

Wahrscheinlichkeitsrechnung und Statistik für Biologen **Versuchsplanung**

Noémie Becker & Dirk Metzler

Sommersemester 2014

Inhaltsverzeichnis

1	Warnung	1
2	Stichprobenlänge	1
2.1	Allgemeines	1
2.2	Einstichproben-Tests	2
2.3	Zweistichproben-Test	5
2.4	Einseitige Tests	6
2.5	Übersicht	7
2.6	Stichprobenlänge ermitteln mit R	8
2.7	F-Test	10
3	Stichprobenwahl	11
3.1	Überspitzte Beispiele	11
3.2	Zufallsstichprobe	12
3.3	Elimination von nicht-interessierenden Einflussgrößen	13
3.4	Blockbildung	14
3.5	Balanced Design vs Non-Balanced Design	15
3.6	Randomisierung	15

1 Warnung

Warnung

Für eine wissenschaftliche Publikation braucht man:

- Signifikanz (\rightsquigarrow Stichprobenlänge groß genug?)
- Geeignete Auswahl der Stichprobe (\rightsquigarrow Randomisierung)

Dies muss bei der [Versuchsplanung](#) beachtet werden!

Warnung

Erst denken, dann arbeiten!

Sonst kann wochen-/monatelange Laborarbeit vergebens sein.

Bei der Versuchsplanung (**BEVOR** man die Daten generiert) muss man u.a. folgende Fragen sinnvoll beantworten:

- „Wie groß muss die Stichprobe sein?“

- „An welchen Versuchsobjekten wird welche Methode angewendet?“ bzw. „Wie wird die Stichprobe gesampelt?“

Um diese Fragen sinnvoll beantworten zu können, muss man sich die statistische Auswertung überlegen, **BEVOR** man die Daten generiert.

2 Stichprobenlänge

2.1 Allgemeines

Allgemeines

Je größer die Stichprobenlänge ist,

- desto wahrscheinlicher wird ein vorhandener Unterschied durch einen statistischen Test angezeigt
- desto kleinere Unterschiede können durch statistische Tests erkannt werden
- desto teurer wird der Versuch.

Es ist also wichtig, eine geeignete Stichprobenlänge zu wählen. Dazu muss man sich überlegen,

- welcher Unterschied durch die anzuwendenden Tests erkannt werden soll,
- wie groß die Variabilität in den Daten in etwa sein wird.

Allgemeines

Man benötigt:

- d = Unterschied, den man mindestens erkennen können möchte. (engl: detection level)
- einen ungefähren Wert s für die Standardabweichung, die man in den Daten erwartet (oft ein Wert aus Vorversuchen).
- $\alpha = \Pr_{H_0}(H_0 \text{ wird (fälschlicherweise) abgelehnt})$. Meist 5%. α ist das Signifikanzniveau. Die Ws α heißt auch Fehler 1.Art.
- $\beta = \Pr_{\text{Alternative}}(H_0 \text{ wird (fälschlicherweise) nicht verworfen})$. Die Wahl von β hängt stark vom Problem ab. $1 - \beta$ ist die Testmacht. Die Ws β heißt auch Fehler 2.Art.

2.2 Einstichproben-Tests

Einstichproben-Tests

Frage: Ist der wahre Mittelwert gleich μ_0 ?

Beispiel: Kältestress-Toleranz bei Fruchtfliegen.



photo (c) André Karwath (Bild zeigt eine *Drosophila melanogaster*)

Einstichproben-Tests

Die Chill-Coma Recovery Time (CCRT) ist die Zeit in Sekunden, nach der die Fliege nach einem Kältekoma wieder aufwacht. In früheren Versuchen wurde bei *Drosophila ananassae* aus Bangkok eine mittlere CCRT von 46 gemessen.

Frage: Ist die CCRT bei *Drosophila ananassae* aus Kathmandu (Nepal) verschieden von 46?

Geplanter Test: (zweiseitiger) Einstichproben t-Test.

Ziel: Finde Unterschiede, die größer als $d = 4$ sind. Signifikanzniveau $\alpha = 5\%$. Testmacht $1 - \beta = 80\%$.

Vorwissen: Standardabweichung bei Vortest war $s = 11.9$

Frage: Bei wie vielen Fliegen muss ich die CCRT messen, um das Ziel zu erreichen?

Einstichproben-Tests

Frage: Stichprobenlänge für CCRT-Versuch?

Lösung: Es soll gelten:

$$n \geq \frac{s^2 \cdot (t_{1-\frac{\alpha}{2}, n-1} + t_{1-\beta, n-1})^2}{d^2}$$

wobei $t_{1-\frac{\alpha}{2}, n-1} \leftarrow \text{qt}(1-\alpha/2, n-1)$ das $(1 - \alpha/2)$ -Quantil und $t_{1-\beta, n-1} \leftarrow \text{qt}(1-\beta, n-1)$ das $(1 - \beta)$ -Quantil der t-Verteilung ist.

Leider kann man nicht einfach einsetzen, da die rechte Seite von n abhängt.

Entweder probiert man herum und sucht das kleinste n wofür die Ungleichung gilt.

Einstichproben-Tests

Oder man beginnt mit

$$n_0 = \frac{s^2 \cdot (z_{1-\frac{\alpha}{2}} + z_{1-\beta})^2}{d^2}$$

wobei $z_{1-\frac{\alpha}{2}} \leftarrow \text{qnorm}(1-\alpha/2)$ das $(1 - \alpha/2)$ -Quantil und $z_{1-\beta} \leftarrow \text{qnorm}(1-\beta)$ das $(1 - \beta)$ -Quantil der Normalverteilung ist. Die benötigte Stichprobenlänge findet man dann durch Iteration:

$$n_1 = \frac{s^2 \cdot (t_{1-\frac{\alpha}{2}, n_0-1} + t_{1-\beta, n_0-1})^2}{d^2}$$
$$n_2 = \frac{s^2 \cdot (t_{1-\frac{\alpha}{2}, n_1-1} + t_{1-\beta, n_1-1})^2}{d^2}$$

usw bis sich nichts mehr ändert.

Einstichproben-Tests

Zurück zum Beispiel:

$$n_0 = \frac{s^2 \cdot (z_{1-\frac{\alpha}{2}} + z_{1-\beta})^2}{d^2} = \frac{11.9^2 (z_{0.975} + z_{0.8})^2}{4^2} = 69.48 \sim 70$$
$$n_1 = \frac{s^2 \cdot (t_{1-\frac{\alpha}{2}, n_0-1} + t_{1-\beta, n_0-1})^2}{d^2} = \frac{11.9^2 (t_{0.975, 69} + t_{0.8, 69})^2}{4^2}$$
$$= 71.47 \sim 72$$
$$n_2 = \frac{s^2 \cdot (t_{1-\frac{\alpha}{2}, n_1-1} + t_{1-\beta, n_1-1})^2}{d^2} = \frac{11.9^2 (t_{0.975, 71} + t_{0.8, 71})^2}{4^2}$$
$$= 71.41 \sim 72$$

Antwort: Die Stichprobenlänge für den CCRT-Versuch sollte mindestens $n \geq 72$ sein.

Einstichproben-Tests

Bemerkung: Bei einer Testmacht von 80% erhält man in ca. 20% der Fälle, in denen sich die wahren Mittelwerte um d unterscheiden (also in 1 von 5 solcher Fälle), keine Signifikanz. Wenn man den Versuch 5 mal durchführt, so erhält man im Schnitt nur 4 mal Signifikanz selbst wenn der wahre Unterschied in etwa d ist.

Theoretischer Hintergrund

Ziel: Wähle n so, dass die Nullhypothese fälschlicherweise höchstens mit Ws β nicht verworfen wird.

Die Nullhypothese wird nicht verworfen falls

$$\frac{|\bar{x} - \mu_0|}{s/\sqrt{n}} \leq t_{1-\frac{\alpha}{2}, n-1}$$

Falls nun die Nullhypothese nicht wahr ist, sondern die wahre Verteilung einen Mittelwert $\mu_1 \geq \mu_0 + d$ hat, so wird die Nullhypothese fälschlicherweise mit Ws

$$\Pr_{\mu_1} \left(\frac{|\bar{x} - \mu_0|}{s/\sqrt{n}} \leq t_{1-\frac{\alpha}{2}, n-1} \right)$$

nicht verworfen. Diese Ws soll kleiner als β sein.

Theoretischer Hintergrund

Nun verwenden wir, dass $\frac{\bar{x} - \mu_1}{s/\sqrt{n}}$ unter \Pr_{μ_1} t-verteilt ist mit **df=n-1**:

$$\begin{aligned} & \Pr_{\mu_1} \left(\frac{|\bar{x} - \mu_0|}{s/\sqrt{n}} \leq t_{1-\frac{\alpha}{2}, n-1} \right) \\ & \leq \Pr_{\mu_1} \left(\frac{\bar{x} - \mu_0}{s/\sqrt{n}} \leq t_{1-\frac{\alpha}{2}, n-1} \right) \\ & = \Pr_{\mu_1} \left(\frac{\bar{x} - \mu_1}{s/\sqrt{n}} \leq \frac{\mu_0 - \mu_1}{s/\sqrt{n}} + t_{1-\frac{\alpha}{2}, n-1} \right) \end{aligned}$$

Dies ist kleiner als β , falls

$$\frac{\mu_0 - \mu_1}{s/\sqrt{n}} + t_{1-\frac{\alpha}{2}, n-1} \leq t_{\beta, n-1} = -t_{1-\beta, n-1}$$

da $t_{\beta, n-1}$ so gewählt ist, dass $\text{pt}(t_{\beta, n-1}, \text{df}=\text{n-1}) = \beta$.

Theoretischer Hintergrund

Dies ist kleiner als β , falls

$$\frac{\mu_0 - \mu_1}{s/\sqrt{n}} + t_{1-\frac{\alpha}{2}, n-1} \leq t_{\beta, n-1} = -t_{1-\beta, n-1}$$

Also muss gelten (bei Multiplikation mit $\mu_0 - \mu_1 < 0$ wird \leq zu \geq)

$$\frac{\sqrt{n}}{s} \geq \frac{-t_{1-\beta, n-1} - t_{1-\frac{\alpha}{2}, n-1}}{\mu_0 - \mu_1} = \frac{t_{1-\beta, n-1} + t_{1-\frac{\alpha}{2}, n-1}}{\mu_1 - \mu_0}$$

Ist $\mu_1 - \mu_0 \approx d$, so muss die Stichprobenlänge mindestens

$$n \geq \frac{s^2 \cdot (t_{1-\frac{\alpha}{2}, n-1} + t_{1-\beta, n-1})^2}{d^2}$$

sein.

Einstichproben-Tests

Beispiel: Ist das Geschlechterverhältnis beim Kuhstärling bei der Geburt gleich 1 : 1?



photo (c) public domain

Einstichproben-Tests

Frage: Ist die relative Häufigkeit von männlichen Kuhstärlingen bei der Geburt gleich $\frac{1}{2}$?

Geplanter Test: (zweiseitiger) Einstichproben z-Test.

Ziel: Finde Unterschiede, die größer als $d = 0.02$ sind. Signifikanzniveau $\alpha = 5\%$. Testmacht $1 - \beta = 80\%$.

Vorwissen: Nicht nötig.

Frage: Bei wie vielen neugeborenen Kuhstärlingen muss das Geschlecht bestimmt werden?

Einstichproben-Tests

Lösung: Das Geschlecht ist Bernoulli-verteilt (2 mögliche Werte) mit Standardabweichung $\sqrt{p(1-p)}$. Allerdings kennen wir p nicht. Vermutlich wird das Geschlechterverhältnis in etwa 1 : 1 sein, also p nahe bei $\frac{1}{2}$. Als Näherung der Standardabweichung verwenden wir deshalb $s = \sqrt{\frac{1}{2}(1 - \frac{1}{2})} = \frac{1}{2}$.

Wähle n mindestens so groß, dass

$$n \geq \frac{s^2 \cdot (z_{1-\frac{\alpha}{2}} + z_{1-\beta})^2}{d^2}$$

Einstichproben-Tests

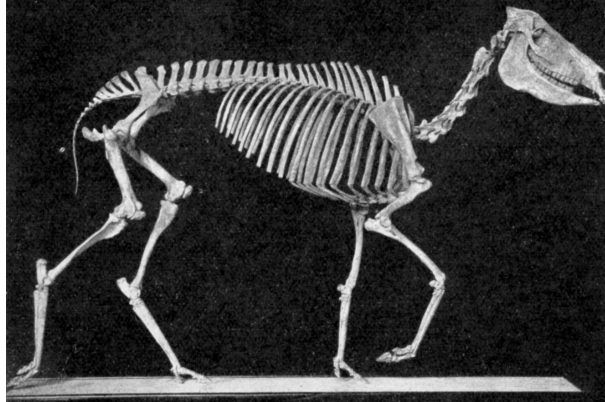
Berechnung:

$$\begin{aligned} n_0 &= \frac{s^2 \cdot (z_{1-\frac{\alpha}{2}} + z_{1-\beta})^2}{d^2} = \frac{\frac{1}{2^2} \cdot (z_{0.975} + z_{0.8})^2}{(0.02)^2} \\ &= 4905.55 \sim 4906 \end{aligned}$$

Die benötigte Stichprobenlänge wäre mindestens 4906! Diese Messreihe wird man vermutlich nicht durchführen wollen.

2.3 Zweistichproben-Test

Beispiel: Backenzähne von Hipparions



Hipparion

Bild: *Panicum miliaceum*
(c) public domain

Beispiel: Backenzähne von Hipparions

Frage: Unterscheidet sich die mesiodistale Länge (mm) der Backenzähne von *Hipparion africanum* und *Hipparion libycum*

Geplanter Test: (zweiseitiger) ungepaarter Zweistichproben t-Test.

Ziel: Finde Unterschiede, die größer als $d = 2.5$ mm sind. Signifikanzniveau $\alpha = 5\%$. Testmacht $1 - \beta = 80\%$.

Vorwissen: Standardabweichung bei *H. africanum* ist in etwa $s_A = 2.2$. Standardabweichung bei *H. libycum* ist in etwa $s_L = 4.3$.

Frage: Bei wie vielen Backenzähnen muss die mesiodistale Länge gemessen werden?

Lösung: In jeder Gruppe muss die Stichprobenlänge mindestens

$$n \geq \frac{(s_A^2 + s_L^2) \cdot (t_{1-\frac{\alpha}{2}, 2*n-2} + t_{1-\beta, 2*n-2})^2}{d^2}$$

sein.

Beispiel: Backenzähne von Hipparions

Berechnung:

$$\begin{aligned} n_0 &= \frac{(s_A^2 + s_L^2) \cdot (z_{1-\frac{\alpha}{2}} + z_{1-\beta})^2}{d^2} \\ &= \frac{(2.2^2 + 4.3^2) \cdot (z_{0.975} + z_{0.8})^2}{2.5^2} = 29.3 \sim 30 \\ n_1 &= \frac{(s_A^2 + s_L^2) \cdot (t_{1-\frac{\alpha}{2}, 2*n_0-2} + t_{1-\beta, 2*n_0-2})^2}{d^2} \\ &= \frac{(2.2^2 + 4.3^2) \cdot (t_{0.975, 58} + t_{0.8, 58})^2}{(2.5)^2} \\ &= 30.3 \sim 31 \\ n_2 &= \frac{(s_A^2 + s_L^2) \cdot (t_{1-\frac{\alpha}{2}, 2*n_1-2} + t_{1-\beta, 2*n_1-2})^2}{d^2} \\ &= 30.28 \sim 31 \end{aligned}$$

Beispiel: Backenzähne von Hipparions

Antwort: Es müssen mindestens 31 Backenzähne von *H. africanum* und 31 Backenzähne von *H. libycum* vermessen werden.

2.4 Einseitige Tests

Wenn man einseitig testen will, so muss man in obigen Formeln $t_{1-\frac{\alpha}{2},n-1}$ durch $t_{1-\alpha,n-1}$ ersetzen.

Beispiel: Wachstumshormon

Frage: Wirkt das zu untersuchende Wachstumshormon signifikant besser als ein Placebo?

Geplanter Test: einseitiger ungepaarter Zweistichproben t-Test.

Ziel: Finde Unterschiede, die größer als $d = 2$ sind. Signifikanzniveau $\alpha = 5\%$. Testmacht $1 - \beta = 80\%$.

Vorwissen: Standardabweichung ist in jeder Gruppe in etwa $s = 4$.

Frage: Wie viele Personen braucht man in der Kontrollgruppe und in der Versuchsgruppe?

Lösung: In jeder Gruppe muss die Stichprobenlänge mindestens

$$n \geq \frac{(s^2 + s^2) \cdot (t_{1-\alpha,2*n-2} + t_{1-\beta,2*n-2})^2}{d^2}$$

sein. Ergebnis: $n = 51$.

2.5 Übersicht

Zweiseitiger Einstichproben t-Test

Geplanter Test: Zweiseitiger Einstichproben t-Test.

Ziel: Finde Unterschiede, die größer als d sind. Signifikanzniveau α . Testmacht $1 - \beta$.

Vorwissen: Standardabweichung bei Vortest war s

Lösung: Es soll gelten:

$$n \geq \frac{s^2 \cdot (t_{1-\frac{\alpha}{2},n-1} + t_{1-\beta,n-1})^2}{d^2}$$

Zweiseitiger ungepaarter Zweistichproben t-Test

Geplanter Test: Zweiseitiger ungepaarter Zweistichproben t-Test.

Ziel: Finde Unterschiede, die größer als d sind. Signifikanzniveau α . Testmacht $1 - \beta$.

Vorwissen: Die Standardabweichungen in den beiden Stichproben sind in etwa s_1 beziehungsweise s_2 .

Lösung: In jeder Gruppe muss die Stichprobenlänge mindestens

$$n \geq \frac{(s_1^2 + s_2^2) \cdot (t_{1-\frac{\alpha}{2},2*n-2} + t_{1-\beta,2*n-2})^2}{d^2}$$

sein.

Zweiseitiger gepaarter Zweistichproben t-Test

Geplanter Test: Zweiseitiger gepaarter Zweistichproben t-Test.

Ziel: Finde Unterschiede, die größer als d sind. Signifikanzniveau α . Testmacht $1 - \beta$.

Vorwissen: Standardabweichung der Differenz der beiden Stichproben ist in etwa s_d .

Lösung: In jeder Gruppe muss die Stichprobenlänge mindestens

$$n \geq \frac{s_d^2 \cdot (t_{1-\frac{\alpha}{2},n-1} + t_{1-\beta,n-1})^2}{d^2}$$

sein.

Einseitiger Einstichproben t-Test

Geplanter Test: Einseitiger Einstichproben t-Test.

Ziel: Finde Unterschiede, die größer als d sind. Signifikanzniveau α . Testmacht $1 - \beta$.

Vorwissen: Standardabweichung bei Vortest war s

Lösung: Es soll gelten:

$$n \geq \frac{s^2 \cdot (t_{1-\alpha, n-1} + t_{1-\beta, n-1})^2}{d^2}$$

Einseitiger ungepaarter Zweistichproben t-Test

Geplanter Test: Einseitiger ungepaarter Zweistichproben t-Test.

Ziel: Finde Unterschiede, die größer als d sind. Signifikanzniveau α . Testmacht $1 - \beta$.

Vorwissen: Die Standardabweichungen in den beiden Stichproben sind in etwa s_1 und s_2 .

Lösung: In jeder Gruppe muss die Stichprobenlänge mindestens

$$n \geq \frac{(s_1^2 + s_2^2) \cdot (t_{1-\alpha, 2*n-2} + t_{1-\beta, 2*n-2})^2}{d^2}$$

sein.

Einseitiger gepaarter Zweistichproben t-Test

Geplanter Test: Einseitiger gepaarter Zweistichproben t-Test.

Ziel: Finde Unterschiede, die größer als d sind. Signifikanzniveau α . Testmacht $1 - \beta$.

Vorwissen: Standardabweichung der Differenz der beiden Stichproben ist in etwa s_d .

Lösung: In jeder Gruppe muss die Stichprobenlänge mindestens

$$n \geq \frac{s_d^2 \cdot (t_{1-\alpha, n-1} + t_{1-\beta, n-1})^2}{d^2}$$

sein.

2.6 Stichprobenlänge ermitteln mit R

In R ermittelt man die benötigte Stichprobenlänge mit

```
power.t.test(n = , delta = , sd = , sig.level = ,
             power = ,
             type = c("two.sample", "one.sample", "paired"),
             alternative = c("two.sided", "one.sided") )
```

Die Argumente sind:

- `n` = Stichprobenlänge (pro Gruppe bzw pro Stichprobe)
- `delta` = d (minimale Differenz, detection level)
- `sd` = s (vermutete Standardabweichung pro Gruppe)
- `sig.level` = α (Signifikanzniveau)
- `power` = $1 - \beta$ (Testmacht)

Genau eines der Argumente `n`, `delta`, `sd`, `sig.level`, `power` muss als NULL übergeben werden. Dieses wird dann berechnet.

Beispiele:

- CCRT bei *D. ananassae*: $d = 4$, $s = 11.9$, $\alpha = 5\%$, $\beta = 0.2$

```
> power.t.test(n=NULL, delta=4, sd=11.9,
+ sig.level=0.05, power=0.8,
+ type="one.sample", alternative="two.sided")
```

One-sample t test power calculation

```
      n = 71.41203
delta = 4
      sd = 11.9
sig.level = 0.05
      power = 0.8
alternative = two.sided
```

- Relative Häufigkeit von männlichen Kuhstärlingen bei der Geburt: $d = 0.02$, $s = \sqrt{\frac{1}{2}(1 - \frac{1}{2})}$, $\alpha = 5\%$, $\beta = 0.2$

```
> power.t.test(n=NULL, delta=0.02, sd=0.5,
+ sig.level=0.05, power=0.8,
+ type="one.sample", alternative="two.sided")
```

One-sample t test power calculation

```
      n = 4907.471
delta = 0.02
      sd = 0.5
sig.level = 0.05
      power = 0.8
alternative = two.sided
```

(wir verwenden `power.t.test` als Approximation, da `power.z.test` nicht existiert)

- Backenzähne von Hipparions: $d = 2.5$, $s = \sqrt{(2.2^2 + 4.3^2)/2}$, $\alpha = 5\%$, $\beta = 0.2$

```
> power.t.test(n=NULL, delta=2.5,
+ sd=sqrt( (2.2^2+4.3^2)/2 ),
+ sig.level=0.05, power=0.8,
+ type="two.sample", alternative="two.sided")
```

Two-sample t test power calculation

```
      n = 30.28929
delta = 2.5
      sd = 3.415406
sig.level = 0.05
      power = 0.8
alternative = two.sided
```

NOTE: n is number in *each* group

- Wachstumshormon: (einseitiger Test) $d = 2$, $s = 4$, $\alpha = 5\%$, $\beta = 0.2$

```
> power.t.test(n=NULL, delta=2, sd=4,
+ sig.level=0.05, power=0.8,
+ type="two.sample", alternative="one.sided")
```

Two-sample t test power calculation

```
      n = 50.1508
    delta = 2
      sd = 4
sig.level = 0.05
  power = 0.8
alternative = one.sided
```

NOTE: n is number in *each* group

Der Befehl `power.t.test()` kann auch dazu benutzt werden, die Testmacht zu berechnen, wenn man sich auf die Stichprobenlänge bereits festgelegt hat. Beispiel: CCRT bei *D. ananassae*: $n = 100$, $d = 4$, $s = 11.9$, $\alpha = 5\%$

```
> power.t.test(n=100, delta=4, sd=11.9,
+ sig.level=0.05, power=NULL,
+ type="one.sample", alternative="two.sided")
```

One-sample t test power calculation

```
      n = 100
    delta = 4
      sd = 11.9
sig.level = 0.05
  power = 0.9144375
alternative = two.sided
```

2.7 F-Test

F-Test

Will man testen, ob die Mittelwerte bei 3 oder mehr Gruppen gleich sind, so verwendet man den F-Test. Um eine Aussage über die Stichprobenlänge treffen zu können, benötigt man die Variabilität innerhalb der Gruppen und die Variabilität zwischen den Gruppen (z.B. aus Vorversuchen).

Die Formel für die benötigte Stichprobe ist hier weniger übersichtlich. Deshalb konzentrieren wir uns auf die Berechnung mit R.

Wir zeigen an folgendem Beispiel, wie man den R-Befehl `power.anova.test()` einsetzt, um die benötigte Stichprobenlänge zu ermitteln.

Beispiel: Blutgerinnungszeit bei Ratten

Frage: Unterscheidet sich die Blutgerinnungszeit bei Ratten unter 4 verschiedenen Behandlungen?

Geplanter Test: F-Test.

Signifikanzniveau: $\alpha = 5\%$

Testmacht: $1 - \beta = 90\%$.

Vorwissen: Standardabweichung innerhalb jeder Gruppe ist in etwa $s_{\text{innerh}} = 2.4$. Beachte: $s_{\text{innerh}}^2 = ss_{\text{innerh}}/df_{\text{innerh}}$. Standardabweichung zwischen den Gruppen ist in etwa $s_{\text{zw}} = 1.2$. Beachte: $s_{\text{zw}}^2 = ss_{\text{zw}}/df_{\text{zw}}$.

Frage: Bei wie vielen Ratten muss die Blutgerinnungszeit gemessen werden?

Beispiel: Blutgerinnungszeit bei Ratten

```
> power.anova.test(groups=4, n=NULL, between.var=1.2^2,
  within.var=2.4^2, sig.level=0.05, power=0.9)
```

Balanced one-way analysis of variance power calculation

```
groups = 4
n = 19.90248
between.var = 1.44
within.var = 5.76
sig.level = 0.05
power = 0.9
NOTE: n is number in each group
```

Antwort: Für jede der 4 Behandlungen braucht man mindestens 20 Ratten.

3 Stichprobenwahl

3.1 Überspitzte Beispiele

Um die Problematik der Stichprobenwahl zu verdeutlichen, beginnen wir überspitzten Beispielen.

- Um die Parteienpräferenz in Deutschland zu messen, stellt ein Wahlforschungsunternehmen die Sonntagsfrage („Was würden Sie wählen, wenn kommenden Sonntag Bundestagswahl wäre“) an 1000 zufällig ausgewählte Bürger aus Garmisch-Partenkirchen.[5mm] **Keine repräsentative Stichprobe!**[5mm] Die Einwohner- und Meinungsstruktur von Garmisch-Partenkirchen ist nicht typisch für Deutschland.

•



photo (c) André Karwath (Bild zeigt eine *Drosophila melanogaster*)

Um die Chill-Coma Recovery Time (CCRT) der europäischen *Drosophila melanogaster* mit der taiwanesischen Population zu vergleichen, werden *Drosophila melanogaster* an jeweils 10 verschiedenen Orten in Frankreich, Spanien und Italien gesammelt.[4mm] **Keine repräsentative Stichprobe!**[4mm] Die CCRT von Fruchtfliegen in Südeuropa ist nicht typisch für die CCRT europäischer Fruchtfliegen.

•

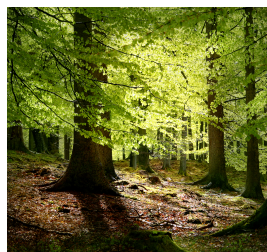


photo (c) Malene Thyssen (Bild zeigt einen Rotbuchenwald in Dänemark)

Um die Blätterdichte in oberbayerischen Wäldern zu messen, wird in 10 zufällig ausgewählten oberbayerischen Wäldern die Blätterdichte entlang des Waldrandes und entlang von Waldwegen gemessen.[3mm] **Keine repräsentative Stichprobe!**[3mm] Am Waldrand und auch entlang von Waldwegen ist die Blätterdichte überdurchschnittlich hoch.

- Waldameisen



photo (c) Oswald Hicker

Es sollen 100 französische Waldameisen gesampelt werden. Dazu wird ein Ameisennest zufällig in Frankreich ausgewählt und hiervon 100 Ameisen genommen.[5mm] **Keine repräsentative Stichprobe der Länge 100!**[5mm] Die erste gesampelte Ameise ist wohl eine typische französische Waldameise. Die weiteren sind aber vermutlich mit der ersten Ameise nahe verwandt. Für eine Stichprobe der Länge 100 braucht man 100 'unabhängige' Ameisen. Kommen die 100 Ameisen aus demselben Ameisennest, so können sie Geschwister sein und sind dann sicherlich nicht unabhängig voneinander.

- 20 zufällig ausgewählte Studenten werden eingeladen, an einem Versuch teilzunehmen. Die ersten 10 Studenten, die am Versuchsort ankommen, bilden die Kontrollgruppe. Die weiteren 10 Studenten bilden die Versuchsgruppe.[5mm] **Die beiden Versuchsgruppen sind nicht identisch verteilt!**[5mm] Die Kontrollgruppe besteht aus pünktlicheren Studenten. Diese Gruppe könnte engagierter am Versuch teilnehmen. Dadurch wird das Ergebnis verfälscht.

Dr. X könnte argumentieren: *Wir haben mit einem zusätzlichen Test gezeigt, dass es die Reihenfolge der Studierenden keinen Einfluss auf die Versuchsergebnisse hatte.*

Was halten Sie von diesem Argument?

Eine solche Argumentation ist aus statistischer Sicht Blödsinn!

- Ein statistischer Test kann niemals zeigen, dass ein Effekt nicht existiert.
- Vermutlich meint Dr. X, dass er einen statistischen Test durchgeführt hat, bei dem es keinen statistisch signifikanten Zusammenhang zwischen Pünktlichkeit und Versuchsergebnis gab.
- Man darf aber aus nicht-Signifikanz niemals schließen, dass es den Effekt nicht gibt.
- Vielleicht ist der Effekt so schwach, dass der Vortest geringe Macht hatte, aber immer noch stark genug um die spätere statistische Analyse zu verfälschen.

3.2 Zufallsstichprobe

Zufallsstichprobe

Eine **Zufallsstichprobe** der Länge n aus einer Gesamtpopulation der Größe N erhält man wie folgt:

- Nummeriere N identische Kugeln von 1 bis N .
- Durchmische die N Kugeln in einem Beutel oder ähnlichem.
- Ziehe (ohne Zurücklegen) n Kugeln.

Die zu den Nummern auf den Kugeln gehörigen Individuen in der Gesamtpopulation bilden dann eine Zufallsstichprobe.

Beispiel

Ziel: Man möchte eine Umfrage unter allen Bachelor-Studenten der Biologie an der LMU München durchführen. Es zu aufwändig ist, alle Studenten zu befragen. Also möchte man 50 Studenten zufällig auswählen, um diese dann zu befragen.

Vorgehen: Die Anzahl N an Studenten ist bekannt. Nun nummerieren wir die Studenten durch und ziehen 50 Nummern rein zufällig. Dies könnte man in R durchführen:

```
sample(1:N, size=50, replace=FALSE)
```

Dieses Vorgehen wird oft als **Lotterieverfahren** bezeichnet.

In Anwendungen ist dies meist nicht möglich, da

- die Größe der Gesamtpopulation meist unbekannt ist (zB: Anzahl an Ameisen, Anzahl an *Drosophila melanogaster*)
- beziehungsweise es bei großen Populationen schwierig ist, den Individuen Nummern zuzuweisen.

Eine **Zufallsstichprobe** ist Teil einer Gesamtpopulation, die durch einen Auswahlprozess mit Zufallsprinzip aus der Gesamtpopulation entnommen wird und stellvertretend, repräsentativ für die Gesamtpopulation ist.

Ein Teil einer Gesamtpopulation kann auch dann als repräsentative Stichprobe angesehen werden, wenn das Auswahlverfahren zwar nicht zufällig, aber von den auszuwertenden Merkmalen stochastisch unabhängig ist.

Anders formuliert: Die Stichprobe muss bezüglich den auszuwertenden Merkmalen typisch für die Gesamtpopulation sein.

Betrachtet man eine „Stichprobe, die gerade zur Hand ist“ und die keine Zufallsstichprobe ist, so darf man Aussagen über die Stichprobe nicht auf die Gesamtpopulation verallgemeinern.

Beispiel

Ziel: Stichprobe von 100 Mäusen.

Beachte: Für die statistische Analyse wird Unabhängigkeit vorausgesetzt. Insbesondere dürfen die Mäuse nicht verwandt sein.

Falsch: 100 Mäuse von demselben Bauernhof. Denn: Von demselben Bauernhof sammelt man mit gewisser Ws verwandte Mäuse. Extremfall: Nimmt man 100 Klone derselben Maus, so ist die tatsächliche Stichprobenlänge gleich 1 (= Anzahl voneinander unabhängiger Mäuse).

Richtig: (Wird jedenfalls in der Literatur akzeptiert)

- Je eine Maus pro Bauernhof.
- Bauernhöfe müssen mindestens 1km voneinander entfernt sein.

Beispiel

Beachte: Sampelt man Mäuse von verschiedenen Bauernhöfen in der Gegend von Memmingen, so ist die Stichprobe nur repräsentativ für die Region Memmingen.

Es darf bezweifelt werden, ob diese Stichprobe repräsentativ für Deutschland oder gar Europa ist.

3.3 Elimination von nicht-interessierenden Einflussgrößen

Nun geht es nicht mehr um Zufallsstichproben, sondern um die Einteilung von Versuchsobjekten in verschiedene Behandlungsgruppen.

Prinzipien der Versuchsplanung

Wir sprechen nun von [Einflussgrößen](#) bzw von [Einflussfaktoren](#) und von [Zielgrößen](#).

Einflussgröße kann so ziemlich alles sein:

- Wurde die Behandlung angewendet: Ja oder Nein?
- Wer hat die Messung durchgeführt?
- Wurde ein großes oder kleines Reagenzglas verwendet?
- Wie waren die Lichtverhältnisse im Labor während des Versuchs?

Prinzip

Nicht interessierende Einflussgrößen sind im Versuch möglichst konstant zu halten.

Prinzipien der Versuchsplanung

Beispiele für die Einhaltung dieses Prinzips:

- Derselbe Experimentator für alle Versuche.
- Weder Experimentator noch Versuchsperson wissen, zu welcher Behandlungsgruppe die Versuchsperson gehört. [Doppelblind](#) (Ausschluss von subjektiven Einflussfaktoren).
- Dieselben oder zumindest baugleiche Materialien und Laborbedingungen bei allen Versuchen.
- Reihenfolge der Behandlungsgruppen ist zufällig. (Also nicht: Versuchsgruppe, Kontrollgruppe, Versuchsgruppe, Kontrollgruppe, ...)

3.4 Blockbildung

Sind die Versuchsobjekte sehr unterschiedlich, so empfiehlt sich eine Zusammenfassung von sehr ähnlichen Versuchsobjekten zu Untergruppen. Die für das Versuchsziel wichtigen Vergleiche werden dann möglichst innerhalb der Blöcke vorgenommen.

Beachte: Die Bildung von Blöcken ist nur dann sinnvoll, wenn die Streuung zwischen den Versuchsobjekten deutlich größer ist als die Streuung zwischen den verschiedenen Behandlungen.

Zweck der Blockbildung ist es, die Genauigkeit blockinterner Vergleiche zu erhöhen.

Beispiel

Frage: Wirkt eine gewisse Diät besser als Placebo?

Problem: Nehmen wir, die Diät verringert das Gewicht tatsächlich im Mittel um 3 kg. Da das Gewicht bei den Versuchspersonen aber sehr stark zwischen 50 kg und 130 kg schwankt, braucht man sehr viele Versuchspersonen, um den kleinen Unterschied festzustellen.

Lösung: Unterteile die Versuchspersonen in Untergruppen gleicher Gruppengröße, so dass die Personen in jeder Untergruppe ähnliches Gewicht haben. Jede Untergruppe wird dann in Diätgruppe und Kontrollgruppe aufgeteilt. Die Gewichtsvergleiche finden dann in jeder Untergruppe statt.

3.5 Balanced Design vs Non-Balanced Design

Balanciertes Design bedeutet, dass jede Gruppe aus gleich vielen Versuchsobjekten besteht. In jeder Behandlungsgruppe hat man also dieselbe Stichprobenlänge.

Im Normalfall bevorzugt man ein balanciertes Versuchs-Design

Vorteil des balancierten Versuchs-Designs: Viele statistische Verfahren setzen balanciertes Design voraus (z.B Tukey's simultane Konfidenzintervalle).

Nachteil des balancierten Versuchs-Designs: Eine balanciertes Design ist nicht immer repräsentativ. Beispiel: Die untypische Gewichtsklasse 140 – 150 kg wird im balancierten Design überrepräsentiert.

3.6 Randomisierung

Randomisierung

Randomisierung ist die zufällige Zuordnung der Behandlungen zu den gegebenen Versuchsobjekten.

Vorgehen: Nummeriere die Versuchsobjekte und wende das Lotterieverfahren an.

Beispiel: Ein Medikament zur Steigerung der Konzentration soll getestet werden an 20 Studenten.

Falsch: Die 10 Studenten, die zuerst im Labor eintreffen, bekommen das Medikament. Die nächsten 10 Studenten bekommen das Placebo. Problem hier: Pünktlichere Studenten können sich vielleicht von vornherein besser konzentrieren.

Richtig: Die Studenten werden von 1 bis 20 durchnummeriert. Die Kontrollgruppe besteht dann aus den Studenten mit Nummern

```
sample(1:20,size=10,replace=FALSE)
19 16  1 13 18 10  2  5  9 14
```

(Natürlich gibt es viele weitere Verfahren, eine Zufallszuordnung zu erreichen.)