

Biostatistik, WS 2010/2011

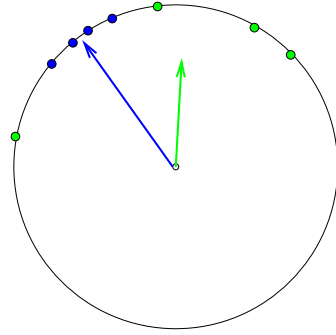
Der t-Test

(Zwei gepaarte Stichproben bzw. eine Stichprobe)

Matthias Birkner

<http://www.mathematik.uni-mainz.de/~birkner/Biostatistik1011/>

10.12.2010



Pfeilspitze: Schwerpunkt der Austrittspunkte bei grünem Licht.
Dasselbe für die "blauen" Austrittspunkte.

Je variabler die Richtungen, desto kürzer der Pfeil!

Fragestellung

Hat die Farbe der monochromatischen Beleuchtung einen Einfluss auf die Orientierung?
 Experiment: Bei 17 Vögeln wurde die Länge des Schwerpunktsvektors sowohl bei **blauem** als auch bei **grünem** Licht bestimmt.

Wir verwenden im Folgenden simulierte Daten, die sich an der Literatur orientieren.



Witischko, W.; Gesson, M.; Stapput, K.; Witischko, R.

Light-dependent magnetoreception in birds: interaction of at least two different receptors. *Naturwissenschaften* 91.3, pp. 130-4, 2004.



Witischko, R.; Ritz, T.; Stapput, K.; Thalau, P.; Witischko, W.

Two different types of light-dependent responses to magnetic fields in birds. *Curr Biol* 15.16, pp. 1518-23, 2005.

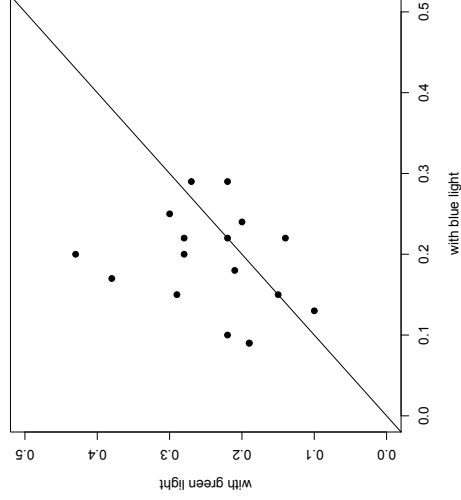


Witischko, R.; Stapput, K.; Bischof, H. J.; Witischko, W.

Light-dependent magnetoreception in birds: increasing intensity of monochromatic light changes the nature of the response. *Front Zool*, 4, 2007.

6/59

Trauerschnäpper:
 Länge des Schwerpunktsvektors
 bei grünem und bei blauem Licht, $n=17$



7/59

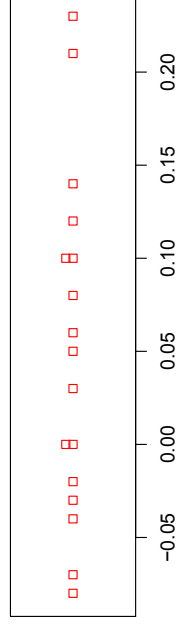
Wie kann ich
statistisch testen,
ob die Farbe
einen Einfluss hat?

8/59

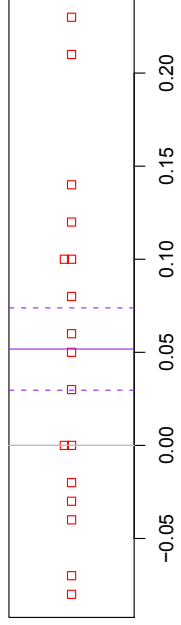
Wir berechnen nun für jeden Vogel den Abstand des Punktes
von der Diagonale,

d.h.

$$x := \text{“Grünwert”} - \text{“Blauwert”}$$



9/59



Kann der wahre Mittelwert $\mu = 0$ sein?

$$\bar{x} = 0.0518$$

$$s = 0.0912$$

$$\text{Standardfehler} = \frac{s}{\sqrt{n}} = \frac{0.0912}{\sqrt{17}} = 0.022$$

10/59

Ist $|\bar{x} - \mu| \approx 0.0518$ eine große Abweichung?

Groß? Groß im Vergleich zu was?

In welcher Vergleichseinheit soll $|\bar{x} - \mu|$ gemessen werden?

Im Vergleich zum
Standardfehler!

$|\bar{x} - \mu|$
gemessen in der Einheit 'Standardfehler'
heißt **t-Statistik**

$$t := \frac{\bar{x} - \mu}{s/\sqrt{n}}$$

11/59

$$t := \frac{\bar{X} - \mu}{s/\sqrt{n}}$$

$t = 1$ bedeutet

1 Standardfehler von μ entfernt
(kommt häufig vor)

$t = 3$ bedeutet

3 Standardfehler von μ entfernt
(kommt selten vor)

In unserem Fall:

$$t = \frac{\bar{X} - \mu}{s/\sqrt{n}} \approx \frac{0.0518}{0.022} \approx 2.34$$

Also: \bar{x} ist mehr als 2.3 Standardfehler von $\mu = 0$ entfernt.

Wie wahrscheinlich ist das, wenn 0 der wahre Mittelwert ist?
anders gefragt:

Ist diese Abweichung signifikant?

Für die Antwort benötigen wir die Verteilung der t-Statistik.

Wir wissen:

$$\frac{\bar{X} - \mu}{\sigma/\sqrt{n}}$$

ist (für große n) ungefähr standardnormalverteilt (dies ist die Grundlage des z-Tests).

Aber:

Die t-Statistik ist jedoch mit s an Stelle von σ definiert (und nicht normalverteilt), die Approximation mit der Normalverteilung ist (speziell für kleine und „mittelgroße“ Stichprobengrößen n) häufig zu grob.

15/59

Allgemein gilt

Sind X_1, \dots, X_n unabhängig aus einer Normalverteilung mit Mittelwert μ gezogen (und beliebiger Varianz $\sigma^2 > 0$), so ist

$$\frac{\bar{X} - \mu}{s/\sqrt{n}}$$

t-verteilt mit $n - 1$ Freiheitsgraden (df=*degrees of freedom*).

Eine t-verteilte Zufallsvariable bezeichnen wir meist mit T .

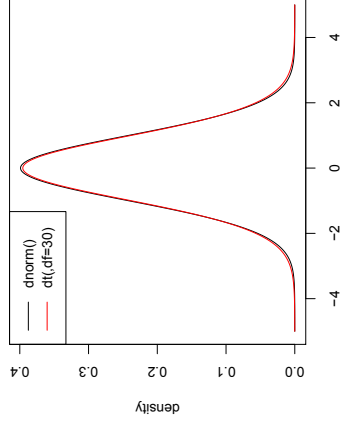
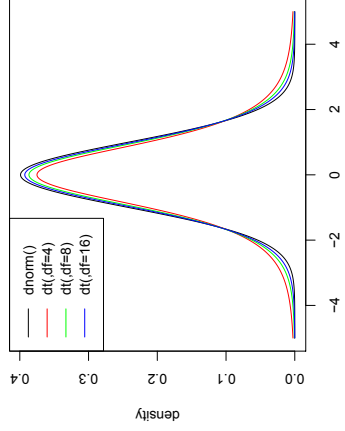
Die t-Verteilung heißt auch **Student-Verteilung**. Sie wurde von William Gosset erforscht und von ihm 1908 veröffentlicht, während er in einer Guinness-Brauerei arbeitete. Da sein Arbeitgeber die Veröffentlichung nicht gestattete, veröffentlichte Gosset sie unter dem Pseudonym *Student*.



W.S. Gossett, 1876–1937

16/59

Dichte der t-Verteilung



17/59

Freiheitsgrade?

Beispiel: Es gibt 5 Freiheitsgrade im Vektor

$$x = (x_1, x_2, x_3, x_4, x_5)$$

da 5 Werte frei wählbar sind. Der Vektor

$$v := x - \bar{x} := (x_1 - \bar{x}, x_2 - \bar{x}, x_3 - \bar{x}, x_4 - \bar{x}, x_5 - \bar{x})$$

hat 4 Freiheitsgrade, denn nach Wahl von v_1, v_2, v_3, v_4 ist v_5 festgelegt wegen $v_1 + \dots + v_4 + v_5 = 0$.

Die Bezeichnung „Student-Verteilung mit $n - 1$ Freiheitsgraden“ ist motiviert durch die Tatsache, dass $t = (\bar{x} - \mu) / (s / \sqrt{n})$, wo

$$s = \sqrt{s^2} = \frac{1}{\sqrt{n-1}} \left((x_1 - \bar{x})^2 + \dots + (x_n - \bar{x})^2 \right)^{1/2}$$

(bis auf Normierung) die Länge eines „Vektors mit $n - 1$ Freiheitsgraden“ ist.

18/59

Wir meinen:

Die Farbe der Beleuchtung
hat einen Einfluss auf die Orientierung

Ein Skeptiker würde erwidern:

Alles nur Zufall

Wir wollen nun zeigen:

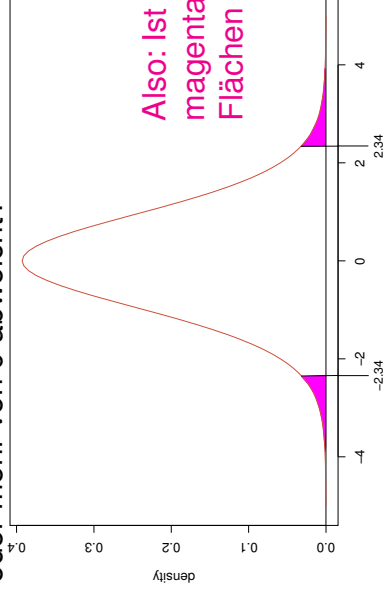
Unter der Annahme 'Kein Einfluss'
ist die Beobachtung sehr unwahrscheinlich

Nullhypothese: $\mu = 0$

20/59

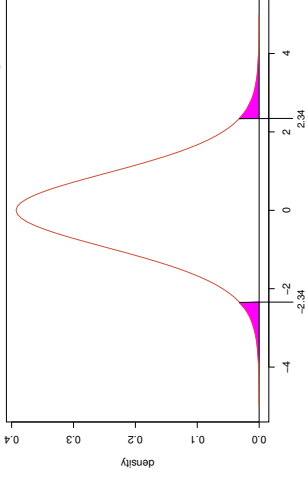
Wie (un)wahrscheinlich ist nun eine **mindestens**
so große Abweichung wie 2.34 Standardfehler?
Idee des statistischen Testens: Wir legen ein Signifikanzniveau
 $\alpha \in (0, 1)$, z.B. $\alpha = 5\%$, fest und fragen:

Wenn das wahre $\mu = 0$ ist (d.h. unter der Nullhypothese), wie
wahrscheinlich ist es, einen t -Wert zu beobachten, der um 2,34
oder mehr von 0 abweicht?



21/59

Beobachtet: $t = 2,34$
Ist der Gesamthalt der magentafarbenen Flächen $\leq \alpha$?



Das 97,5%-Quantil der Student-Verteilung mit 16 Freiheitsgraden ist ca. $2,12 \leq 2,34$, also: Ja.

Genauer: Es gilt $\mathbb{P}(|T| \geq 2,34) \approx 0,0325 (\leq 0,05)$.

Man nennt $\mathbb{P}(|T| \geq t)$ den p -Wert (bei beobachtetem Wert t).

Wir halten fest:

$$p - \text{Wert} = 0,03254 (\leq 0,05)$$

Wenn die **Nullhypothese** “alles nur Zufall” (hier $\mu = 0$) gilt, dann ist eine mindestens so große Abweichung sehr unwahrscheinlich.

Sprechweise:

Wir verwerfen die Nullhypothese auf dem 5%-Signifikanzniveau.

Oder:

Die Differenz zwischen grün und blau ist auf dem 5%-Niveau signifikant.

Die Nullhypothese wurde also auf dem 5%-Niveau verworfen.
Welche Aussagen sind wahr/sinnvoll?

- Die Nullhypothese ist falsch. **Die Nullhypothese ist falsch.**
- Die Nullhypothese ist mit 95%-iger Ws falsch.
Die Nullhypothese ist mit 95%-iger Ws falsch.
- Falls die Nullhypothese wahr ist, beobachtet man ein so extremes Ergebnis nur in höchstens 5% der Fälle. **Falls die Nullhypothese wahr ist, beobachtet man ein so extremes Ergebnis nur in höchstens 5% der Fälle. ✓**
- Die Orientierung der Vögel ist bei blau und grün verschieden.
Die Orientierung der Vögel ist bei blau und grün verschieden.
- Die Orientierung bei grün und blau war in dem Experiment auf dem 5%-Niveau signifikant verschieden. **Die Orientierung bei grün und blau war in dem Experiment auf dem 5%-Niveau signifikant verschieden. ✓**

24/59

Man könnte auch ein anderes Signifikanzniveau α wählen.
Dann müsste man zeigen, dass der p-Wert kleiner als α ist.

Wichtig: Wähle zuerst das Signifikanzniveau und ermittle erst dann den p-Wert! Das Signifikanzniveau je nach p-Wert zu wählen ist geschummelt.

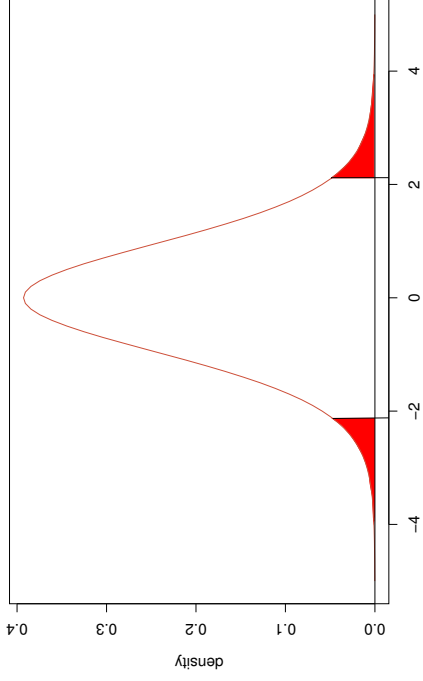
**In der Literatur wird üblicherweise
5% als Signifikanzniveau gewählt.**

Beachte:

**Falls die Nullhypothese zutrifft,
ist die Wahrscheinlichkeit,
dass wir sie zu Unrecht auf dem 5%-Niveau verwerfen,
höchstens 5%.**

25/59

Wir verwerfen also die Nullhypothese auf 5%-Niveau, wenn der Wert der t -Statistik in den roten Bereich fällt:



(hier am Beispiel der t -Verteilung mit $df=16$ Freiheitsgraden)

26/59

Welche t -Werte sind auf dem 5%-Niveau signifikant?

Anzahl Freiheitsgrade	$ t \geq \dots$
5	2.57
10	2.23
20	2.09
30	2.04
100	1.98
" ∞ "	1.96

Diese sog. kritischen Werte (für das zweiseitige 5%-Signifikanzniveau beim t -Test) sind die 97,5%-Quantile der t -Verteilung mit der jeweils angegebenen Anzahl Freiheitsgrade.

In der Praxis entnimmt man sie einem Statistik-Computerprogramm oder einer Tabelle.

27/59

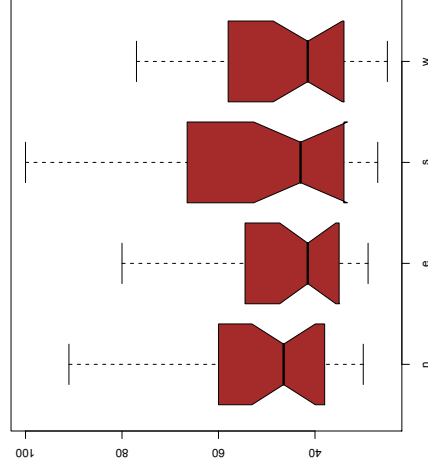
Bei $n = 28$ Bäumen wurden die Korkdicken [mm] in den vier Himmelsrichtungen gemessen:

	n	o	s	w
	72	66	76	77
	60	53	66	63
	5	57	64	58
	41	29	36	38
	32	32	35	36
	30	35	34	26
	39	39	31	27
	•	•	•	•
	•	•	•	•

(Wir verwenden wieder simulierte Daten, die aber Daten aus echten Studien nachempfunden sind, auch im Ergebnis.)

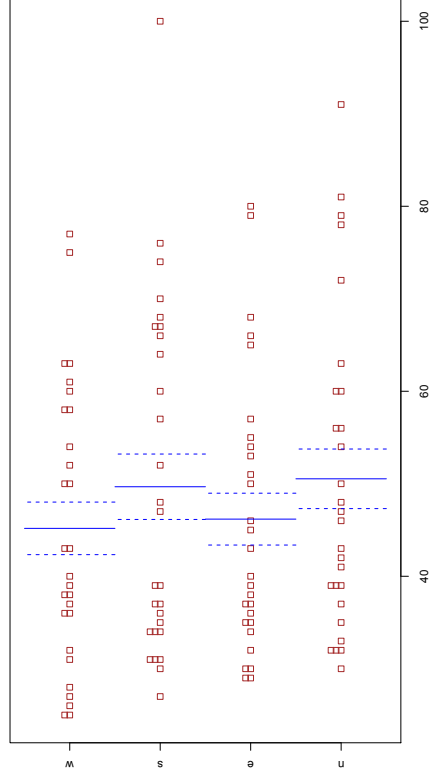
31/59

Korkdicken nach Himmelsrichtung getrennt



Kann da was signifikant unterschiedlich sein???

32/59



Stripchart der Korkdicken je nach Himmelsrichtung und
Mittelwerten \pm Standardfehler

Kann da was signifikant unterschiedlich sein???

33/59

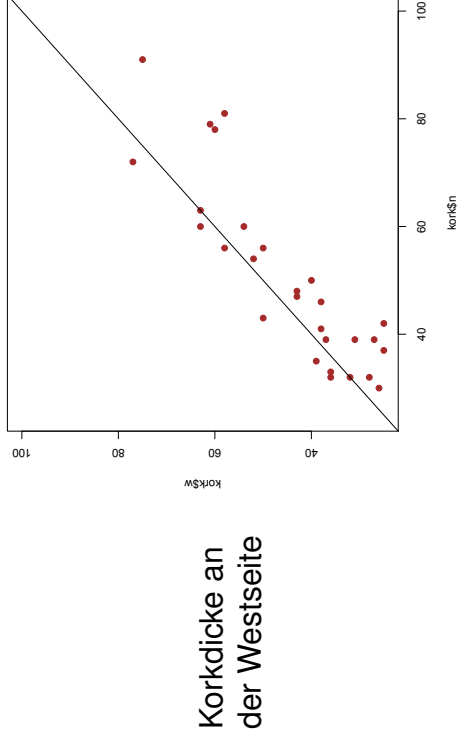
Haben wir irgend etwas übersehen?

Wir haben bisher vernachlässigt,
welche Werte von demselben Baum kommen!
Die Bäume unterscheiden sich sehr in ihrer Größe und Dicke.

Vergleiche also jeweils Paare von Korkdicken,
die von demselben Baum kommen!
(\rightsquigarrow gepaarter t-Test)

34/59

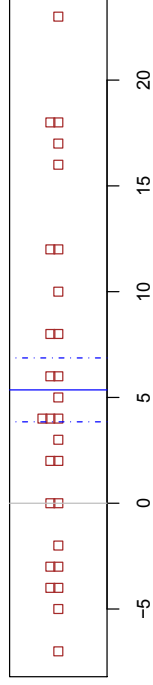
Korkdicken [mm] bei $n = 28$ Bäumen



Korkdicke an der Nordseite

35/59

Differenz der Korkdicken an der Nord- und der Westseite für
 $n = 28$ Bäume



und $\text{Mittelwert} \pm \text{Standardfehler}$
Ist die Differenz signifikant von 0
verschieden?

36/59

$x :=$ (Korkdicke Nordseite) – (Korkdicke Westseite)

$$\bar{x} \approx 5,36$$

$$s_x \approx 7,99$$

$$\frac{s_x}{\sqrt{n}} \approx 1,51$$

$$t\text{-Wert} = \frac{\bar{x}}{s_x/\sqrt{n}} \approx 3,547$$

$$\text{Anzahl Freiheitsgrade: } df = n - 1 = 27$$

Das 97,5%-Quantil der Student- t -Verteilung mit 27 Freiheitsgraden ist $\approx 2,05$,

$$|3,547| > 2,05.$$

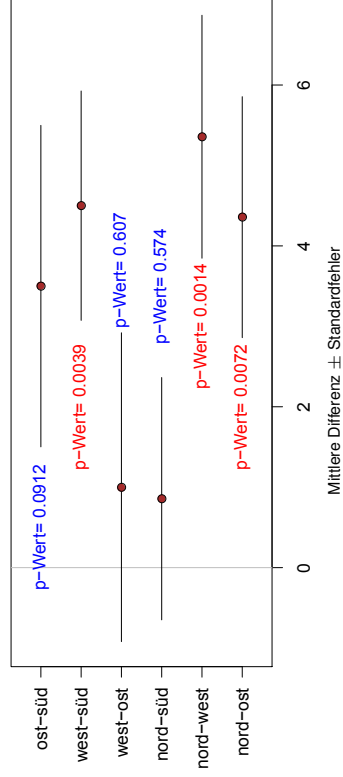
Die Differenz der Korkdicke an Nord- und Westseite ist auf dem 5%-Niveau signifikant verschieden.

37/59

Bemerkung: Die 6 möglichen Paarungen im Kork-Beispiel

Paarung	N-O	N-W	N-S	W-O	W-S	O-S
Mittlere Differenz \bar{x}	4,36	5,36	0,86	1,00	4,50	3,50
Standardfehler s/\sqrt{n}	1,50	1,51	1,51	1,92	1,43	2,00
$t = \frac{\bar{x}-0}{s/\sqrt{n}}$	2,91	3,55	0,57	0,52	3,16	1,75
$ t \geq 2,052?$	ja	ja	nein	nein	ja	nein

(Das 97,5%-Quantil der Student-Verteilung mit 27 Freiheitsgraden ist 2,052.)



38/59

Zusammenfassung gepaarter t-Test

Gegeben: gepaarte Beobachtungen

$$(Y_1, Z_1), (Y_2, Z_2), \dots, (Y_n, Z_n)$$

Nullhypothese $H_0: \mu_Y = \mu_Z$

Signifikanzniveau: α (meist $\alpha = 5\%$)

Test: **gepaarter t-Test** (genauer: zweiseitiger gepaarter t-Test)

Berechne Differenz $X := Y - Z$

Berechne Teststatistik

$$t := \frac{\bar{X}}{s(X)/\sqrt{n}}$$

Verwirf Nullhypothese, falls $|t| \geq (1 - \alpha/2)$ -Quantil der t -Verteilung mit $n - 1$ Freiheitsgraden

(d.h. wenn der p-Wert $\mathbb{P}(|T_{n-1}| \geq |t|) \leq \alpha$ ist)

40/59

Zusammenfassung Ein-Stichproben t-Test

Gegeben: Beobachtungen

$$X_1, X_2, \dots, X_n$$

Nullhypothese $H_0: \mu_X = c$ (Den Wert c kennt man, oft $c = 0$)

Signifikanzniveau: α (meist $\alpha = 5\%$)

Test: **t-Test**

Berechne Teststatistik

$$t := \frac{\bar{X} - c}{s(X)/\sqrt{n}}$$

Verwirf Nullhypothese, falls $|t| \geq (1 - \alpha/2)$ -Quantil der t -Verteilung mit $n - 1$ Freiheitsgraden

(d.h. wenn der p-Wert $\mathbb{P}(|T_{n-1}| \geq |t|) \leq \alpha$ ist)

41/59

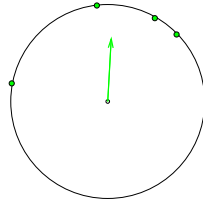
Beispiel: Codon Bias

- Wir beobachten 101844 mal CCT und 106159 mal CCA
- Wenn beide eigentlich gleich wahrscheinlich sind, erwarten wir 104001.5 von jedem.
- Die Beobachtung weicht um 2156 von diesem Erwartungswert ab
- z-Test: Die Wahrscheinlichkeit einer mindestens so großen Abweichung ist kleiner als 10^{-20}
- Also sind CCT und CCA wohl nicht gleich wahrscheinlich.

43/59

Beispiel: Zugvogelorientierung

- Wie variabel ist die Abflugrichtung bei grünem und bei blauem Licht.
- Wir messen die Variabilität durch die Länge des Schwerpunktsvektors.



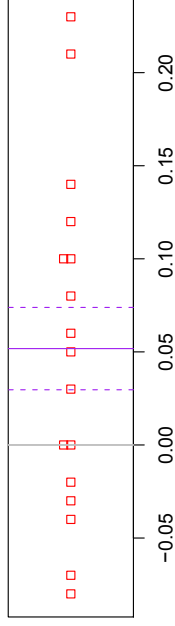
- Quantifiziere Unterschied durch $X = (\text{Länge grün}) - (\text{Länge blau})$.
- Wenn das Licht keinen Einfluss hat, gilt $\mathbb{E}X = 0$.

44/59

Beispiel: Zugvogelorientierung

X =(Länge grün)– (Länge blau)

- Wenn das Licht keinen Einfluss hat, gilt $\mathbb{E}X = 0$.
- Wir beobachten aber $\bar{X} = 0.0518$ und $SEM=0.022$



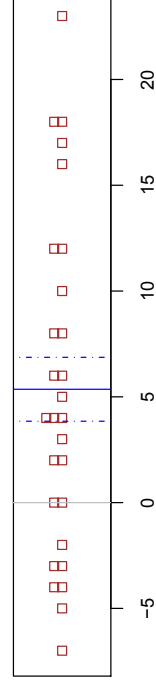
- t -Test: p -Wert dieser Abweichung ist ca. 3.3%.
- Vermutlich hat die Lichtfarbe also doch einen Einfluss

45/59

Beispiel: Dicke des Korks

X =(Korkdicke an der Nordseite)– (Korkdicke an der Westseite)

- Wenn die Seite keine Rolle spielt, ist $\mathbb{E}X = 0$.
- Wir sehen aber $\bar{X} = 5.36$ und $SEM= 1.51$



- t -Test: p -Wert dieser Abweichung ist ca. 0.14%.
- Also hat die Himmelsrichtung wohl doch einen Einfluss.

46/59

Prinzip des statistischen Testens

- Wir möchten belegen, dass eine Abweichung in den Daten vermutlich nicht allein auf Zufallsschwankung beruht.
- Dazu spezifizieren wir zunächst eine **Nullhypothese H_0** , d.h. wir konkretisieren, was "allein auf Zufall beruhen" bedeutet.
- Dann versuchen wir zu zeigen: Wenn H_0 gilt, dann sind **Abweichungen**, die mindestens so groß sind wie die beobachtete, sehr unwahrscheinlich.
- Wenn uns das gelingt, verwerfen wir H_0 .
- Was wir als **Abweichung** auffassen, sollte klar sein, bevor wir die Daten sehen.

47/59

Nullhypothesen

- H_0 bei Codon-Bias: CCT und CCA haben jeweils W'keit $\frac{1}{2}$
Außerdem: alle Positionen entscheiden unabhängig zwischen CCT und CCA
- H_0 bei Vogelorientierung und Korkdicken: $\mathbb{E}X = 0$.
Außerdem: X normalverteilt, X_i unabhängig.

48/59

Abweichungen und p -Werte

- Codon Bias: Anzahl CCT weicht um 2156 vom Mittelwert ab.
Wegen der Binomialverteilungsannahme gehen wir von festem σ aus und berechnen mit dem z-Test den p -Wert: Die Wahrscheinlichkeit, dass eine $\text{bin}(n, \frac{1}{2})$ -verteilte Zufallsgröße um mindestens 2156 von $n/2$ abweicht.
- Vogelorientierung und Korkdicke:

$$t\text{-Wert} = \frac{\bar{X}}{s/\sqrt{n}}$$

p -Wert: W'keit, dass t -Wert bei $n - 1$ mindestens so stark von 0 abweicht wie beobachtet.

49/59

Zweiseitig oder einseitig testen?

In den meisten Fällen will man testen, ob zwei Stichproben sich signifikant unterscheiden.

↪ **zweiseitiger Test**

In manchen Fällen

- kann man von vornherein ausschließen, dass die erste Stichprobe kleinere Werte als die zweite Stichprobe hat. Dann will man testen, ob die erste Stichprobe signifikant größer ist.
- will man nur testen, ob die erste Stichprobe signifikant größer ist.
- will man nur testen, ob die erste Stichprobe signifikant kleiner ist.

↪ **einseitiger Test**

50/59

Beispiel für einseitigen Test:

Man will zeigen,
dass ein Wachstumshormon wirkt,
also kein Placebo ist.

Dazu müssen die Größen Y in der behandelten Gruppe
signifikant größer sein
als die Größen Z in der Kontrollgruppe.
Die zu entkräftende Nullhypothese wäre hier:

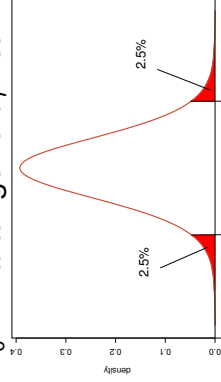
$$\text{Nullhypothese } \mu_Y \leq \mu_Z$$

Definiere die Differenz $X := Y - Z$.

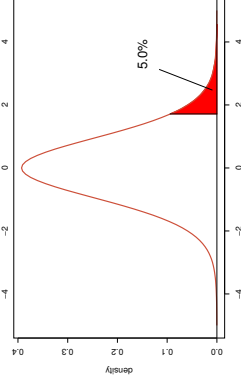
51/59

Zweiseitig oder einseitig testen?

Wir beobachten einen Wert x , der deutlich größer als der H_0 -Erwartungswert μ ist.



$$p\text{-Wert} = \mathbb{P}_{H_0}(|X - \mu| \geq |x - \mu|)$$



$$p\text{-Wert} = \mathbb{P}_{H_0}(X - \mu \geq x - \mu)$$

52/59

Reine Lehre des statistischen Testens

- Formuliere eine **Nullhypothese** H_0 , z.B. $\mu = 0$.
- Lege ein **Signifikanzniveau** α fest; üblich ist $\alpha = 0.05$.
- Lege ein Ereignis \mathcal{A} fest, so dass

$$\Pr_{H_0}(\mathcal{A}) = \alpha$$

(oder zumindest $\mathbb{P}_{H_0}(\mathcal{A}) \leq \alpha$).

z.B. $\mathcal{A} = \{\bar{X} > q\}$ oder $\mathcal{A} = \{|\bar{X} - \mu| > r\}$

- **ERST DANN:** Betrachte die Daten und überprüfe, **ob** \mathcal{A} eintritt.
- Dann ist die Wahrscheinlichkeit, dass H_0 verworfen wird, wenn H_0 eigentlich richtig ist ("**Fehler erster Art**") , lediglich α .

53/59

Verstöße gegen die reine Lehre

“Beim zweiseitigen Testen kam ein p -Wert von 0.06 raus. Also hab ich einseitig getestet, da hat's dann funktioniert.”

genauso problematisch:

“Beim ersten Blick auf die Daten habe ich sofort gesehen, dass \bar{x} größer ist als μ_{H_0} . Also habe ich gleich einseitig getestet”

54/59

Wichtig

Die Entscheidung, ob einseitig oder zweiseitig getestet wird, darf nicht von den konkreten Daten abhängen, die zum Test verwendet werden.

Allgemeiner: Ist \mathcal{A} das Ereignis, dass zum Verwerfen von H_0 führt (falls es eintritt), so muss die Festlegung von H_0 stattfinden bevor man die Daten betrachtet hat.

Die **Wahl von \mathcal{A}** sollte von der **Alternative H_1** abhängen, also davon, was wir eigentlich zeigen wollen, indem wir H_0 durch einen Test verwerfen. Es muss gelten:

$$\mathbb{P}_{H_0}(\mathcal{A}) \leq \alpha$$

und

$$\mathbb{P}_{H_1}(\mathcal{A}) = \text{möglichst groß,}$$

damit die **W'keit eines Fehlers zweiter Art**, dass also H_0 nicht verworfen wird, obwohl H_1 zutrifft, möglichst klein ist.

Beispiele

- Wenn wir von Anfang an unsere Vermutung belegen wollten, dass sich die Trauerschnäpper bei grünem Licht stärker auf eine Richtung konzentrieren als bei blauem, dürfen wir einseitig testen.
- Wenn dann aber noch so deutlich herauskommt, dass die Richtungswahl bei blauem Licht enger konzentriert war, so ist das dann strenggenommen nicht als signifikant zu betrachten.
- Wenn wir von Anfang an die Vermutung belegen wollten, dass der Kork an der Nordseite des Baumes dicker war, dürfen wir einseitig testen.
- Wenn dann aber noch so deutlich herauskommt, dass der Kork im Westen dicker ist, ist das strenggenommen nicht signifikant.

57/59

Angenommen, H_0 wird auf dem 5%-Niveau verworfen. Welche Aussage gilt dann?

- Die Nullhypothese ist falsch. **Die Nullhypothese ist falsch.**
- H_0 ist mit 95%-iger Wahrscheinlichkeit falsch.
 H_0 ist mit 95%-iger Wahrscheinlichkeit falsch.
- Falls die Nullhypothese wahr ist, beobachtet man ein so extremes Ergebnis nur in 5% der Fälle. **Falls die Nullhypothese wahr ist, beobachtet man ein so extremes Ergebnis nur in 5% der Fälle. ✓**

58/59

Angenommen, H_0 konnte durch den Test nicht verworfen werden. Welche Aussagen sind dann richtig?

- Wir müssen die Alternative H_1 verwerfen.
~~Wir müssen die Alternative H_1 verwerfen.~~
- H_0 ist wahr. ~~H_0 ist wahr.~~
- H_0 ist wahrscheinlich wahr. ~~H_0 ist wahrscheinlich wahr.~~
- Es ist ungefährlich, davon auszugehen, dass H_0 zutrifft.
~~Es ist ungefährlich, davon auszugehen, dass H_0 zutrifft.~~
- Auch wenn H_0 zutrifft, ist es nicht sehr unwahrscheinlich, dass unsere Teststatistik einen so extrem erscheinenden Wert annimmt. ~~Auch wenn H_0 zutrifft, ist es nicht sehr unwahrscheinlich, dass unsere Teststatistik einen so extrem erscheinenden Wert annimmt.~~
- Die Nullhypothese ist in dieser Hinsicht mit den Daten verträglich. ~~Die Nullhypothese ist in dieser Hinsicht mit den Daten verträglich.~~