlichen Abdichten sämtlicher Fugen die Luftwechselrate weniger als erwartet verminderte. Daraus schlussfolgerte er einen erheblichen Luftaustausch durch die Ziegelwände hindurch. Vermutlich kam er nicht darauf, den Kamin eines im Raum befindlichen Ofens abzudichten. Luftaustausch durch die Zimmerwände hindurch sei, so Pettenkofer, ein wesentlicher Beitrag zur Reinigung der Raumluft.

1. Bestimmen Sie das y bzw. den Endpunkt/Outcome für das Experiment "Atmende Wand". Wie messen Sie den Endpunkt? **(1 Punkt)**
2. Bestimmen Sie mögliche Einflussfaktoren auf den Endpunkt! Begründen Sie Ihre Entscheidung! **(2 Punkt)**
3. Beschreiben Sie das Modell in der Form $y \sim x$! **(1 Punkt)**
4. Wie realisieren Sie Wiederholung in dem Experiment? Skizzieren Sie das Experiment! **(2 Punkte)**
5. Wie könnten Sie das Experiment graphisch darstellen? Was befindet sich auf der x-Achse und was auf der y-Achse? Beschriften Sie die Achsen entsprechend! **(2 Punkte)**
6. Basiert die heutige Wissenschaft auch auf dem Positivismus? Begründen Sie Ihre Antwort! **(2 Punkte)**

Der t-Test

Mehr Informationen zu den Aufgaben in den folgenden Kapiteln aus dem Skript Bio Data Science.

- Kapitel 22 - Der t-Test

71 Aufgabe

(12 Punkte)

Nach einem Experiment mit zwei Pestiziden (*RoundUp* und *GoneEx*) ergibt sich die folgende Datentabelle mit dem gemessenen Trockengewicht (*drymatter*) von Weizen.

pesticide	drymatter
RoundUp	18
GoneEx	17
RoundUp	20
RoundUp	19
RoundUp	18
GoneEx	12
GoneEx	18
GoneEx	11
GoneEx	20
GoneEx	15
GoneEx	18
RoundUp	19

1. Formulieren Sie die wissenschaftliche Fragestellung! **(1 Punkt)**
2. Formulieren Sie das statistische Hypothesenpaar! **(2 Punkte)**
3. Bestimmen Sie die Teststatistik T_{calc} eines Student t-Tests für den Vergleich der beiden Pestizide. Geben Sie den Rechenweg und die Formeln mit an! **(5 Punkte)**
4. Treffen Sie mit $T_{\alpha=5\%} = 2.04$ und dem berechneten T_{calc} eine Aussage zur Nullhypothese! **(2 Punkte)**
5. Wenn Sie keinen Unterschied zwischen den beiden Pestiziden erwarten würden, wie große wäre dann die Teststatistik T_{calc} ? Begründen Sie Ihre Antwort! **(2 Punkte)**

72 Aufgabe

(13 Punkte)

Das Gewicht von Küken wurde *vor* der Behandlung mit STARTex und 1 Woche *nach* der Behandlung gemessen. Es ergibt sich die folgende Datentabelle.

animal_id	before	after
1	13	18
2	16	20
3	16	20
4	11	14
5	15	12
6	16	14
7	14	8
8	14	13
9	11	17

1. Formulieren Sie die Fragestellung! **(1 Punkt)**
2. Formulieren Sie das statistische Hypothesenpaar! **(2 Punkte)**
3. Bestimmen Sie die Teststatistik T_{calc} eines gepaarten t-Tests für den Vergleich der beiden Zeitpunkte. Geben Sie den Rechenweg und die Formeln mit an! **(4 Punkte)**
4. Treffen Sie mit $T_{\alpha=5\%} = 2.04$ und dem berechneten T_{calc} eine Aussage zur Nullhypothese! **(2 Punkte)**
5. Wenn Sie keinen Unterschied zwischen den beiden Zeitpunkten erwarten würden, wie große wäre dann die Teststatistik T_{calc} ? Begründen Sie Ihre Antwort! **(2 Punkte)**
6. Schätzen Sie $Pr(D|H_0)$ ab. Begründen Sie Ihre Antwort! **(2 Punkte)**

73 Aufgabe

(10 Punkte)

Sie erhalten folgende  Ausgabe der Funktion `t.test()`.

```
##  
## Two Sample t-test  
##  
## data: weight by group  
## t = 2.15718, df = 14, p-value = 0.048859  
## alternative hypothesis: true difference in means between group high and group low is not equal to 0  
## 95 percent confidence interval:  
## 0.021332578 7.407238850  
## sample estimates:  
## mean in group ctrl mean in group trt2  
## 17.714286 14.000000
```

1. Formulieren Sie das statistische Hypothesenpaar! **(2 Punkte)**
2. Liegt ein signifikanter Unterschied zwischen den Gruppen vor? Begründen Sie Ihre Antwort! **(2 Punkte)**
3. Skizzieren Sie eine Abbildung in der Sie T_{calc} , $Pr(D|H_0)$, $A = 0.95$, sowie $T_{\alpha=5\%} = |2.14|$ einzeichnen! **(4 Punkte)**
4. Beschriften Sie die Abbildung entsprechend! **(2 Punkte)**

74 Aufgabe

(8 Punkte)


Sie erhalten folgende  Ausgabe der Funktion `t.test()`.

```
##  
## Two Sample t-test  
##  
## data: weight by group  
## t = -0.0449272, df = 14, p-value = 0.9648  
## alternative hypothesis: true difference in means between group high and group low is not equal to 0  
## 95 percent confidence interval:  
## -4.6418233 4.4513471  
## sample estimates:  
## mean in group ctrl mean in group low  
## 14.333333 14.428571
```

1. Formulieren Sie die wissenschaftliche Fragestellung! **(2 Punkte)**
2. Liegt ein signifikanter Unterschied zwischen den Gruppen vor? Begründen Sie Ihre Antwort! **(2 Punkte)**
3. Skizzieren Sie das sich ergebende 95% Konfidenzintervall! **(2 Punkte)**
4. Beschriften Sie die Abbildung und das 95% Konfidenzintervall entsprechend! **(2 Punkte)**

75 Aufgabe

(8 Punkte)

Sie erhalten folgende  Ausgabe der Funktion `t.test()`.

```
##  
## Welch t-Test  
##  
## data: weight by group  
## t = -2.86109, df = 12, p-value = 0.014325  
## alternative hypothesis: true difference in means between group high and group low is not equal to 0  
## 95 percent confidence interval:  
## -6.18493991 -0.83728231  
## sample estimates:  
## mean in group ctrl mean in group trt2  
## 15.600000 19.111111
```

1. Formulieren Sie die wissenschaftliche Fragestellung! **(2 Punkte)**
2. Liegt ein signifikanter Unterschied zwischen den Gruppen vor? Begründen Sie Ihre Antwort! **(2 Punkte)**
3. Skizzieren Sie die sich ergebenden Boxplots! Welche Annahmen an die Daten haben Sie getroffen? Begründen Sie Ihre Antwort! **(4 Punkte)**

76 Aufgabe

(8 Punkte)

Sie erhalten folgende  Ausgabe der Funktion `t.test()`.

```
##  
## Paired t-test  
##  
## data: waterintake by infusion  
## t = -1.50756, df = 6, p-value = 0.18239  
## alternative hypothesis: true mean difference is not equal to 0  
## 95 percent confidence interval:  
## -7.4945649 1.7802791  
## sample estimates:  
## mean difference  
## -2.8571429
```

1. Formulieren Sie das statistische Hypothesenpaar! **(2 Punkte)**
2. Liegt ein signifikanter Unterschied zwischen den Gruppen vor? Begründen Sie Ihre Antwort! **(2 Punkte)**
3. Skizzieren Sie den sich ergebenden Datensatz mit $n = 4$! Beobachtungen! **(2 Punkte)**
4. Skizzieren Sie die sich ergebenden Boxplots! Welche Annahmen an die Daten haben Sie getroffen? Begründen Sie Ihre Antwort! **(2 Punkte)**

Die ANOVA

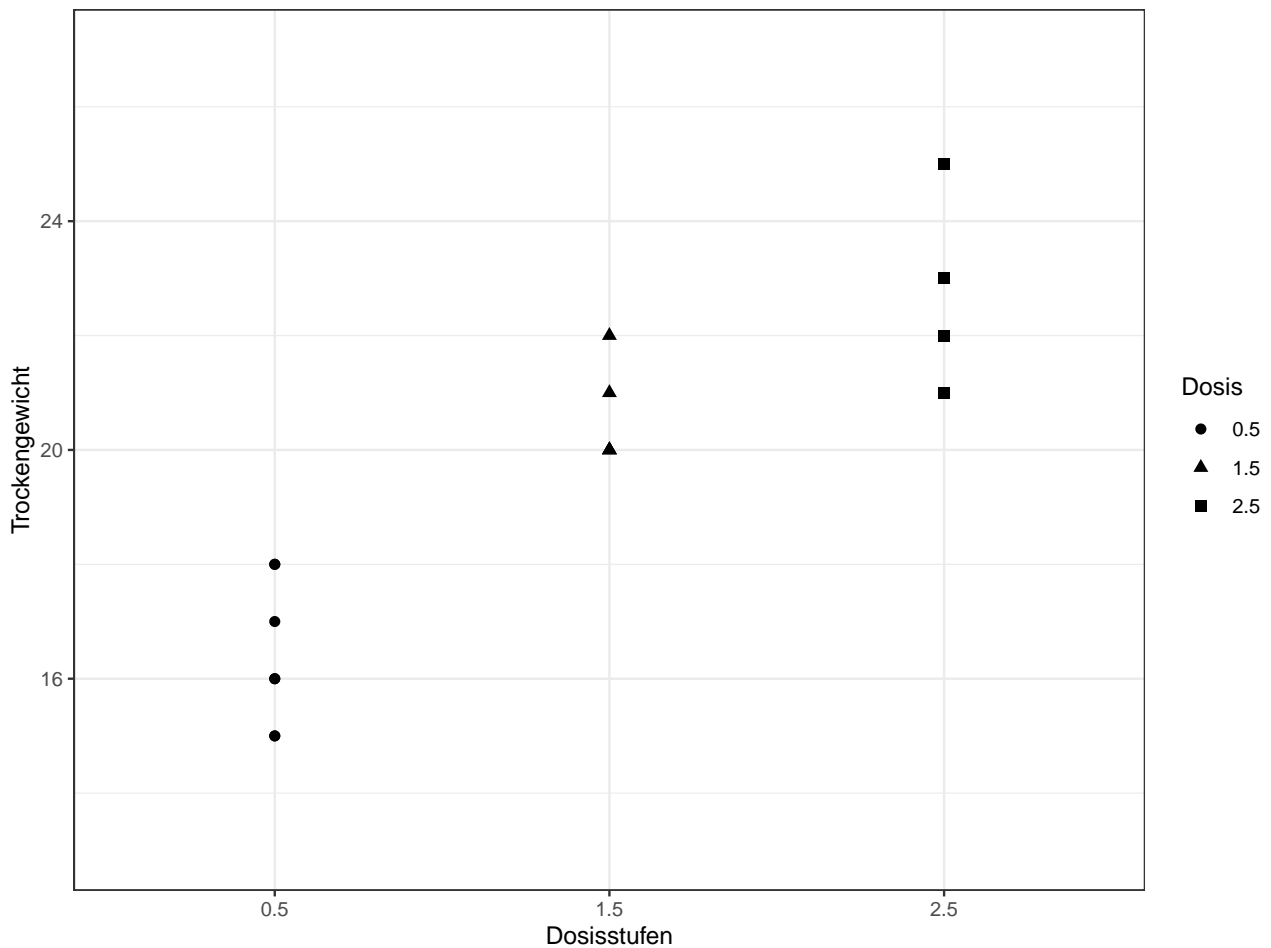
Mehr Informationen zu den Aufgaben in den folgenden Kapiteln aus dem Skript Bio Data Science.

- Kapitel 23 - Die ANOVA

77 Aufgabe

(8 Punkte)

In einem Experiment wurde der Ertrag von Erbsen unter drei verschiedenen Pestizid-Dosen 0.5 g/l, 1.5 g/l und 2.5 g/l gemessen. Unten stehend sehen Sie die Visualisierung des Datensatzes.



1. Zeichnen Sie folgende statistischen Masszahlen in die Abbildung ein! (6 Punkte)
 - Total (grand) mean: β_0
 - Mittelwerte der Dosen: $\bar{y}_{0.5}$, $\bar{y}_{1.5}$ und $\bar{y}_{2.5}$
 - Effekt der einzelnen Level der Dosen: $\beta_{0.5}$, $\beta_{1.5}$, und $\beta_{2.5}$
 - Residuen oder Fehler: ϵ
2. Schätzen Sie den p-Wert einer einfaktoriellen ANOVA ab. Liegt ein *vermutlicher* signifikanter Unterschied zwischen den Dosisstufen vor? Begründen Sie Ihre Antwort! (2 Punkte)

78 Aufgabe

(13 Punkte)

Der Datensatz PlantGrowth enthält das Gewicht der Pflanzen (*weight*), die unter einer Kontrolle und zwei verschiedenen Behandlungsbedingungen erzielt wurden – dem Faktor *group* mit den Faktorstufen *ctrl*, *trt1*, *trt2*.

1. Füllen Sie die unterstehende einfaktorielle ANOVA Ergebnistabelle aus mit den gegebenen Informationen von Df und Sum Sq! **(4 Punkte)**
2. Schätzen Sie den p-Wert der Tabelle mit der Information von $F_{\alpha=5\%} = 3.35$ ab. Begründen Sie Ihre Antwort! **(2 Punkte)**

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
group	2	3.94			
Residuals	27	9.88			

3. Was bedeutet ein signifikantes Ergebnis in einer einfaktoriellen ANOVA im Bezug auf die möglichen Unterschiede zwischen den Gruppen? Beziehen Sie sich auf den obigen Fragetext bei Ihrer Antwort! **(2 Punkte)**
4. Berechnen Sie *einen* Student t-Test mit für den *vermutlich* signifikantesten Gruppenvergleich anhand der untenstehenden Tabelle mit $T_{\alpha=5\%} = 2.03$. Begründen Sie Ihre Auswahl! **(3 Punkte)**

group	n	mean	sd
ctrl	10	5.01	0.55
trt1	10	4.68	0.79
trt2	10	5.56	0.41

5. Gegebenen der ANOVA Tabelle war das Ergebnis des t-Tests zu erwarten? Begründen Sie Ihre Antwort! **(2 Punkte)**

79 Aufgabe

(8 Punkte)

Der Datensatz *Crop* enthält das Trockengewicht der Maispflanzen (*drymatter*), die unter drei verschiedenen Düngerbedingungen erzielt wurden – dem Faktor *trt* mit den Faktorstufen *low*, *mid*, *high*. Sie erhalten folgenden Output in .

```
## Analysis of Variance Table
##
## Response: drymatter
##          Df    Sum Sq Mean Sq F value    Pr(>F)
## trt        2    3.16501  1.582503  3.31222 0.051709 .
## Residuals 27   12.89999  0.477777
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

1. Stellen Sie die statistische H_0 und H_A Hypothese für die obige einfaktorielle ANOVA auf! **(2 Punkte)**
2. Interpretieren Sie das Ergebnis der einfaktoriellen ANOVA! **(2 Punkt)**
3. Berechnen Sie den Effektschätzer η^2 . Was sagt Ihnen der Wert von η^2 aus? **(2 Punkte)**
4. Skizzieren Sie eine Abbildung, der dem obigen Ergebnis der einfaktoriellen ANOVA näherungsweise entspricht! **(2 Punkte)**

80 Aufgabe

(12 Punkte)

Der Datensatz *ToothGrowth* enthält Daten aus einer Studie zur Bewertung der Wirkung von Vitamin C auf das Zahnwachstum bei Meerschweinchen. Der Versuch wurde an 60 Schweinen durchgeführt, wobei jedes Tier eine von drei Vitamin-C-Dosen *dose* (0.5 mg/Tag, 1 mg/Tag und 2 mg/Tag) über eine von zwei Verabreichungsmethoden *supp* erhielt (Orangensaft oder Ascorbinsäure). Die Zahnlänge wurde als normalverteiltes Outcome gemessen.


1. Füllen Sie die unterstehende zweifaktorielle ANOVA Ergebnistabelle aus mit den gegebenen Informationen von Df und Sum Sq! **(4 Punkte)**
2. Schätzen Sie den p-Wert der Tabelle mit der Information von den kritischen F-Werten mit $F_{supp} = 4.02$ und $F_{dose} = 3.17$ sowie $F_{supp:dose} = 3.17$ ab. Begründen Sie Ihre Antwort! **(4 Punkte)**

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
supp	1	210.5			
dose	2	2413.38			
supp:dose	2	108.05			
Residuals	54	717.62			

3. Was bedeutet ein signifikantes Ergebnis in einer zweifaktoriellen ANOVA im Bezug auf die möglichen Unterschiede zwischen den Gruppen? Beziehen Sie sich dabei einmal auf den Faktor *supp* und einmal auf den Faktor *dose*! **(2 Punkte)**
4. Was sagt der Term *supp:dose* aus? Interpretieren Sie das Ergebnis des abgeschätzten p-Wertes! **(2 Punkte)**

81 Aufgabe

(8 Punkte)

Der Datensatz *PigGain* enthält Daten aus einer Studie zur Bewertung der Wirkung vom Vitamin Selen auf das Wachstum bei Mastschweinen. Der Versuch wurde an 74 Mastschweinen durchgeführt, wobei jedes Tier eine von drei Selen-Dosen (0.5 ng/Tag, 1 ng/Tag und 5 ng/Tag) über eine von zwei Verabreichungsmethoden erhielt (Wasser oder Festnahrung). Sie erhalten folgenden Output in .

```
## Analysis of Variance Table
##
## Response: len
##      Df    Sum Sq  Mean Sq  F value    Pr(>F)
## supp     1      5.426     5.426  0.34046  0.56149428
## dose      2 2644.412 1322.206 82.96486 < 0.000000000000000222 ***
## supp:dose  2  307.848   153.924  9.65832  0.00020324 ***
## Residuals 68 1083.712    15.937
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

1. Stellen Sie die statistische H_0 und H_A Hypothese für die obige zweifaktorielle ANOVA für den Faktor *supp* auf! **(2 Punkte)**
2. Interpretieren Sie das Ergebnis der zweifaktoriellen ANOVA. Gehen Sie im besonderen auf den Term *supp : dose* ein! **(2 Punkte)**
3. Zeichnen Sie eine Abbildung, der dem obigen Ergebnis der zweifaktoriellen ANOVA näherungsweise entspricht! **(4 Punkte)**

82 Aufgabe

(8 Punkte)

In der untenstehenden Tabelle ist die Formel für den F-Test aus der ANOVA und die Formel für den t-Test dargestellt. In der ANOVA berechnen Sie die F-Statistik F_{calc} und in dem t-Test die T-Statistik T_{calc} .

$$F_{calc} = \frac{MS_{treatment}}{MS_{error}} \quad T_{calc} = \frac{\bar{y}_1 - \bar{y}_2}{s_p \cdot \sqrt{2/n_g}}$$

1. Erklären Sie den konzeptionellen Zusammenhang zwischen der F_{calc} Statistik und T_{calc} Statistik! **(2 Punkte)**
2. Visualisieren Sie eine nicht signifikante F_{calc} Statistik sowie eine signifikante F_{calc} Statistik anhand von $MS_{treatment}$ und MS_{error} ! Beschriften Sie die Abbildung! **(2 Punkte)**
3. Erklären Sie an der Formel des F-Tests sowie an der Abbildung warum das Minimum der F-Statistik 0 ist! **(2 Punkte)**
4. Wenn die F-Statistik 0 ist, spricht dies eher für oder gegen die Nullhypothese? Begünden Sie Ihre Antwort! **(2 Punkte)**

83 Aufgabe

(6 Punkte)

Sie rechnen eine zweifaktorielle ANOVA und erhalten einen signifikanten Interaktionseffekt zwischen den beiden Faktoren f_1 und f_2 . Der Faktor f_1 hat drei Level. Der Faktor f_2 hat dagegen nur zwei Level.

1. Visualisieren Sie in zwei getrennten Abbildungen eine mittlere und eine keine Interaktion zwischen den Faktoren f_1 und f_2 ! **(2 Punkte)**
2. Erklären Sie den Unterschied zwischen den beiden Stärken der Interaktion! **(2 Punkte)**
3. Wenn eine signifikante Interaktion in den Daten vorliegt, wie ist dann das weitere Vorgehen bei einem Posthoc-Test? **(2 Punkte)**

84 Aufgabe

(7 Punkte)

Sie rechnen eine einfaktorielle ANOVA mit einem Faktor f_1 mit drei Leveln. Nachdem Sie die einfaktorielle ANOVA gerechnet haben, erhalten Sie einen p-Wert von 0.078 und eine F Statistik mit $F_{calc} = 1.2$. Als Sie sich die Boxplots der Behandlungen anschauen, stellen Sie fest, dass es eigentlich einen Mittelwertsunterschied zwischen dem zweiten und dritten Level geben müsste. Die IQR-Bereiche überlappen sich nicht und die Mediane liegen auch weit vom globalen Mittel entfernt.

1. Erklären Sie die Annahme der Normalverteilung und die Annahme der Varianzhomogenität für eine ANOVA an einer passenden Abbildung! **(2 Punkte)**
2. Visualisieren Sie die Berechnung von F_{calc} am obigen Beispiel! **(2 Punkte)**
3. Erklären Sie das Ergebnis der obigen einfaktoriellen ANOVA unter der Berücksichtigung der Annahmen an eine ANOVA! **(3 Punkte)**

Der χ^2 -Test & Der diagnostische Test

Mehr Informationen zu den Aufgaben in den folgenden Kapiteln aus dem Skript Bio Data Science.

- Kapitel 28 - Der χ^2 -Test
- Kapitel 29 - Der diagnostische Test

85 Aufgabe

(10 Punkte)

Nach einem Experiment ergibt sich die folgende 2x2 Datentabelle mit einem Pestizid (ja/nein), dargestellt in den Zeilen. Im Weiteren mit dem infizierten Pflanzenstatus (ja/nein) in den Spalten. Insgesamt wurden $n = 128$ Pflanzen untersucht.

	Erkrankt (ja)	Erkrankt (nein)	
Pestizid (ja)	38	19	
Pestizid (nein)	27	44	

1. Ergänzen Sie die Tabelle um die Randsummen! (1 Punkt)
2. Formulieren Sie die Fragestellung! (1 Punkt)
3. Formulieren Sie das Hypothesenpaar! (2 Punkte)
4. Berechnen Sie die Teststatistik eines Chi-Quadrat-Test auf der 2x2 Tafel. Geben Sie Formeln und Rechenweg mit an! (4 Punkte)
5. Treffen Sie eine Entscheidung im Bezug zu der Nullhypothese gegeben einem $T_k = 3.841$! (1 Punkt)
6. Skizzieren Sie eine 2x2 Tabelle mit $n = 30$ Pflanzen in dem *vermutlich* die Nullhypothese nicht abgelehnt werden kann! (1 Punkt)

86 Aufgabe

(7 Punkte)

Gegeben sind folgende Randsummen in einer 2x2 Kreuztabelle aus einem Experiment mit $n = 141$ Sauen. In dem Experiment wurde gemessen, ob eine Sau nach einer Behandlung mit einem Medikament (ja/nein) mehr als 30 Ferkel pro Jahr bekommen konnte (ja/nein).

	>30 Ferkel (ja)	≤30 Ferkel (nein)	
Medikament (ja)			92
Medikament (nein)			49
	82	59	141

1. Ergänzen Sie die Felder innerhalb der 2x2 Kreuztabelle in dem Sinne, dass *kein* signifikanter Effekt zu erwarten wäre! **(2 Punkte)**
2. Erklären und Begründen Sie Ihr Vorgehen an der Formel des Chi-Quadrat-Tests mit

$$\chi^2 = \sum \frac{(O-E)^2}{E}.$$

Sie können dies an einem Beispiel erklären! **(2 Punkte)**

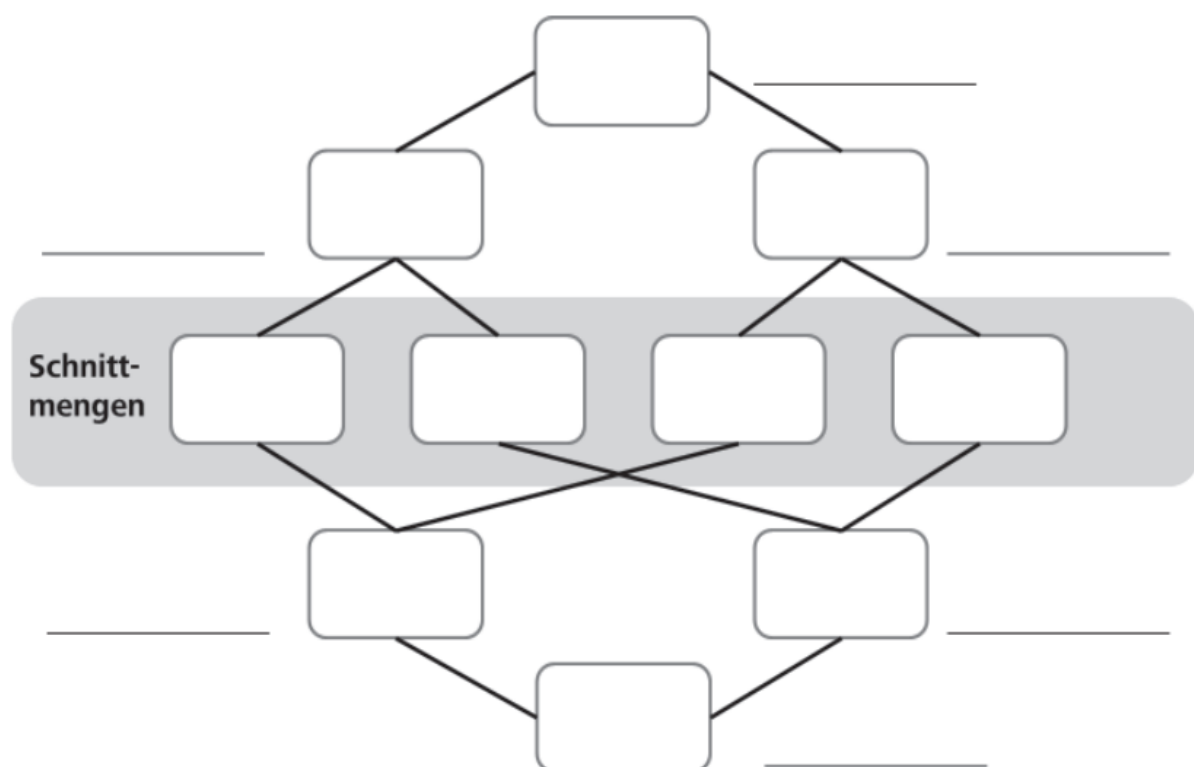
3. Was ist die Mindestanzahl an Beobachtungen je Zelle? Wenn in einer der Zellen weniger Beobachtungen auftreten, welchen Test können Sie anstatt des „normalen“ Chi-Quadrat-Tests anwenden? **(2 Punkte)**
4. Warum hat die obige Vierfeldertafel einen Freiheitsgrad von $df = 1$? **(1 Punkt)**

87 Aufgabe

(11 Punkte)

Die Prävalenz von Klauenseuche bei Wollschweinen wird mit 4% angenommen. In 75% der Fälle ist ein Test positiv, wenn das Wollschwein erkrankt ist. In 8% der Fälle ist ein Test positiv, wenn das Wollschwein *nicht* erkrankt ist und somit gesund ist. Sie werten 2000 Wollschweine mit einem diagnostischen Test auf Klauenseuche aus.

1. Füllen und beschriften Sie den untenstehenden Doppelbaum! Beschriften Sie auch die Äste des Doppelbaumes, mit denen Ihnen bekannten Informationen! **(8 Punkte)**
2. Berechnen Sie die Wahrscheinlichkeit $Pr(K^+|T^+)$! **(2 Punkte)**
3. Was sagt Ihnen die Wahrscheinlichkeit $Pr(K^+|T^+)$ aus? **(1 Punkt)**



88 Aufgabe

(9 Punkte)

Beim diagnostischen Testen erhalten Sie *True Positives (TP)*, *True Negatives (TN)*, *False Positives (FP)* und *False Negatives (FN)*. Erklären Sie den Zusammenhang wie folgt.

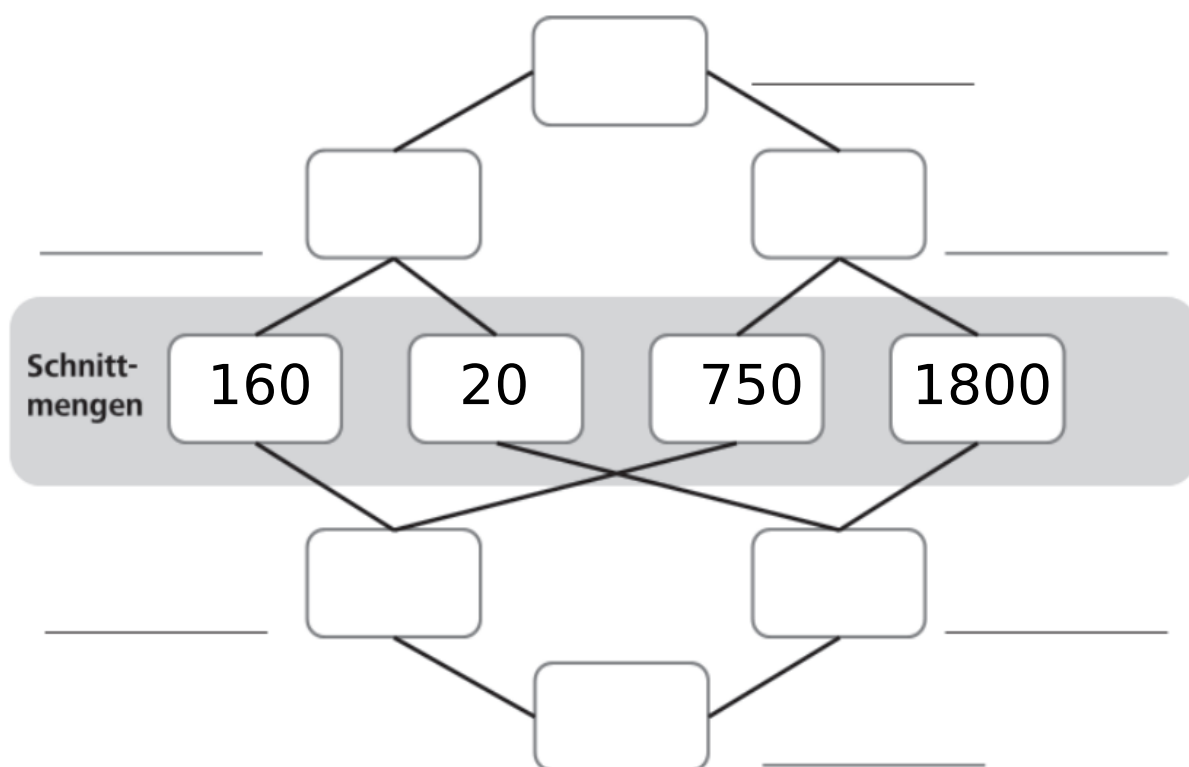
1. Visualisieren Sie *TP*, *TN*, *FP* und *FN* in einer Abbildung. Beschriften Sie die Abbildung und die Achsen entsprechend! **(6 Punkte)**
2. Tragen Sie *TP*, *TN*, *FP* und *FN* in eine 2x2 Kreuztabelle ein. Beschriften Sie die Tabelle entsprechend! **(3 Punkte)**

89 Aufgabe

(12 Punkte)

Folgender diagnostischer Doppelbaum nach der Testung auf Klauenseuche bei Fleckvieh ist gegeben.

1. Füllen und beschriften Sie den untenstehenden Doppelbaum! (4 Punkte)
2. Berechnen Sie die Wahrscheinlichkeit $Pr(K^+|T^+)$! (2 Punkte)
3. Berechnen Sie die Prävalenz für Klauenseuche! (2 Punkte)
4. Berechnen Sie die Sensitivität und Spezifität des diagnostischen Tests für Klauenseuche! Erstellen Sie dafür zunächst eine 2x2 Kreuztabelle aus dem ausgefüllten Doppelbaum! (4 Punkte)



Simple linear Regression

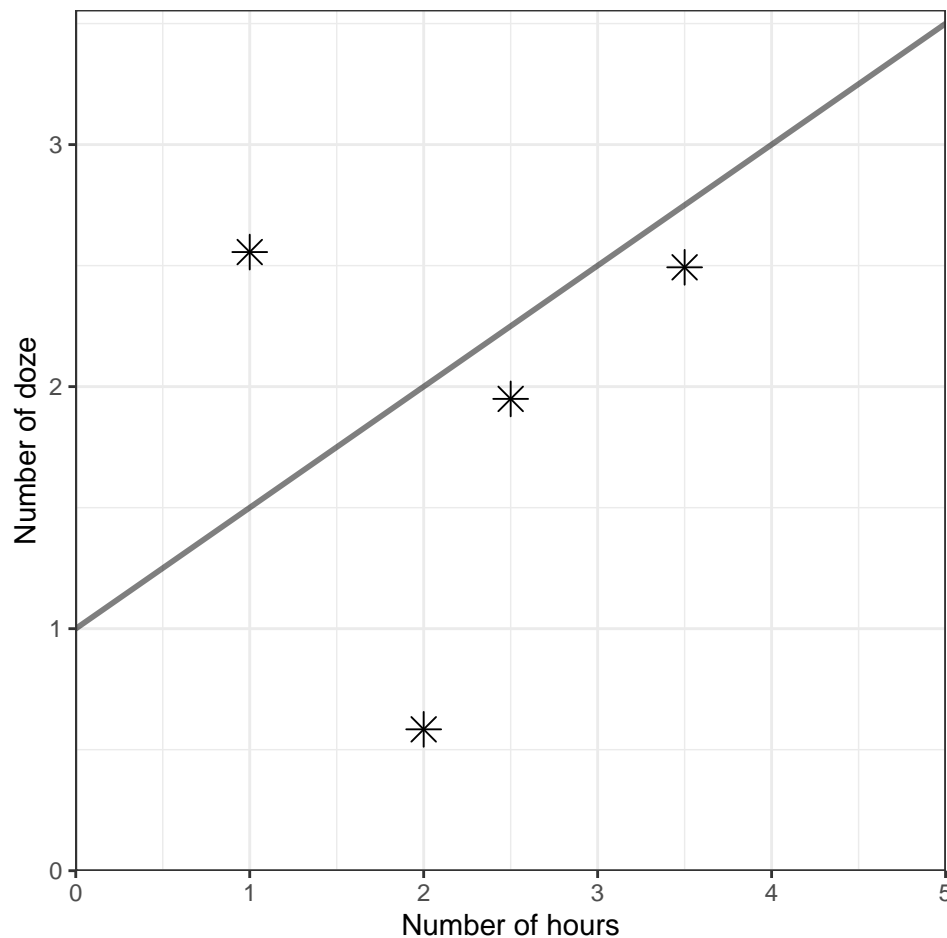
Mehr Informationen zu den Aufgaben in den folgenden Kapiteln aus dem Skript Bio Data Science.

- [Kapitel 32 - Simple lineare Regression](#)
- [Kapitel 33 - Maßzahlen der Modelgüte](#)
- [Kapitel 34 - Korrelation](#)

90 Aufgabe

(7 Punkte)


In einer Studie zur „Arbeitssicherheit auf dem Feld“ wurde gemessen wie viele Stunden auf einem Feld gefahren wurden und wie oft der Fahrer dabei drohte einzunicken. Es ergab sich folgende Abbildung.



1. Erstellen Sie die Regressionsgleichung aus der obigen Abbildung in der Form $y \sim \beta_0 + \beta_1 \cdot x$! **(2 Punkte)**
2. Beschriften Sie die Gerade mit den Parametern der linearen Regressionsgleichung! **(2 Punkte)**
3. Liegt ein Zusammenhang zwischen der Anzahl an gefahrenen Runden und der Müdigkeit vor? Begründen Sie Ihre Antwort! **(2 Punkte)**
4. Wenn kein Zusammenhang zu beobachten wäre, wie würde die Gerade aussehen? **(1 Punkt)**

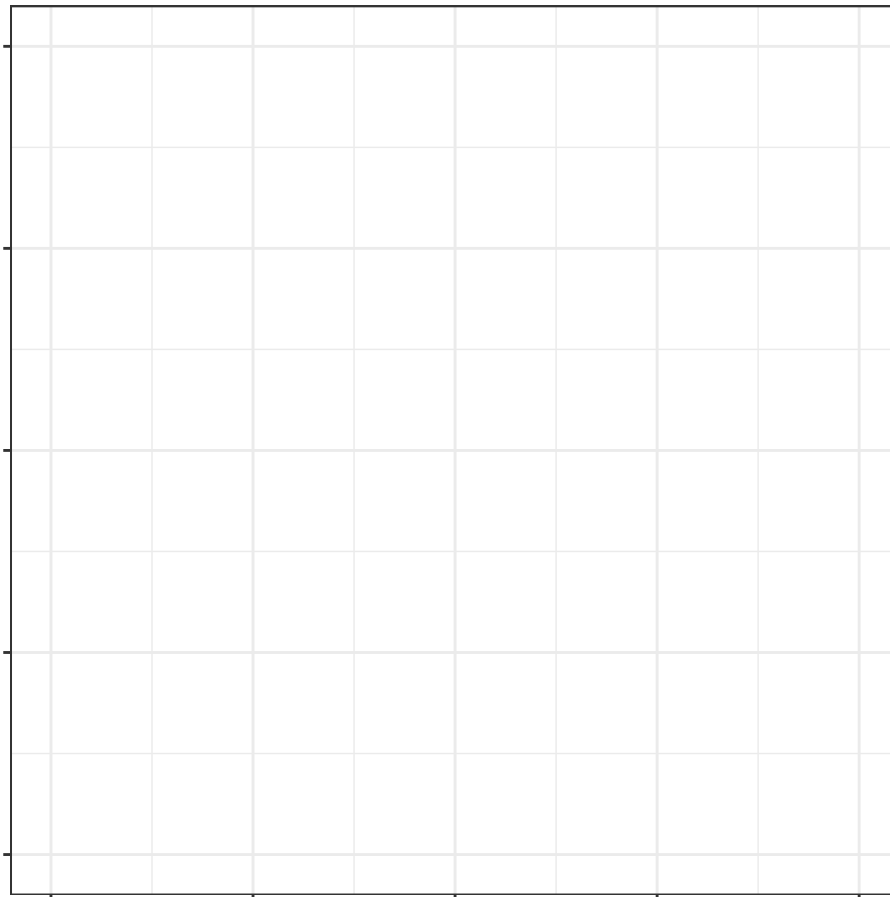
91 Aufgabe

(10 Punkte)

In einem Stallexperiment mit $n = 103$ Ferkeln wurde der Gewichtszuwachs unter bestimmten Lichtverhältnissen gemessen. Sie erhalten den  Output der Funktion `tidy()` einer simplen Gaussian linearen Regression sieben Wochen nach der ersten Messung.

term	estimate	std.error	t statistic	p-value
(Intercept)	25.02	0.99		
light	1.50	0.10		

1. Berechnen Sie die t Statistik für *(Intercept)* und *light*! **(2 Punkte)**
2. Schätzen Sie den p-Wert für *(Intercept)* und *light* mit $T_k = 1.96$ ab. Was sagt Ihnen der p-Wert aus? Begründen Sie Ihre Antwort! **(3 Punkte)**
3. Zeichnen Sie die Gerade aus der obigen Tabelle in die untenstehende Abbildung! **(1 Punkt)**
4. Beschriften Sie die Abbildung und die Gerade mit den statistischen Kenngrößen! **(2 Punkte)**
5. Formulieren Sie die Regressionsgleichung! **(2 Punkte)**




92 Aufgabe

(6 Punkte)

Sie erhalten folgende R Ausgabe der Funktion `lm()`.

```
##
## Call:
## lm(formula = rsp ~ trt, data = data_tbl)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -4.00  -1.25   0.00   2.00   4.00
##
## Coefficients:
##              Estimate Std. Error t value      Pr(>|t|)
## (Intercept)  19.0000     1.0954  17.3445 0.000000008597 ***
## trtB         8.0000     1.4343   5.5777  0.0002348 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 2.4495 on 10 degrees of freedom
## Multiple R-squared:  0.75676, Adjusted R-squared:  0.73243
## F-statistic: 31.111 on 1 and 10 DF,  p-value: 0.00023485
```

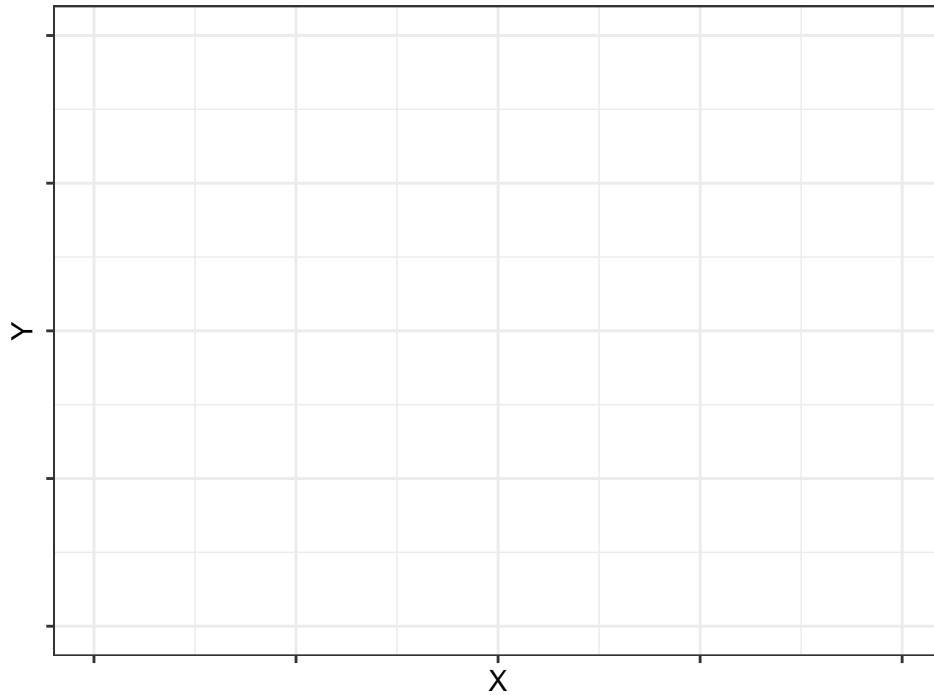
1. Ist die Annahme der Normalverteilung an das Outcome *rsp* erfüllt? Begründen Sie die Antwort! **(2 Punkte)**
2. Wie groß ist der Effekt des *Trt*? Liegt ein signifikanter Effekt vor? Geben Sie die Formel und den Rechenweg mit an! **(2 Punkte)**
3. Schreiben Sie das Ergebnis der  Ausgabe in einen Satz nieder, der die Information zum Effekt und der Signifikanz enthält! **(2 Punkte)**

93 Aufgabe

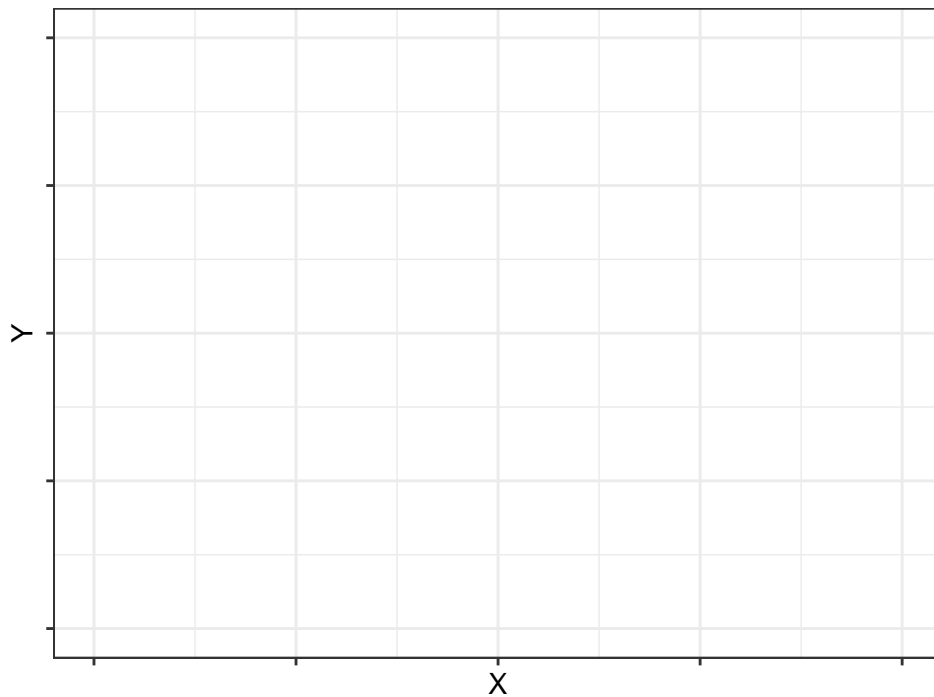
(6 Punkte)

1. Skizzieren Sie in die unten stehenden, freien Abbildungen ein kausales und ein prädiktives Modell mit $n = 5$ Beobachtungen! **(4 Punkte)**
2. Beachten Sie bei der Erstellung der Skizze, ob ein Effekt von X vorliegt oder nicht! **(2 Punkte)**

Causal model with no effect of X



Predictive model with no effect of X



94 Aufgabe

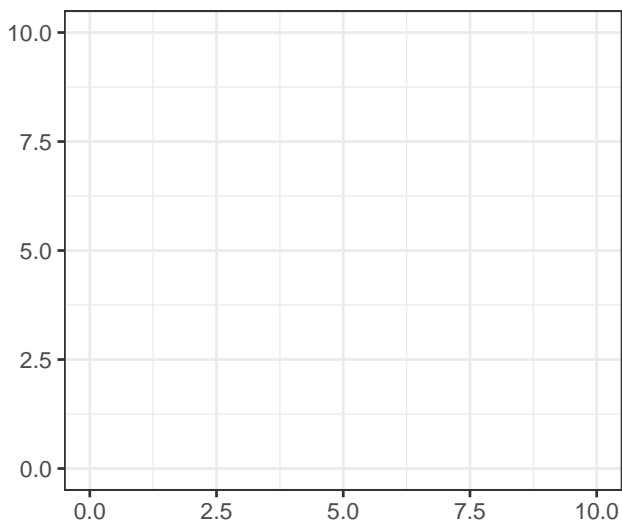
(9 Punkte)

Im folgenden sehen Sie drei leere Scatterplots. Füllen Sie diese Scatterplots nach folgenden Anweisungen.

1. Zeichnen Sie für die angegebene ρ -Werte eine Gerade in die entsprechende Abbildung! **(3 Punkte)**
2. Zeichnen Sie für die angegebenen R^2 -Werte die entsprechende Punktwolke um die Gerade. **(3 Punkte)**
3. Sie rechnen ein statistisches Modell. Was sagen Ihnen die R^2 -Werte über das jeweilige Modell? **(3 Punkte)**

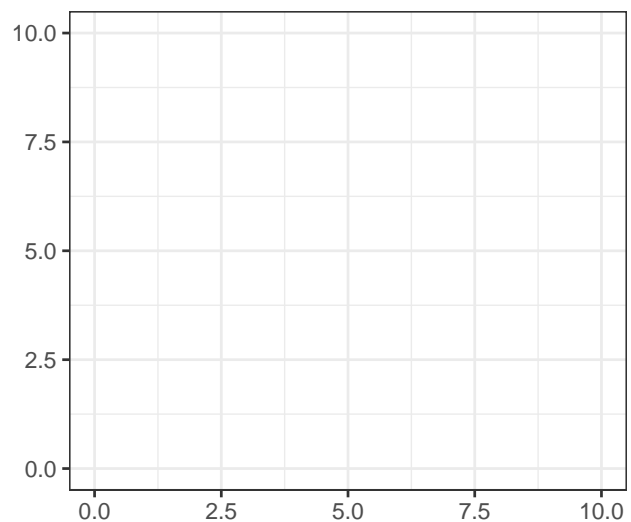
Pearsons $\rho = 1$

$R^2 = 1$



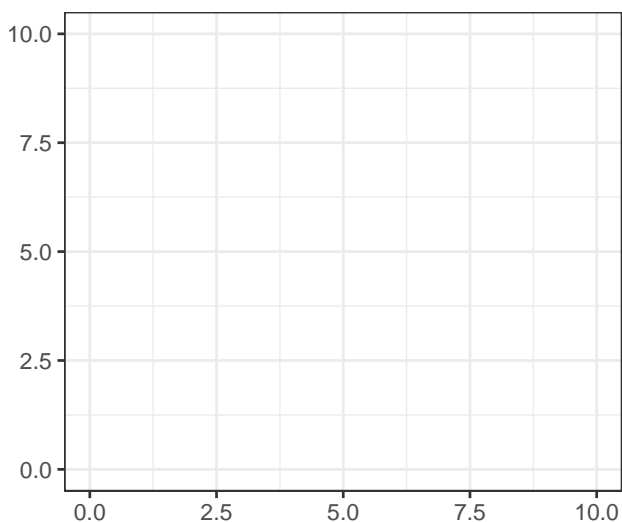
Pearsons $\rho = -0.25$

$R^2 = 0.5$



Pearsons $\rho = -1$

$R^2 = 0.75$



95 Aufgabe

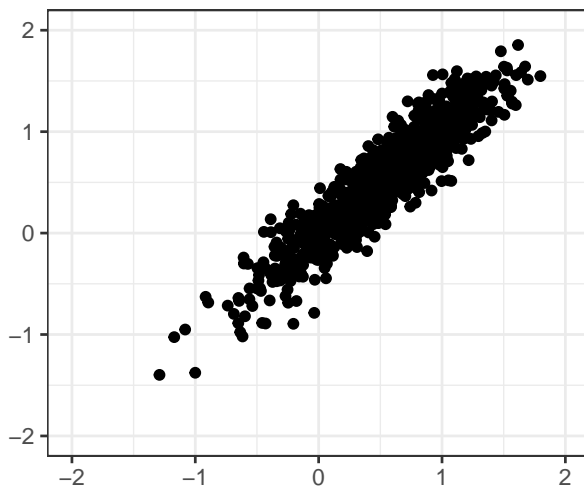
(10 Punkte)

Im folgenden sehen Sie vier Scatterplots. Ergänzen Sie die Überschriften der jeweiligen Scatterplots.

1. Schätzen Sie die ρ -Werte in der entsprechenden Abbildung! **(4 Punkte)**
2. Schätzen Sie die R^2 -Werte in der entsprechenden Punktwolke um die Gerade! **(4 Punkte)**
3. Sie rechnen ein statistisches Modell. Was sagen Ihnen die R^2 -Werte über das jeweilige Modell? **(2 Punkte)**

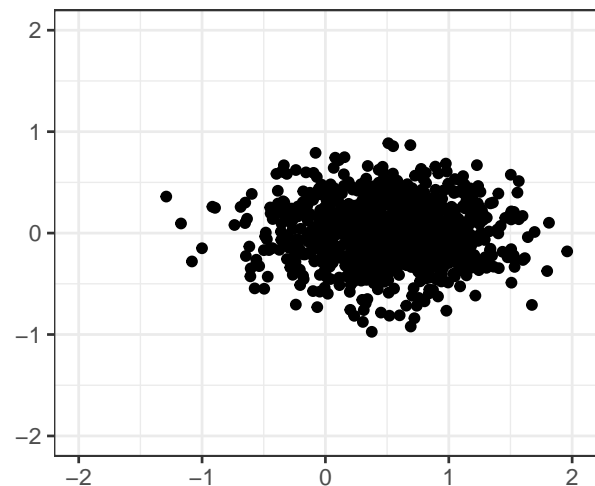
Pearsons $\rho =$

$R^2 =$



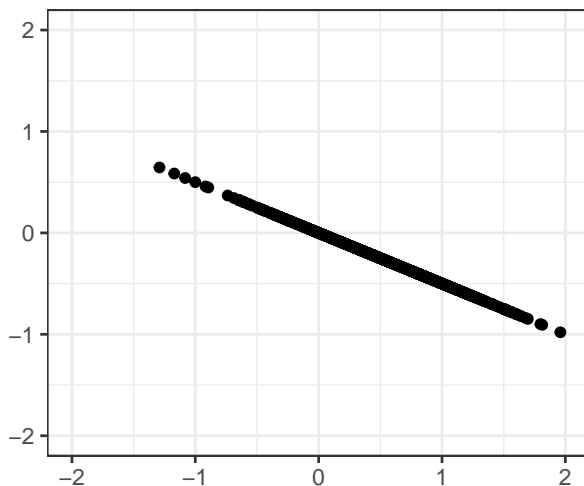
Pearsons $\rho =$

$R^2 =$



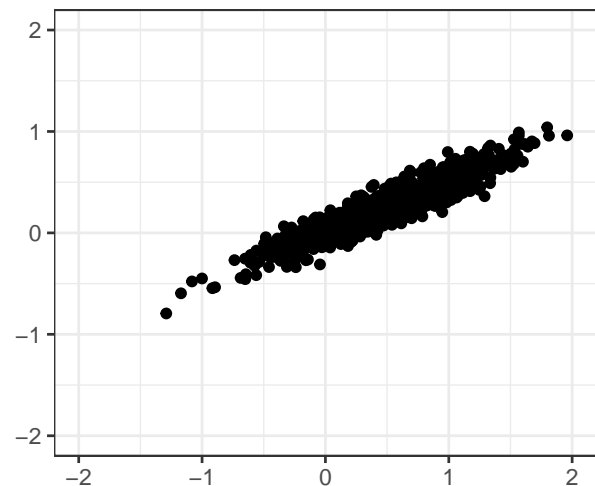
Pearsons $\rho =$

$R^2 =$



Pearsons $\rho =$

$R^2 =$

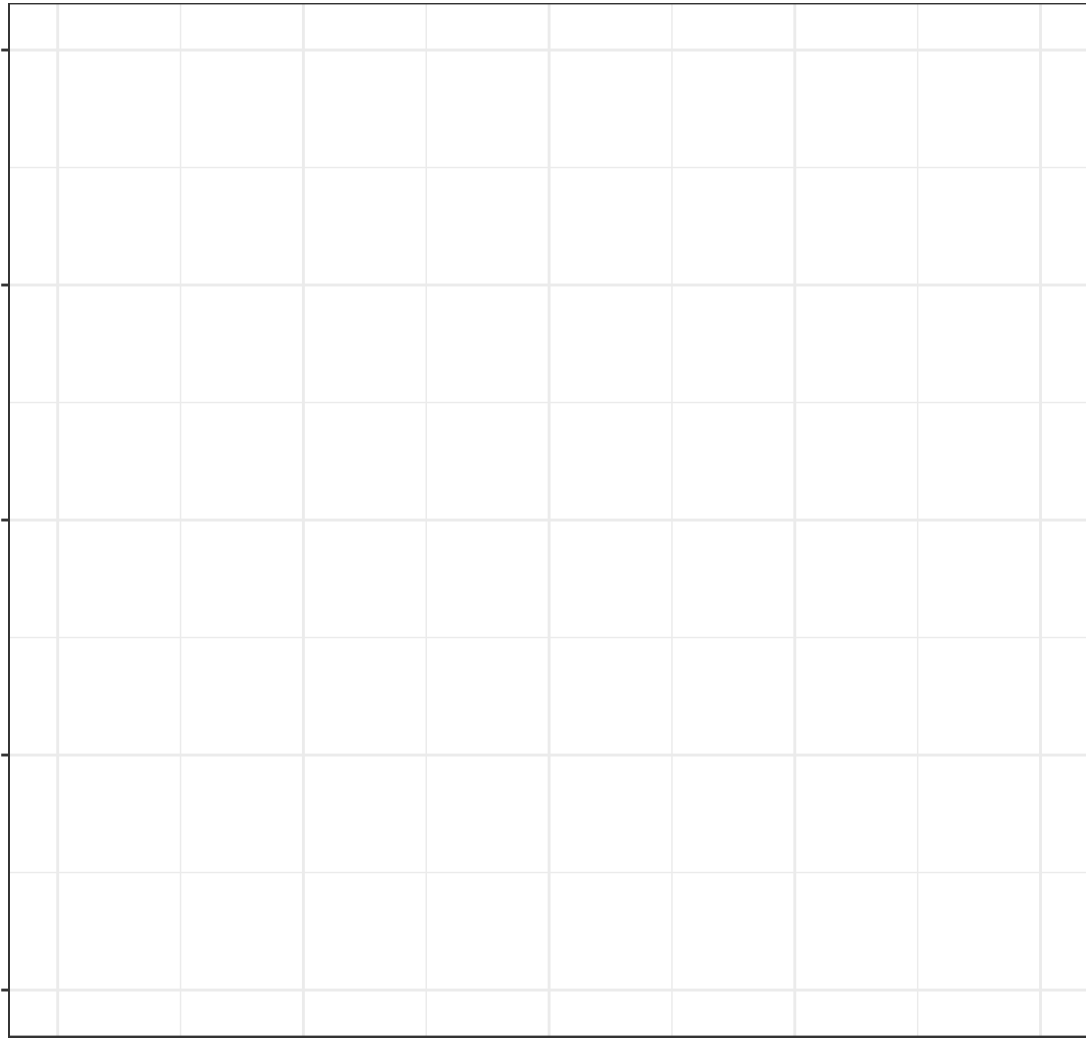


96 Aufgabe

(4 Punkte)

1. Skizzieren Sie in die unten stehende, freie Abbildung die Abbildung, die sich nach der Überschrift ergibt! **(2 Punkte)**
2. Beschriften Sie die Achsen entsprechend! **(2 Punkte)**

Residual plot with 3 outlier fulfilling the normality assumption.



97 Aufgabe

(4 Punkte)

1. Skizzieren Sie in die unten stehende, freie Abbildung die Abbildung, die sich nach der Überschrift ergibt! **(2 Punkte)**
2. Beschriften Sie die Achsen entsprechend! **(2 Punkte)**

QQ plot with residuals validating the normality assumption.



98 Aufgabe

(10 Punkte)

Sie rechnen eine lineare Regression um nach einem Feldexperiment den Zusammenhang zwischen Trockengewicht kg/m^2 (*drymatter*) und Wassergabe l/m^2 (*water*) bei Spargel zu bestimmen. Sie erhalten folgende Datentabelle.

.id	drymatter	water	.fitted	.resid
1	15.6	3.5	14.1	
2	24.3	12.5	28.2	
3	37.1	18.4	37.4	
4	19.8	6.4	18.7	
5	29.0	12.7	28.5	
6	15.4	6.2	18.3	
7	6.4	-1.3	6.8	
8	30.2	12.6	28.3	
9	30.0	13.6	29.9	
10	30.9	12.7	28.5	
11	14.5	5.6	17.4	
12	29.4	13.3	29.4	
13	19.0	4.6	15.9	

1. Ergänzen Sie die Werte in der Spalte `.resid` in der obigen Tabelle. Geben Sie den Rechenweg und Formel mit an! **(4 Punkte)**
2. Zeichnen Sie den sich aus der obigen Tabelle ergebenden Residualplot. Beschriften Sie die Abbildung! **(4 Punkte)**
3. Gibt es auffällige Werte anhand des Residualplots? Begründen Sie Ihre Antwort! **(2 Punkte)**

Multiple lineare Regression

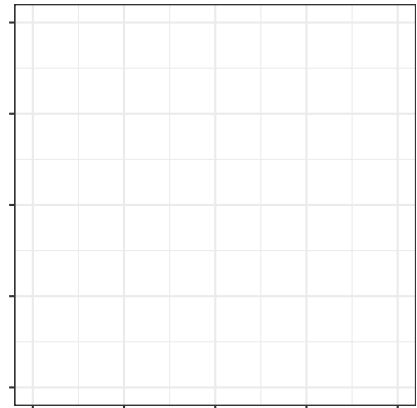
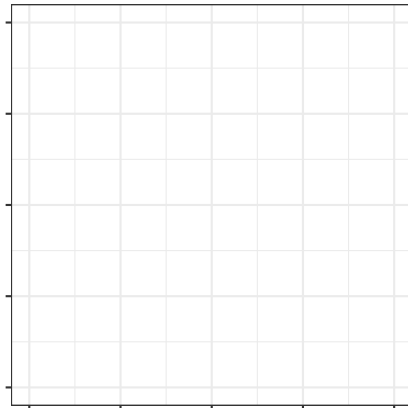
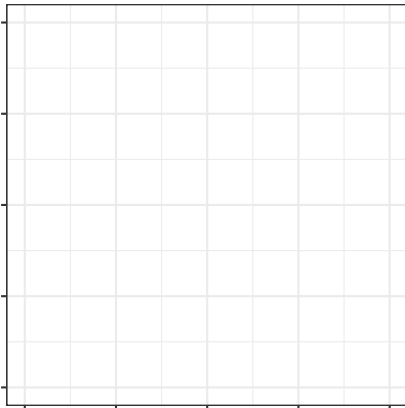
Mehr Informationen zu den Aufgaben in den folgenden Kapiteln aus dem Skript Bio Data Science.

- Kapitel 35 - Multiple lineare Regression
- Kapitel 40 - Gaussian Regression
- Kapitel 41 - Poisson Regression
- Kapitel 43 - Logistische Regression
- Kapitel 44 - Lineare gemischte Modelle

99 Aufgabe

(10 Punkte)



1. Zeichnen Sie in die drei untenstehenden, leeren Abbildungen die Zeile des Regressionskreuzes der Poissonverteilung. Wählen Sie die Beschriftung der y-Achse sowie der x-Achse entsprechend aus! **(6 Punkte)**
2. Welchen Effektschätzer erhalten Sie aus der entsprechend linearen Regression? Geben Sie ein Beispiel! **(2 Punkte)**
3. Wenn Sie keinen Effekt erwarten, welchen *Zahlenraum* nimmt dann der Effektschätzer ein? Geben Sie ein Beispiel! **(2 Punkte)**



100 Aufgabe

(9 Punkte)

Ein Feldexperiment wurde mit $n = 200$ Pflanzen durchgeführt. Folgende Einflussvariablen (x) wurden erhoben: fertilizier, S und region. Als mögliche Outcomevariablen stehen Ihnen nun folgende gemessene Endpunkte zu Verfügung: drymatter, yield, count, quality_score und dead.

1. Wählen Sie ein Outcome was zu der Verteilungsfamilie *Gaussian* gehört! **(1 Punkt)**
2. Schreiben Sie das Modell in der Form $y \sim x$ wie es in  üblich ist *ohne Interaktionsterm*! **(2 Punkte)**
3. Schreiben Sie das Modell in der Form $y \sim x$ wie es in  üblich ist und ergänzen Sie *einen* Interaktionsterm nach Wahl! **(2 Punkte)**
4. Zeichnen Sie eine *schwache* Interaktion in die Abbildung unten für den Endpunkt *yield*. Ergänzen Sie eine aussagekräftige Legende. Wie erkennen Sie eine Interaktion? Begründen Sie Ihre Antwort! **(4 Punkte)**




101 Aufgabe



(8 Punkte)

In verschiedenen Flüssen (*stream*) wurde die Anzahl an Knochenhechten (*longnose*) gezählt. Daneben wurden noch andere Eigenschaften der entsprechenden Flüsse gemessen. Es ergibt sich folgender Auszug aus den Daten.

stream	longnose	temp	acerage	maxdepth	so4
CABIN_JOHN_CR	21	15.0	8612	103	16.09
HERRINGTON_RUN	1	24.5	7056	60	9.82
MORGAN_RUN	76	16.9	9765	130	5.74
HAINES_BR	98	16.8	1967	50	26.44

Sie rechnen nun eine Poisson lineare Regression auf den Daten und erhalten folgenden  Output.

```
##
## Call:
## glm(formula = reformulate(response = "longnose", termlabels = wanted_vec),
##      family = quasipoisson, data = data_tbl)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -9.5434  -3.6392  -1.9086   2.1575  16.8042
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  2.304994027  0.923167174   2.4968 0.015941 *
## temp         0.013662071  0.042506415   0.3214 0.749265
## acerage      0.000041398  0.000014371   2.8806 0.005872 **
## maxdepth     0.008333758  0.004151677   2.0073 0.050247 .
## so4          -0.001398368  0.027152794  -0.0515 0.959136
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for quasipoisson family taken to be 36.134413)
##
## Null deviance: 2008.56  on 53  degrees of freedom
## Residual deviance: 1486.81  on 49  degrees of freedom
## AIC: NA
##
## Number of Fisher Scoring iterations: 5
```


1. Warum wurde hier eine Poisson bzw. Quasipoisson-Verteilung gewählt? Begründen Sie Ihre Antwort mit dem  Output! **(2 Punkte)**
2. Können Sie die *Estimate* der einzelnen Einflussvariablen direkt interpretieren? Begründen Sie Ihre Antwort! **(2 Punkte)**
3. Interpretieren Sie die *signifikanten* Effekte auf die Anzahl an Knochenhechten! **(2 Punkte)**
4. Erklären Sie am  Output wie sich die *t value* Spalte errechnet! **(2 Punkte)**

102 Aufgabe


(10 Punkte)

Auf einer Erdbeerplantage treten unerwartet häufig infizierte Erdbeerpflanzen auf. In einem Versuch sollen verschiedene Einflussfaktoren auf den Infektionsstatus betrachtet werden. Dafür wurde für jede Erdbeerpflanze gemessen, wieviel Wasser die Pflanze erhalten hat oder ob die Pflanze ein neuartiges Lichtregime erhalten hatte. Zusätzlich wurde die Anzahl an Nematoden im Boden bestimmt. Es ergibt sich folgender Auszug aus den Daten.

	infected	water	light	nematodes
0		12.18	0	5
1		9.15	1	1
1		8.80	0	1
1		9.34	1	1

Sie rechnen nun eine logistische lineare Regression auf den Daten und erhalten folgenden  Output.

```
## # A tibble: 3 x 4
##   term          std.error statistic p.value
##   <chr>          <dbl>      <dbl>   <dbl>
## 1 (Intercept)    0.371      0.768   0.443
## 2 nematodes      0.148     -1.07   0.283
## 3 light          0.472      0.853   0.394
```


1. Die Spalte *estimate* wurde gelöscht. Berechnen Sie die Werte der Spalte *estimate* aus den  Output! (2 Punkte)
2. Welche Einflussfaktoren sind protektiv, welche ein Risiko? Berechnen Sie hierfür zunächst das OR aus der Spalte *estimate*! (4 Punkte)
3. Interpretieren Sie die Spalte *estimate* im Bezug auf den Infektionsstatus der Erdbeerpflanzen! (2 Punkte)
4. Was ist der Unterschied zwischen einem OR und einem RR? Geben Sie ein numerisches Beispiel! (2 Punkte)

103 Aufgabe

(7 Punkte)

Maispflanzen sollen auf die ertragssteigernde Wirkung von verschiedenen Einflussfaktoren untersucht werden. Gemessen wurde als Outcome die Trockenmasse in kg/m^2 . Dafür wurde für jede Maispflanze gemessen wieviel Wasser (l/m^2) die Pflanze erhalten hat oder ob die Pflanze ein neuartiges Lichtregime (0 = alt, 1 = neu) erhalten hatte. Zusätzlich wurde die Anzahl an Nematoden im Boden bestimmt sowie der Eisen- und Phosphorgehalt ($\mu\text{g/kg}$) des Bodens. Es ergibt sich folgender Auszug aus den Daten.

water	light	P	Fe	drymatter	nematodes
10.21	0	12.31	99.75	70.19	12
10.61	1	9.20	100.75	73.42	12
8.14	1	9.53	106.78	73.93	9
9.44	0	9.37	96.88	68.60	13

Sie rechnen nun eine Gaussian lineare Regression auf den Daten und erhalten folgenden  Output.

```
##
## Call:
## lm(formula = reformulate(response = "drymatter", termlabels = wanted_vec),
##     data = data_tbl)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -7.89156 -2.03981 -0.02329  2.01449  7.75142
##
## Coefficients:
##              Estimate Std. Error t value      Pr(>|t|)
## (Intercept)  72.58567    3.33148  21.7878 <0.0000000000000002 ***
## water        -0.34342    0.30420  -1.1289    0.2622
## nematodes     0.04228    0.14658   0.2884    0.7737
## light         0.71388    0.84532   0.8445    0.4008
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 3.5123 on 82 degrees of freedom
## Multiple R-squared:  0.028148, Adjusted R-squared:  -0.007408
## F-statistic: 0.79165 on 3 and 82 DF, p-value: 0.502
```


1. Welche der Einflussfaktoren sind signifikant? Begründen Sie Ihre Antwort! **(3 Punkte)**
2. Interpretieren Sie die Spalte *estimate* im Bezug auf den Ertrag in Trockenmasse der Maispflanzen! **(2 Punkte)**
3. Sind die Residuals approximativ Normalverteilt? Begründen Sie Ihre Antwort! **(2 Punkte)**

104 Aufgabe

(11 Punkte)

In einem Experiment zur Steigerung der Milchleistung (*gain*) von Kühen wurden zwei Arten von Musik in den Ställen gespielt. Zum einen ruhige Musik (*calm*) und eher flotte Musik (*pop*). Die Messungen wurden an jeder Kuh (*subject*) wiederholt durchgeführt. Darüber hinaus wurden verschiedene Ställe (*barn*) mit der Musik bespielt.

```
## Linear mixed model fit by REML ['lmerMod']
## Formula: gain ~ attitude + (1 | subject) + (1 | barn)
## Data: data_tbl
##
## REML criterion at convergence: 794.4
##
## Scaled residuals:
##      Min       1Q   Median       3Q      Max
## -2.43600 -0.56787 -0.10670  0.56216  3.30631
##
## Random effects:
## Groups   Name      Variance Std.Dev.
## barn     (Intercept) 242.39   15.569
## subject  (Intercept) 4021.22  63.413
## Residual                    649.76  25.490
## Number of obs: 83, groups:  barn, 7; subject, 6
##
## Fixed effects:
##              Estimate Std. Error t value
## (Intercept) 201.7662    26.8384  7.5178
## attitudepop -17.8009     5.6009 -3.1782
##
## Correlation of Fixed Effects:
##              (Intr)
## attitudepop -0.103
```

1. Ist die Annahme der Normalverteilung an das Outcome *gain* erfüllt? Begründen Sie Ihre Antwort! **(2 Punkte)**
2. Wie groß ist der Effekt der Musikart *attitude*? Liegt ein signifikanter Effekt vor? Schätzen Sie den p-Wert mit einem kritischen t-Wert von $T_k = 1.96$ ab. Begründen und visualisieren Sie Ihre Antwort und Entscheidung! **(3 Punkte)**
3. Was ist der Unterschied zwischen einem „random“ und „fixed“ Effekt. Gehen Sie in der Begründung Ihrer Antwort auf dieses konkrete Beispiel ein! **(3 Punkte)**
4. Wie groß ist die Varianz, die durch die zufälligen Effekte erklärt wird? **(1 Punkt)**
5. Schreiben Sie das Ergebnis der  Ausgabe in einen Satz nieder, der die Information zum Effekt und der Signifikanz enthält! **(2 Punkte)**

Nicht parametrische Tests

Mehr Informationen zu den Aufgaben in den folgenden Kapiteln aus dem Skript Bio Data Science.

- Kapitel 25 - Der Wilcoxon-Mann-Whitney-Test
- Kapitel 26 - Der Kruskal-Wallis-Test

105 Aufgabe

(12 Punkte)

Die Anzahl an Nematoden wurde vor und nach einer Behandlung mit einem bioaktiven Dünger gezählt. Es ergibt sich folgende Datentabelle.

Vorher	Nachher	Differenz	Vorzeichen	Rang	Positiv Rang	Negativ Rang
12	13					
12	12					
9	8					
9	12					
8	9					
7	13					
10	9					
8	12					
9	15					
10	12					
10	11					
12	13					
9	13					
10	15					
14	12					

1. Ergänzen Sie die obige Tabelle mit den notwendigen Informationen, die Sie benötigen um einen Wilcoxon-Vorzeichen-Rang-Test zu rechnen! **(4 Punkte)**
2. Bestimmen Sie die Teststatistik W mit $W = \min(T_-; T_+)$ und berechnen Sie den erwarteten Wert $\mu_W = \frac{n \cdot (n + 1)}{4}$! **(2 Punkte)**
3. Berechnen Sie anschließend den z-Wert mit $z = \frac{W - \mu_W}{\sigma_W}$! **(2 Punkte)**
4. Liegt mit einer Signifikanzschwelle von $z = 1.96$ ein Unterschied zwischen den beiden Zeitpunkten vor? Begründen Sie Ihre Antwort! **(2 Punkte)**
5. Berechnen Sie die Effektstärke mit $r = \left| \frac{z}{\sqrt{n}} \right|$ und interpretieren Sie die Effektstärke! **(2 Punkte)**

106 Aufgabe

(8 Punkte)

Nach einer Behandlung mit RootsGoneX wurde die mittlere Anzahl an Wurzeln an der invasiven Lupine (*Lupinus polyphyllus*) gezählt. Es ergab sich folgender Datensatz an mittleren Wurzelanzahl.

Treatment	Count
RootsGoneX	11.2
RootsGoneX	9.4
RootsGoneX	5.2
RootsGoneX	10.1
RootsGoneX	9.9
RootsGoneX	8.1
RootsGoneX	11.4
Kontrolle	13.6
Kontrolle	6.1
Kontrolle	15.2
Kontrolle	14.8
Kontrolle	12.3
Kontrolle	16.2
Kontrolle	10.7

Rechnen Sie einen Mann-Whitney-U-Test auf den obigen Daten.

1. Bestimmen Sie hierfür U_c mit $U_c = n_1 n_2 + \frac{n_1(n_1 + 1)}{2} - R_1$! (4 Punkte)

2. Geben Sie eine Aussage über die Signifikanz von U_c durch $z = \frac{U_c - \frac{n_1 n_2}{2}}{\sqrt{\frac{n_1 n_2 (n_1 + n_2 + 1)}{12}}}$ und dem kritischen Wert von $z = 1.96$. Begründen Sie Ihre Antwort! (2 Punkte)

3. Berechnen Sie die Effektstärke mit $r = |\frac{z}{\sqrt{n}}|$ und interpretieren Sie die Effektstärke! (2 Punkte)

107 Aufgabe

(9 Punkte)

Die Anzahl an Blüten der Vanillepflanze pro Box wurde nach der Gabe von zusätzlichen Phosphorlösung (Kontrolle, Dosis 20 und Dosis 40) bestimmt. Es ergeben sich folgende nach der Anzahl der Blüten geordnete Daten.

Treatment	Count	Rang Kontrolle	Rang Dosis 20	Rang Dosis 40
Dosis 20	3.6			
Kontrolle	4.3			
Dosis 20	5.3			
Dosis 40	7.6			
Dosis 40	8.0			
Kontrolle	9.1			
Kontrolle	9.1			
Dosis 40	10.9			
Dosis 20	11.2			
Kontrolle	11.3			
Kontrolle	12.0			
Dosis 40	12.5			
Dosis 20	13.1			
Kontrolle	13.2			
Dosis 20	14.4			
Dosis 20	14.6			
Kontrolle	17.5			
Dosis 40	17.8			

Rechnen Sie einen Kruskal-Wallis-Test auf den obigen Daten.

- Bestimmen Sie hierfür H_c mit $H_c = \frac{12}{n(n+1)} \left(\frac{R_1^2}{n_1} + \frac{R_2^2}{n_2} + \frac{R_3^2}{n_3} \right) - 3(n+1)!$ (6 Punkte)
- Geben Sie eine Aussage über die Signifikanz von H_c durch den kritischen Wert von $H = 5.99!$ (1 Punkt)
- Wie lautet die Nullhypothese die Sie mit dem Kruskal-Wallis-Test überprüfen? (1 Punkt)
- Was sagt ein signifikantes Ergebnis des Kruskal-Wallis-Test in Bezug auf die einzelnen Gruppenvergleiche aus? (1 Punkt)

Multiple Gruppenvergleiche

Mehr Informationen zu den Aufgaben in den folgenden Kapiteln aus dem Skript Bio Data Science.

- Kapitel 20.3 Adjustierung für multiple Vergleiche

108 Aufgabe

(6 Punkte)


In einem Experiment zur Dosiswirkung wurden verschiedene Dosisstufen mit einer Kontrollgruppe verglichen. Es wurden vier t-Test für den Mittelwertsvergleich gerechnet und es ergab sich folgende Tabelle mit den rohen p-Werten.

Vergleich	Raw p-val	Adjusted p-val	Reject H_0
dose 10 - ctrl	0.030		
dose 15 - ctrl	0.012		
dose 20 - ctrl	0.020		
dose 40 - ctrl	0.060		

1. Füllen Sie die Spalte „adjustierte p-Werte“ mit den adjustierten p-Werten nach Bonferroni aus! **(2 Punkte)**
2. Entscheiden Sie, ob nach der Adjustierung die Nullhypothese weiter abgelehnt werden kann. Tragen Sie Ihre Entscheidung in die obige Tabelle ein. Begründen Sie Ihre Antwort! **(2 Punkte)**
3. Erklären Sie warum die p-Werte bei multiplen Vergleichen adjustiert werden müssen? **(2 Punkte)**

109 Aufgabe

(9 Punkte)

In einem Experiment mit vier Dosisstufen (A, B, C und D) erhalten Sie folgende Matrix als  Ausgabe mit den rohen, unadjustierten p -Werten.

```
##           A           B           C           D
## A 1.00000000 0.00000001 0.00000005 0.00000000
## B 0.00000001 1.00000000 0.6354730 0.0123340
## C 0.00000005 0.6354730 1.00000000 0.0037075
## D 0.00000000 0.0123340 0.0037075 1.00000000
```

Im Weiteren erhalten Sie folgende Informationen über die Fallzahl n , den Mittelwert $mean$ und die Standardabweichung sd in den jeweiligen Dosisstufen.

trt	n	mean	sd
A	9	12.59	1.40
B	9	6.33	2.59
C	9	6.77	1.80
D	9	3.88	1.86

1. Zeichnen Sie in eine Abbildung, die sich ergebenden Barplots! **(2 Punkte)**
2. Adjustieren Sie die rohen p -Werte nach Bonferroni. Begründen Sie Ihre Antwort! **(3 Punkte)**
3. Ergänzen Sie das *Compact letter display (CLD)* zu der Abbildung! **(2 Punkte)**
4. Interpretieren Sie das *Compact letter display (CLD)*! **(2 Punkte)**