

Mathematik für Biologen

Prof. Dr. Rüdiger W. Braun

<http://blog.ruediger-braun.net>

Heinrich-Heine-Universität Düsseldorf

31. Januar 2014

- 1 Exakter Test nach Fisher
 - Mendelsche Erbgeltn als Beispiel
 - Test auf Übereinstimmung zweier Verteilungen
- 2 Der F -Test zum Vergleich zweier Varianzen
- 3 ANOVA
 - Beispielhafte Fragestellung
 - Idee der Varianzanalyse
 - Gruppenmittelwerte
 - Gesamtmittelwert
 - Zerlegung der Varianz
 - Teststatistik
 - Zusammenfassung der Varianzanalyse
 - Beispiele

Mendelsche Erbgeltn

- Bei den Mendelschen Erbversuchen tritt das Merkmal *Blütenfarbe* in drei Ausprägungen auf, nämlich weiß, rosa und rot
- weiß und rot haben dieselbe Wahrscheinlichkeit, rosa die doppelte
- 4 Blüten werden beobachtet, alle sind rosa
- Ist diese Beobachtung zum Signifikanzniveau $\alpha = 0.05$ mit den Mendelschen Regeln vereinbar?

Interpretation als Vergleich zweier Verteilungen

- Modellannahme: Die Mendelschen Regeln gelten für die untersuchte Situation
- Das entspricht der Verteilung

Nummer	Ausprägung	Wahrscheinlichkeit
1	weiß	25%
2	rosa	50%
3	rot	25%

- Zu vergleichen mit der tatsächlichen Verteilung der Blütenfarben in dem Kollektiv
- Der Stichprobenumfang ist 4
- Das ist für praktische Zwecke zu wenig, lässt sich aber gut von Hand rechnen

Mendelsche Erbgelien, Fortsetzung

- Ordne die möglichen Ergebnisse mit aufsteigender Wahrscheinlichkeit an
- Entscheidungsstrategie am Beispiel $\alpha = 0.05$

Lehne H_0 ab, wenn die Beobachtung zu den 5% unwahrscheinlichsten Ereignissen gehört

Test auf Übereinstimmung zweier Verteilungen

- Unabhängige Zufallsvariable X_1, \dots, X_n , die alle mit Wahrscheinlichkeit p_1 den Wert w_1 , mit Wahrscheinlichkeit p_2 den Wert w_2, \dots , mit Wahrscheinlichkeit p_s den Wert w_s annehmen
- Vergleichswahrscheinlichkeiten $\pi_1, \pi_2, \dots, \pi_s$ mit $\pi_1 + \pi_2 + \dots + \pi_s = 1$
- Nullhypothese und Alternative:

$$H_0 : p_1 = \pi_1, p_2 = \pi_2, \dots, p_s = \pi_s$$

$$H_1 : \text{mindestens ein } p_j \neq \pi_j$$

Test auf Übereinstimmung zweier Verteilungen: Summenvariable

- Summenvariable

$$Y_1 = \text{Anzahl aller } X_j \text{ mit } X_j = w_1$$

$$Y_2 = \text{Anzahl aller } X_j \text{ mit } X_j = w_2$$

$$\vdots$$

$$Y_s = \text{Anzahl aller } X_j \text{ mit } X_j = w_s$$

- Erwartungswerte unter H_0

$$E(Y_1) = n \cdot \pi_1$$

$$E(Y_2) = n \cdot \pi_2$$

$$\vdots$$

$$E(Y_s) = n \cdot \pi_s$$

Exakter Test nach Fisher

- Bestimme für jede mögliche Kombination von Werten von Y_1, \dots, Y_s deren Wahrscheinlichkeit
- Ordne diese Wahrscheinlichkeiten aufsteigend in einer Liste
- Der kritische Bereich, in dem H_0 abgelehnt wird, besteht aus den obersten Zeilen dieser Liste
- Man nimmt genau so viele Zeilen, dass die erlaubte Fehlerwahrscheinlichkeit erster Art nicht überschritten, aber möglichst gut ausgeschöpft wird

Beispiel Mendel: Formalisierung

- $s = 3$
- X_1 ist der Zahlencode der Blütenfarbe der ersten Blüte, X_2 dasselbe für die zweite Blüte, ...
- Y_1 bezeichnet die Anzahl der weißen, Y_2 die der rosafarbenen und Y_3 die der roten Blüten
- Dann $Y_1 + Y_2 + Y_3 = 4$
- Im Beispiel $Y_1 = 0$, $Y_2 = 4$, $Y_3 = 0$
- Rechne sämtliche Einzelwahrscheinlichkeiten aus

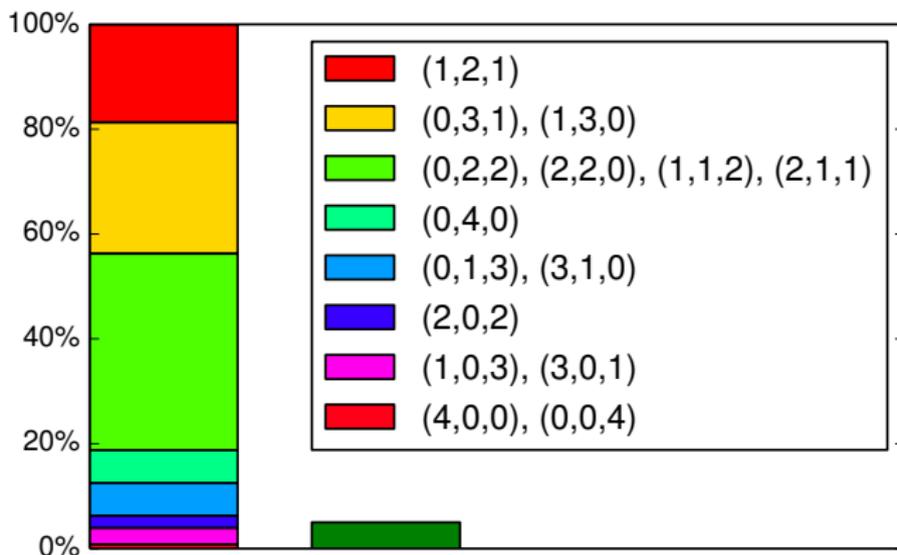
Beispiel Mendel: Wahrscheinlichkeiten der Einzelereignisse

$$\begin{aligned}
 P(Y_1 = k_1, Y_2 = k_2, Y_3 = k_3) &= \binom{4}{k_1} \cdot \binom{4 - k_1}{k_2} \cdot \left(\frac{1}{4}\right)^{k_1} \cdot \left(\frac{1}{2}\right)^{k_2} \cdot \left(\frac{1}{4}\right)^{k_3} \\
 &= \frac{4! \cdot (4 - k_1)!}{k_1! \cdot (4 - k_1)! \cdot k_2! \cdot (4 - k_1 - k_2)!} \cdot \left(\frac{1}{4}\right)^{k_1} \cdot \left(\frac{1}{2}\right)^{k_2} \cdot \left(\frac{1}{4}\right)^{k_3} \\
 &= \frac{4!}{k_1! \cdot k_2! \cdot k_3!} \cdot \left(\frac{1}{4}\right)^{k_1} \cdot \left(\frac{1}{2}\right)^{k_2} \cdot \left(\frac{1}{4}\right)^{k_3}
 \end{aligned}$$

Beispiel Mendel: Tabelle der W'keiten der Einzelereignisse

k_1	k_2	k_3	$P(X_1 = k_1, X_2 = k_2, X_3 = k_3)$	kumulierte Summe
0	0	4	0.0039	0.0039
4	0	0	0.0039	0.0078
1	0	3	0.0156	0.0234
3	0	1	0.0156	0.0391
2	0	2	0.0234	0.0625
0	1	3	0.0312	0.0938
3	1	0	0.0312	0.1250
0	4	0	0.0625	0.1875
0	2	2	0.0938	0.2812
1	1	2	0.0938	0.3750
2	1	1	0.0938	0.4688
2	2	0	0.0938	0.5625
0	3	1	0.1250	0.6875
1	3	0	0.1250	0.8125
1	2	1	0.1875	1.0000

Beispiel Mendel: Balkendiagramm



Der linke Balken zeigt die kumulierten Werte aus der Tabelle, der rechte die 5%-Schwelle

Beispiel Mendel: Ergebnis

- In den folgenden Fällen kann die Nullhypothese zum Signifikanzniveau $\alpha = 0.05$ abgelehnt werden
 - 4 weiße oder 4 rote Blüten
 - keine rosa, aber 3 weiße oder 3 rote Blüten
- Der p -Wert des beobachteten Ereignisses “4 rosa Blüten” beträgt 18.75%

Exakter Test nach Fisher

Warum vergleicht man zwei Varianzen?

- Um Unterschiedlichkeit zweier Verteilungen nachzuweisen
- Um Voraussetzungen eines anderen Tests zu prüfen
- Um eine ANOVA zu rechnen
- ANOVA= "Analysis of Variance": Ein Test, mit welchem man den Einfluss der Gruppenzugehörigkeit auf einen Parameter prüfen kann, indem man Varianzen vergleicht

F-Test zum Vergleich zweier Varianzen

- X_1, \dots, X_{n_1} und Y_1, \dots, Y_{n_2} bezeichnen zwei Gruppen von Messwerten
- Verteilungsvoraussetzungen:
 - Die X_j sind verteilt gemäß $N(\mu_1, \sigma_1^2)$, wobei μ_1 und σ_1 unbekannt sind
 - Die Y_j sind verteilt gemäß $N(\mu_2, \sigma_2^2)$, wobei μ_2 und σ_2 unbekannt sind
- Ziel: σ_1 und σ_2 sollen verglichen werden

F-Test, Teststatistik

- x_j und y_j seien Realisierungen.
- Bestimme arithmetische Mittelwerte und Stichprobenstreuungen

$$\bar{x} = \frac{1}{n_1} \sum_{j=1}^{n_1} x_j$$

$$s_x = \sqrt{\frac{1}{n_1 - 1} \sum_{j=1}^{n_1} (x_j - \bar{x})^2}$$

$$\bar{y} = \frac{1}{n_2} \sum_{j=1}^{n_2} y_j$$

$$s_y = \sqrt{\frac{1}{n_2 - 1} \sum_{j=1}^{n_2} (y_j - \bar{y})^2}$$

- Die Teststatistik ist

$$t = \frac{s_x^2}{s_y^2}$$

Die F -Verteilung

- Die F -Verteilung mit $n_1 - 1$ und $n_2 - 1$ Freiheitsgraden ist definiert als die Verteilung derjenigen Zufallsvariablen Z , die erklärt ist durch

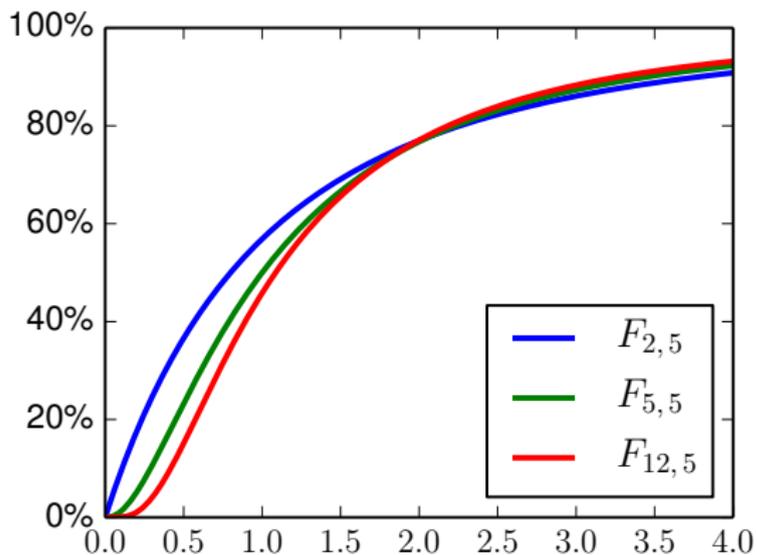
$$Z = \frac{(n_2 - 1) \cdot \sum_{j=1}^{n_1} (X_j - \bar{X})^2}{(n_1 - 1) \cdot \sum_{j=1}^{n_2} (Y_j - \bar{Y})^2}$$

wenn X_1, \dots, X_{n_1} und Y_1, \dots, Y_{n_2} unabhängige, standardnormalverteilte Zufallsvariablen sind.

- Für f_1 bzw. f_2 Freiheitsgrade sind die Quantile der F -Verteilung für α nahe 1 tabelliert.
- Für α nahe 0 benutzt man die Formel

$$F_{f_1, f_2, \alpha} = \frac{1}{F_{f_2, f_1, 1-\alpha}}$$

Verteilungsfunktionen von F -Verteilungen



Quantile $f_{f_1, f_2, 0.95}$ der F -Verteilungen

f_2	f_1					
	1	2	3	4	5	6
1	161.45	199.50	215.71	224.58	230.16	233.99
2	18.51	19.00	19.16	19.25	19.30	19.33
3	10.13	9.55	9.28	9.12	9.01	8.94
4	7.71	6.94	6.59	6.39	6.26	6.16
5	6.61	5.79	5.41	5.19	5.05	4.95
6	5.99	5.14	4.76	4.53	4.39	4.28
7	5.59	4.74	4.35	4.12	3.97	3.87
8	5.32	4.46	4.07	3.84	3.69	3.58
9	5.12	4.26	3.86	3.63	3.48	3.37
10	4.96	4.10	3.71	3.48	3.33	3.22

Beispiele:

- $f_{4, 5, 0.95} = 5.19$
- $f_{5, 4, 0.95} = 6.26$

F -Test, Entscheidungsregel

- Das Signifikanzniveau sei α
- Wir bestimmen die Quantile der F -Verteilung (genauer s. u.)

$$F_{n_1-1, n_2-1, 1-\alpha/2} \quad \text{beim zweiseitigen Test}$$

$$F_{n_1-1, n_2-1, 1-\alpha} \quad \text{bei einem einseitigen Test}$$

- Entscheidung

$H_0 = \{\sigma_1 = \sigma_2\}$: Die Nullhypothese H_0 wird abgelehnt, wenn

$$t > F_{n_1-1, n_2-1, 1-\alpha/2} \quad \text{oder} \quad t < \frac{1}{F_{n_2-1, n_1-1, 1-\alpha/2}}$$

$H_0 = \{\sigma_1 \leq \sigma_2\}$: Die Nullhypothese H_0 wird abgelehnt, wenn

$$t > F_{n_1-1, n_2-1, 1-\alpha}$$

$H_0 = \{\sigma_1 \geq \sigma_2\}$: Die Nullhypothese H_0 wird abgelehnt, wenn

$$t < \frac{1}{F_{n_2-1, n_1-1, 1-\alpha}}$$

F-Test, Beispiel

- Beim t -Test für unverbundene Stichproben hatten wir im Beispiel "Bodenbakterium" zwei Datensätze erhalten
- Der erste hatte $n_x = 10$ und $s_x = 7.972$
- Der zweite hatte $n_y = 9$ und $s_y = 6.280$
- Können wir zum Signifikanzniveau $\alpha = 0.10$ ausschließen, dass beide Verteilungen dieselbe Varianz aufweisen?
- Teststatistik

$$t = \frac{s_x^2}{s_y^2} = \frac{63.55}{39.44} = 1.611$$

- Das benötigte Quantil ist $f_{9,8,0.95} = 3.388$
- H_0 kann nicht abgelehnt werden

ANOVA

Beispielhafte Fragestellung: Unterrichtsmethoden

- Jeweils 4 bis 5 Schüler wurden nach einer von 4 Methoden in Statistik unterrichtet. Hat die Wahl der Unterrichtsmethode überhaupt einen Einfluss auf den Lernerfolg?
- Daten: (Der Erfolg wurde auf einer Skala von 0 bis 8 gemessen)

Unterrichtsmethode			
1	2	3	4
2	3	6	5
1	4	8	5
3	3	7	5
1	5	6	3
	0	8	2

Unterrichtsmethoden, Fortsetzung

- Die Nullhypothese ist, dass alle Methoden bis auf zufällige Abweichungen dasselbe Ergebnis liefern
- Wir könnten je zwei Unterrichtsmethoden mit einem t -Test für unverbundene Stichproben testen
- Das sind $\binom{4}{2} = 6$ Paarvergleiche

Unterrichtsmethoden, Fortsetzung

- Für jeden einzelnen Paarvergleich sei α die Fehlerwahrscheinlichkeit erster Art
- Wir nehmen an, dass die Paarvergleiche unabhängig sind
- Mit Wahrscheinlichkeit $p = 1 - \alpha$ vermeidet ein einzelner Paarvergleich den Fehler erster Art
- Mit Wahrscheinlichkeit $(1 - \alpha)^6$ vermeiden ihn alle
- Das Signifikanzniveau ist

$$1 - (1 - \alpha)^6$$

- Für $\alpha = 0.05$ erhält man

$$1 - 0.95^6 = 0.2649$$

- Die Bonferroni-Korrektur löst dieses Problem im Prinzip, aber auf Kosten der Power

Quellen der Variabilität: zufällige Effekte

Messfehler: Die Körpergröße einer Person wird fünfmal gemessen. Die Ergebnisse werden voneinander abweichen. Durch Sorgfalt und geeignete Messmethoden kann der Messfehler beeinflusst werden. Ganz auszuschalten ist er nicht.

Natürliche Variabilität: Innerhalb einer Altersgruppe sind die Probanden unterschiedlich groß. Diese Variabilität ist unvermeidlich.

Diese beiden Quellen der Variabilität machen die statistische Betrachtung überhaupt erst erforderlich. Man fasst sie häufig als zufällige Effekte zusammen.

Quellen der Variabilität: Gruppenunterschiede

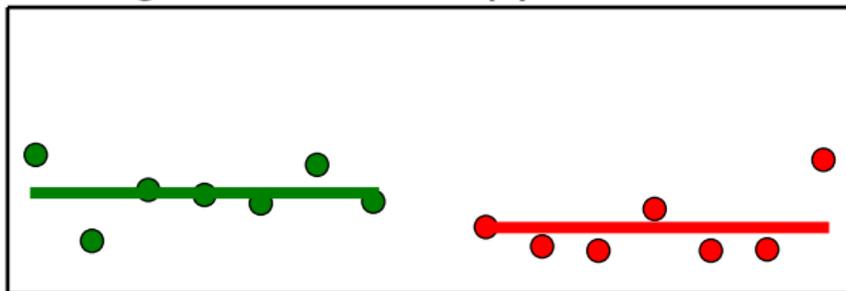
- Ziel ist die Untersuchung des Einflusses eines *Faktors* auf das Messergebnis. Beispiele für Faktoren sind etwa Alter, Unterrichtsmethode, Sonneneinstrahlung.
- Der Einfachheit halber gehen wir von nur einem Faktor aus.
- Der Faktor besitzt endlich viele Faktorstufen, in den Beispielen etwa
 - Alter: 10–14, 15–19, 20–24, ≥ 25 Jahre
 - Unterrichtsmethode: Methoden 1 bis 4
 - Sonneneinstrahlung: sonnig, halbschattig, schattig
- Zu jeder Faktorstufe werden mehrere Individuen ausgewählt und zu einer *Gruppe* zusammengefasst.
- Falls es Unterschiede zwischen den Gruppen gibt, so tragen sie zur Varianz bei.

Idee der Varianzanalyse

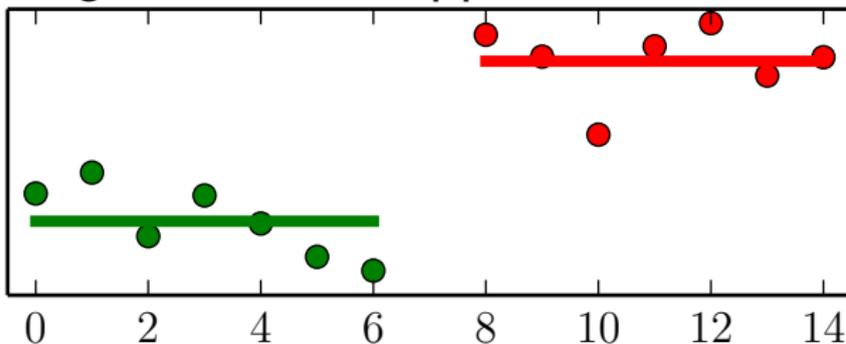
- Es ist rechnerisch möglich, die empirische Varianz aufzuteilen in den von den zufälligen Effekten und den von den Gruppenunterschieden verursachten Anteil.
- Dabei geht man davon aus, dass die zufälligen Effekte in allen Gruppen gleich wirken.
- Falls sich die empirische Varianz der zufälligen Effekte signifikant von dem durch die Gruppenunterschiede verursachten Anteil unterscheidet, dann ist der Einfluss des Faktors auf den Messwert nachgewiesen.
- Zum Vergleich dieser beiden Varianzen wird ein F -Test eingesetzt.

Idee der Varianzanalyse

Kein signifikanter Gruppenunterschied



Signifikanter Gruppenunterschied



Gruppenmittelwerte

- Es gibt k Faktorstufen und zu jeder dieser Faktorstufen eine Gruppe.
- Die j -te Gruppe hat n_j Elemente, welche mit

$$x_{j,1}, x_{j,2}, \dots, x_{j,n_j}$$

bezeichnet werden.

- Das arithmetische Mittel der Daten in der j -ten Gruppe ist

$$\bar{x}_j = \frac{1}{n_j} \sum_{i=1}^{n_j} x_{j,i}$$

Es ist der *Gruppenmittelwert*.

Gruppenmittelwerte im Beispiel

	Unterrichtsmethode			
	1	2	3	4
	$x_{1,1} = 2$	$x_{2,1} = 3$	$x_{3,1} = 6$	$x_{4,1} = 5$
	$x_{1,2} = 1$	$x_{2,2} = 4$	$x_{3,2} = 8$	$x_{4,2} = 5$
	$x_{1,3} = 3$	$x_{2,3} = 3$	$x_{3,3} = 7$	$x_{4,3} = 5$
	$x_{1,4} = 1$	$x_{2,4} = 5$	$x_{3,4} = 6$	$x_{4,4} = 3$
		$x_{2,5} = 0$	$x_{3,5} = 8$	$x_{4,5} = 2$
	$\bar{x}_1 = 1.75$	$\bar{x}_2 = 3.00$	$\bar{x}_3 = 7.00$	$\bar{x}_4 = 4.00$

Gruppengrößen: $n_1 = 4, n_2 = n_3 = n_4 = 5$

Gesamtmittelwert

- Der Gesamtstichprobenumfang ist $N = n_1 + \dots + n_k$, im Beispiel also $N = 19$
- Der Gesamtmittelwert ist der Mittelwert über alle Daten

$$\bar{\bar{x}} = \frac{1}{N} \sum_{j=1}^k \sum_{i=1}^{n_j} x_{j,i}$$

- Man kann ihn als gewichtetes arithmetisches Mittel der Gruppenmittelwerte berechnen

$$\bar{\bar{x}} = \frac{1}{N} \sum_{j=1}^k n_j \cdot \bar{x}_j$$

- Im Beispiel

$$\bar{\bar{x}} = \frac{1}{19} (4 \cdot 1.75 + 5 \cdot 3.00 + 5 \cdot 7.00 + 5 \cdot 4.00) = \frac{77}{19} = 4.053$$

Zerlegung der Varianz

- Für jeden Gruppenmittelwert \bar{x}_j sei a_j die Differenz zwischen \bar{x}_j und $\bar{\bar{x}}$, also

$$\bar{x}_j = \bar{\bar{x}} + a_j$$

- Für jeden einzelnen Messwert $x_{j,i}$ sei $e_{j,i}$ die Differenz zwischen $x_{j,i}$ und \bar{x}_j , also

$$x_{j,i} = \bar{x}_j + e_{j,i} = \bar{\bar{x}} + a_j + e_{j,i}$$

Gruppenmittelwerte im Beispiel

j	Unterrichtsmethode			
	1	2	3	4
$x_{j,1}$	1.75+0.25	3.00+0.00	7.00-1.00	4.00+1.00
$x_{j,2}$	1.75-0.75	3.00+1.00	7.00+1.00	4.00+1.00
$x_{j,3}$	1.75+1.25	3.00+0.00	7.00+0.00	4.00+1.00
$x_{j,4}$	1.75-0.75	3.00+2.00	7.00-1.00	4.00-1.00
$x_{j,5}$		3.00-3.00	7.00+1.00	4.00-2.00
\bar{x}_j	4.05-2.30	4.05-1.05	4.05+2.95	4.05-0.05

$$\bar{\bar{x}} = 4.05$$

Berechnung der Quadratsumme

- Für die bei der Berechnung der empirischen Varianz auftauchenden Differenzen bedeutet das

$$x_{j,i} - \bar{\bar{x}} = a_j + e_{j,i} = (\bar{x}_j - \bar{\bar{x}}) + (x_{j,i} - \bar{x}_j)$$

- Zur Berechnung des Quadrats ziehen wir die zweite binomische Formel heran

$$(x_{j,i} - \bar{\bar{x}})^2 = (\bar{\bar{x}} - \bar{x}_j)^2 + 2(\bar{\bar{x}} - \bar{x}_j) \cdot (\bar{x}_j - x_{j,i}) + (\bar{x}_j - x_{j,i})^2$$

Zerlegung der Varianz, Fortsetzung

Das muss zuerst über alle i und dann noch über alle j summiert werden. Wir fangen mit der Summe über alle i an. Dabei ist j die Nummer einer festen Gruppe

$$\begin{aligned}\sum_{i=1}^{n_j} (x_{j,i} - \bar{x})^2 &= n_j \cdot (\bar{x} - \bar{x}_j)^2 + 2(\bar{x} - \bar{x}_j) \cdot \underbrace{\sum_{i=1}^{n_j} (\bar{x}_j - x_{j,i})}_{=0} \\ &\quad + \sum_{i=1}^{n_j} (\bar{x}_j - x_{j,i})^2 \\ &= n_j \cdot (\bar{x} - \bar{x}_j)^2 + \sum_{i=1}^{n_j} (\bar{x}_j - x_{j,i})^2\end{aligned}$$

Zerlegung der Varianz, Fortsetzung

Aufsummiert über alle j

$$\sum_{j=1}^k \sum_{i=1}^{n_j} (x_{j,i} - \bar{x})^2 = \sum_{j=1}^k n_j \cdot (\bar{x} - \bar{x}_j)^2 + \sum_{j=1}^k \sum_{i=1}^{n_j} (\bar{x}_j - x_{j,i})^2$$

Abkürzungen

Quadratsumme, welche die Gesamtvariabilität repräsentiert:

$$\text{SQT} = \sum_{j=1}^k \sum_{i=1}^{n_j} (x_{j,i} - \bar{\bar{x}})^2$$

Quadratsumme, welche die Variabilität zwischen den Faktor-Stufen repräsentiert:

$$\text{SQZ} = \sum_{j=1}^k n_j \cdot (\bar{\bar{x}} - \bar{x}_j)^2$$

Quadratsumme, welche den Anteil der zufälligen Effekte repräsentiert:

$$\text{SQI} = \sum_{j=1}^k \sum_{i=1}^{n_j} (\bar{x}_j - x_{j,i})^2$$

Freiheitsgrade

- Abgekürzt lautet die gefundene Gleichung

$$SQT = SQZ + SQI$$

- Wenn der Faktor keinen Einfluss hat, dann kann man aus jeder dieser drei Größen die Varianz schätzen. Man muss nur die Zahl der Freiheitsgrade kennen.
- Der Gesamtversuch hat $N - 1$ Freiheitsgrade,

$$MQT = \frac{SQT}{N - 1}$$

ist die empirische Varianz.

Freiheitsgrade, Fortsetzung

- In **SQI** stecken k Schätzer, nämlich die Gruppenmittelwerte. Deswegen ist

$$\text{MQI} = \frac{\text{SQI}}{N - k}$$

ebenfalls ein Schätzer für die Varianz, wenn die Gruppe keinen Einfluss auf die Verteilung hat.

- $k - 1$ Freiheitsgrade bleiben übrig

$$\text{MQZ} = \frac{\text{SQZ}}{k - 1}$$

ist ebenfalls ein Schätzer für die Varianz.

Entscheidungsregel

- Die Teststatistik ist

$$t = \frac{\text{MQZ}}{\text{MQI}}$$

- Die Nullhypothese, dass der Faktor keinen Einfluss besitzt, wird abgelehnt, wenn

$$t > f_{k-1, N-k, 1-\alpha}$$

- Hierbei ist $f_{k-1, N-k, 1-\alpha}$ ein Quantil der F -Verteilung.

Anova, Zusammenfassung

- Gegeben k -Gruppen von Messwerten

$$X_{1,1}, X_{1,2}, \dots, X_{1,n_1}$$

$$X_{2,1}, X_{2,2}, \dots, X_{2,n_2}$$

...

$$X_{k,1}, X_{k,2}, \dots, X_{k,n_k}$$

- Dann ist $N = n_1 + \dots + n_k$ der Gesamtstichprobenumfang.
- Verteilungsvoraussetzungen: Die j -te Gruppe von Messwerten ist verteilt gemäß $N(\mu_j, \sigma^2)$ für unbekannte Werte μ_j und σ . Dabei hängt σ nicht von der Gruppe ab.
- Ziel: Die μ_j sollen miteinander verglichen werden.

Anova, Fortsetzung

- $x_{1,1}$ usw. seien die Realisierungen.
- Bestimme arithmetische Mittelwerte und Schätzer für die Varianzen

$$\bar{x}_j = \frac{1}{n_j} \sum_{i=1}^{n_j} x_{j,i}$$

$$\bar{\bar{x}} = \frac{1}{N} \sum_{j=1}^k n_j \cdot \bar{x}_j$$

$$\text{SQZ} = \sum_{j=1}^k n_j \cdot (\bar{\bar{x}} - \bar{x}_j)^2$$

$$\text{MQZ} = \frac{\text{SQZ}}{k-1}$$

$$\text{SQI} = \sum_{j=1}^k \sum_{i=1}^{n_j} (\bar{x}_j - x_{j,i})^2$$

$$\text{MQI} = \frac{\text{SQI}}{N-k}$$

- Die Teststatistik ist

$$t = \frac{\text{MQZ}}{\text{MQI}}$$

Anova, Fortsetzung

- Die Nullhypothese ist $H_0 = \{\mu_1 = \mu_2 = \dots = \mu_k\}$.
- Das Signifikanzniveau sei α
- Das folgende Quantil einer F -Verteilung wird benötigt

$$f_{k-1, N-k, 1-\alpha}$$

- Entscheidung: Die Nullhypothese wird abgelehnt, wenn

$$t > f_{k-1, N-k, 1-\alpha}$$

Beispiel Unterrichtsmethode

- $SQT = 98.04$, also $MQT = \frac{SQT}{19} = 5.208$
- Die Streuung über alle Daten ist 2.282
- $SQZ = 28.75$, also $MQZ = \frac{SQZ}{3} = 9.583$
- $SQI = 70.20$, also $MQI = \frac{SQI}{15} = 4.680$
- Der zufällige Streuungsanteil ist 2.163

Entscheidung im Beispiel Unterrichtsmethoden

- Im Beispiel ist die Teststatistik gleich

$$t = \frac{\text{MQZ}}{\text{MQI}} = \frac{9.583}{4.680} = 2.048$$

- Das Quantil ist $f_{3, 15, 0.95} = 3.287$
- Zum Signifikanzniveau $\alpha = 0.05$ kann nicht nachgewiesen werden, dass die Unterrichtsmethode einen Einfluss auf den Lernerfolg hat

Beispiel: Huflattich

- Die Aufnahme von Mg-Ionen wird bei Huflattich untersucht. Dazu werden je 6 Pflanzen in drei Nährlösungen mit gleicher Konzentration an Mg-Ionen, aber unterschiedlichen Konzentrationen an K- und Ca-Ionen herangezogen.
- Der Mg-Ionengehalt in den Blättern ist der folgende

	Nährlösung		
	1	2	3
	208	184	182
	175	161	193
	196	155	166
	181	185	145
	201	203	135
	166	166	151
n_j	6	6	6
\bar{x}_j	187.8	175.7	162.0

Beispiel, Fortsetzung

Wir wollen zum Signifikanzniveau $\alpha = 0.05$ prüfen, ob die Konzentration an K- und Ca-Ionen in der Nährlösung den Mg-Ionengehalt in den Blättern beeinflusst.

$$\bar{\bar{x}} = \frac{1}{18}(6 \cdot 187.8 + 6 \cdot 175.7 + 6 \cdot 162.0) = 175.2$$

$$\begin{aligned} \text{SQZ} &= 6 \cdot (175.2 - 187.8)^2 + 6 \cdot (175.2 - 175.7)^2 + 6 \cdot (175.2 - 162.0)^2 \\ &= 2004 \end{aligned}$$

$$\text{SQI} = 5490$$

$$\text{MQZ} = \frac{2004}{2} = 1002$$

$$\text{MQI} = \frac{5490}{15} = 366.0$$

$$t = \frac{1002}{366.0} = 2.738$$

Beispiel, Fortsetzung

- Folgendes Quantil einer F -Verteilung wird benötigt

$$f_{2,15,0.95} = 3.68$$

- Der Wert der Teststatistik ist $t = 2.738$. Die Nullhypothese kann nicht abgelehnt werden.
- Andererseits ist der p -Wert ungefähr 0.10
- Daher unterstützt der Test weder die Nullhypothese noch die Alternative. Die Ausgangsfrage bleibt offen.
- Die Biologin und der Biologe müssen entweder den Stichprobenumfang erhöhen oder ein besseres Verständnis der zu Grunde liegenden biochemischen Vorgänge entwickeln, um damit einen präziseren Test zu konzipieren.