

Statistische Methoden der Datenanalyse

Markus Schumacher, Stan Lai, Florian Kiss

Übung XI

21.1.2014, 22.1.2014

Anwesenheitsaufgaben

Aufgabe 53 Vergleich von Messungen einer gaussverteilten Variablen

Betrachten Sie den Fall, Sie hätten einen Satz von N Messungen einer gaussverteilten Variablen $\vec{x} = (x_1, x_2, \dots, x_N)$ aufgenommen, wobei \vec{x} gemäß $f_G(x; \mu_0, \sigma_0)$ verteilt sei. In dem vorliegenden Beispiel sollen Sie zwei verschiedene Hypothesentests betrachten, um sowohl den Mittelwert als auch Varianz Ihrer Messungen mit der erwarteten Verteilung $f_G(x; \mu_0, \sigma_0)$ zu vergleichen.

Das Makro `/home/slai/StatisticsCourse/PS11/aufgabe_53_anfang.C` beinhaltet Code, welcher einen Satz von M Experimenten generiert, jeweils mit N Messungen einer gaussverteilten Variablen. Jede Messung wird in ein Histogramm gefüllt, welches am Ende angezeigt wird.

- (i) Vergleichen Sie zuerst den Mittelwert der generierten Messdaten mit der Gaussverteilung, welche Sie dazu verwendeten, die Messungen zu erstellen.
- a) Nehmen Sie an, Sie kennen den Mittelwert μ wie auch σ der Gaussfunktion, die zum Generieren der Daten benutzt wurde. Um zu prüfen, ob Ihre Daten den Mittelwert $\mu = \mu_0$ besitzen, berechnen Sie für jedes Experiment die Teststatistik

$$t = \frac{\bar{x} - \mu_0}{\sigma/\sqrt{N}}$$

und füllen Sie diese in ein bereitgestelltes Histogramm. Diese Variable sollte nach der gaussischen WDF $f_G(t; 0, 1)$ verteilt sein. Überzeugen Sie sich davon, indem Sie die Methode

```
hist.Fit("gaus");
```

verwenden, um eine Gaussverteilung an Ihr Histogramm von t anzupassen. Stimmen die angepassten Parameter mit der Erwartung überein?

- b) Nehmen Sie nun an, Sie würden lediglich den Mittelwert μ , jedoch nicht die Breite der den Messungen zugrunde liegenden Gaussverteilung kennen. Folglich prüfen Sie, ob Ihre Daten den Mittelwert $\mu = \mu_0$ besitzen, indem Sie für jedes Experiment die Teststatistik

$$t' = \frac{\bar{x} - \mu_0}{s/\sqrt{N}}$$

berechnen, wobei die Standardabweichung s gegeben ist durch

$$s^2 = \frac{1}{N-1} \sum (x_i - \bar{x})^2 = \frac{N}{N-1} (\overline{x^2} - \bar{x}^2)$$

Füllen Sie Ihre Werte von t' in ein Histogramm. Die Variable t' sollte entsprechend einer Studentischen t -Verteilung $f_t(t; N-1)$ mit $N-1$ Freiheitsgraden verteilt sein. Überzeugen Sie sich davon, dass dies der Fall ist, indem Sie eine Studentische t -Verteilung an Ihr Histogramm anpassen. Verwenden Sie

```
TF1 tFit=TF1("tFit", "[1]*TMath::Student(x, [0])", -50., 50.);
```

um eine Anpassungsfunktion in Form einer Studentischen t -Verteilung bereitzustellen. Der `[0]`-te Parameter steht für die Anzahl der Freiheitsgrade. Passen Sie diese Funktion an Ihr Histogramm von t' an (`thist.Fit("tFit");`).

- (ii) Als nächstes vergleichen Sie die Breite der generierten Daten mit der der Gaussverteilung, die zum Generieren der Messungen benutzt wurde.

Gehen Sie davon aus, dass Ihnen der Mittelwert μ , jedoch nicht die Breite der Gaussverteilung bekannt ist. Um nun zu prüfen, ob die Daten die Breite $\sigma = \sigma_0$ aufweisen, berechnen Sie für jedes Experiment die Teststatistik

$$t'' = \frac{(N-1)s^2}{\sigma_0^2}$$

und füllen Sie diese in ein weiteres Histogramm. Die Variable t'' sollte nach einer χ^2 WDF $f_{\chi^2}(t''; N-1)$ mit $N-1$ Freiheitsgraden verteilt sein. Überzeugen Sie sich davon, dass dies der Fall ist, indem Sie eine χ^2 -Verteilung an Ihr Histogramm anpassen. Verwenden Sie

```
TF1 tChi2Fit = TF1("tChi2Fit", "[0]*(1.0/(TMath::Power(2,[1]/2.0)
*TMath::Gamma([1]/2.0))
*TMath::Power(x,([1]/2.0)-1.0)
*TMath::Exp(-x/2.0)", 0., 50.);
```

um eine χ^2 -Verteilung zur Anpassung bereitzustellen. Der [0]-te Parameter steht für den Normierungsfaktor und der [1]-te für die Anzahl an Freiheitsgraden. Passen Sie diese Funktion mittels `tthist.Fit("tChi2Fit");` an Ihr Histogramm von t'' an.

- (iii) Wie verändern sich die Verteilungen von t , t' und t'' , wenn Sie einen systematischen Fehler hinzufügen, der alle Messungen um einen Wert von 1 erhöht, in der Art

$$\vec{x} \rightarrow \vec{x}' = (x_1 + 1, x_2 + 1, \dots, x_N + 1)?$$

- (iv) Wie verändern sich die Verteilungen von t , t' und t'' , wenn Sie einen gaussischen Fehler mit Standardabweichung $\sigma = 0.1$ als zusätzliche Verschmierung zu allen Messungen hinzufügen?

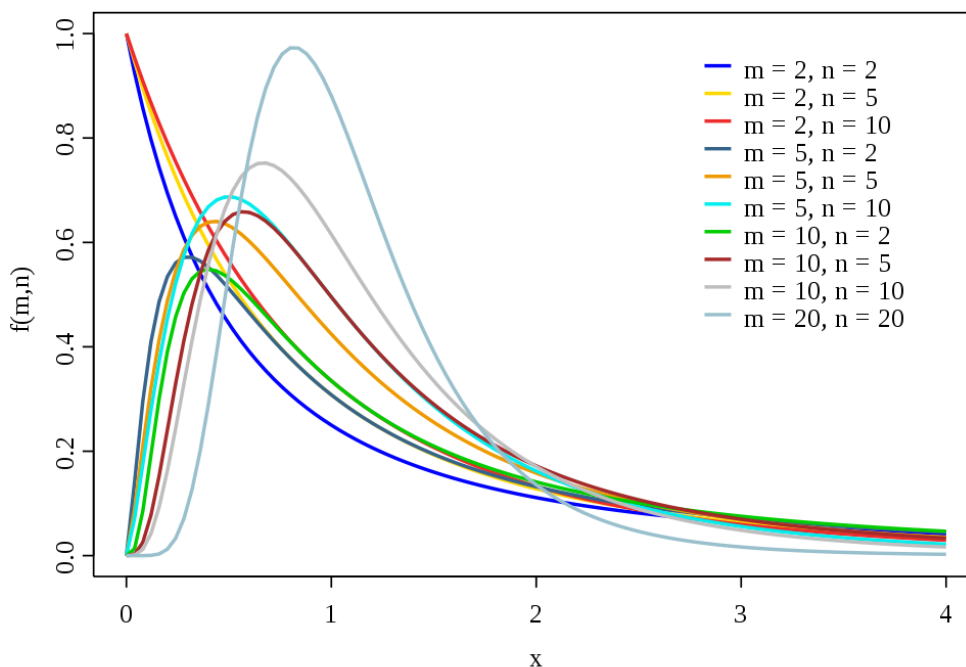


Abbildung 1: Die F-Verteilung

Hausaufgaben

Aufgabe 54 *Studentsche t -Verteilung - Teil II*

4 Punkte

Die Studentsche t -Verteilung kann dazu benutzt werden, um auf einem Datensatz eine Nullhypothese H_0 zu testen.

Gegeben sei eine Stichprobe vom Umfang n aus einer Gaussverteilung $N(\mu, \sigma^2)$. Falls σ bekannt ist, ist die Verteilung für

$$t = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}} \quad (1)$$

eine Gaussverteilung $N(0,1)$. Wenn σ^2 jedoch nicht bekannt ist, dann ist t gegeben durch:

$$t = \frac{\bar{x} - \mu}{s/\sqrt{n}} \quad (2)$$

mit der Stichprobenvarianz $s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$. In diesem Fall ist t nach der Studentschen t -Verteilung mit $n - 1$ Freiheitsgraden verteilt.

Betrachten Sie als Beispiel die Messung eines monoenergetischen Strahls von Teilchen mit Impuls $P_0 = 24.90 \text{ GeV}/c$. Dieser trifft auf eine Blasenkammer und durch Messung der Krümmung entlang der Teilchen spur wird der inverse Impuls $1/P_i$ bestimmt. Nehmen Sie an, dass $1/P$ für 20 Teilchen durch zwei verschiedene Detektoren A und B mit den Ergebnissen $1/P_A = (40.12 \pm 0.46) \times 10^{-3} (\text{GeV}/c)^{-1}$ und $1/P_B = (40.25 \pm 0.25) \times 10^{-3} (\text{GeV}/c)^{-1}$ gemessen wurde.

Um zu testen, ob beide Messungen mit der Bestimmung des inversen Impulses der einfallenden Teilchen, $1/P_0$, konsistent sind, sollten Sie diese beiden Hypothesen betrachten:

$$H_0 : \frac{1}{P_i} = \frac{1}{P_0}$$
$$H_1 : \frac{1}{P_i} \neq \frac{1}{P_0}$$

- (i) Was sind, unter Hinzunahme von Gleichung 2, die Werte von t für beide Messungen?
- (ii) Wie viele Freiheitsgrade hat jede Messung?
- (iii) Nutzen Sie die zur Verfügung gestellte Tabelle 3, um die Grenze der kritischen Region mit einer Signifikanz von $\alpha = 0.05$ zu finden. Bedenken Sie hierbei, dass Sie einen beidseitigen Test durchführen. Wieso muss dieser Test auf zwei Seiten durchgeführt werden?
- (iv) In Bezug auf den inversen Impuls der einfallenden Teilchen: Sind beide Messungen damit konsistent?

Aufgabe 55 *F-Test für zwei verschiedene Messungen*

5 Punkte

Betrachten Sie zwei unabhängige Chi-Quadrat verteilte Variablen, u_1 und u_2 , mit ν_1 und ν_2 Freiheitsgraden, d.h. u_1 ist gemäß $\chi^2(\nu_1)$ verteilt und u_2 nach gemäß $\chi^2(\nu_2)$. Dann ist die Variable F , definiert durch

$$F \equiv \frac{u_1/\nu_1}{u_2/\nu_2} \quad 0 \leq F \leq \infty; \nu_1, \nu_2 > 0 \quad (3)$$

verteilt nach der WDF

$$f(F; \nu_1, \nu_2) = \frac{\Gamma\left(\frac{1}{2}(\nu_1 + \nu_2)\right)}{\Gamma\left(\frac{1}{2}\nu_1\right)\Gamma\left(\frac{1}{2}\nu_2\right)} \left(\frac{\nu_1}{\nu_2}\right)^{\frac{1}{2}\nu_1} \frac{F^{\frac{1}{2}\nu_1 - 1}}{\left(1 + \frac{\nu_1 F}{\nu_2}\right)^{\frac{1}{2}(\nu_1 + \nu_2)}} \quad (4)$$

welche ‘‘F-Verteilung f#ur (ν_1, ν_2) Freiheitsgrade’’ genannt wird (siehe Abb. 1).

F#ur zwei Datens#atze x_1, x_2, \dots, x_n , gaussverteilt nach $N(\mu_1, \sigma_1^2)$, und y_1, y_2, \dots, y_n , gaussverteilt nach $N(\mu_2, \sigma_2^2)$, wobei die Mittelwerte μ_1 und μ_2 beider Verteilungen bekannt sind, ist die Gr#o#e

$$F = \frac{s_1}{s_2} = \frac{\frac{1}{n-1} \sum_{i=1}^n (x_i - \mu_1)^2}{\frac{1}{m-1} \sum_{i=1}^m (y_i - \mu_2)^2} \quad (5)$$

durch die F-Verteilung mit (n, m) Freiheitsgraden verteilt, wenn $\sigma_1^2 = \sigma_2^2$. Daher kann das Verh#altnis s_1/s_2 benutzt werden, um die Hypothese, dass beide Verteilungen die selbe Varianz ($H_0 : \sigma_1^2 = \sigma_2^2$) aufweisen, gegen die Hypothese, dass beide Varianzen verschieden sind ($H_1 : \sigma_1^2 > \sigma_2^2$), zu testen.

Kehren wir noch einmal zur Messung des Teilchenimpulses aus Aufgabe 54 zur#uck. Beide Messungen haben den selben Mittelwert $\mu_1 = \mu_2 = \mu_0$. Betrachten Sie hier nun folgende beiden Hypothesen f#ur die Varianzen des inversen Impulses:

$$H_0 : \frac{1}{\sigma_1^2} = \frac{1}{\sigma_2^2}$$

$$H_1 : \frac{1}{\sigma_1^2} < \frac{1}{\sigma_2^2}$$

- (i) Berechnen Sie den Wert von F f#ur beide Messungen.
- (ii) Wie viele Freiheitsgrade haben die Messungen?
- (iii) Was ist der kritische Wert von F bei einer Signifikanz von 5%? Sollte der Test auf einer oder auf zwei Seiten durchgef#uhrt werden? (Nutzen Sie Tabelle 2)
- (iv) Sind daher die Pr#azisionen der beiden Messungen miteinander konsistent?

TABLE 8.2
F DISTRIBUTION CRITICAL VALUES FOR 5% SIGNIFICANCE

	1	2	3	4	5	6	7	8	9	10	11	15	20	50	100
1	161	199	216	225	230	234	237	239	241	242	243	246	248	252	253
2	18.5	19.0	19.2	19.2	19.3	19.3	19.4	19.4	19.4	19.4	19.4	19.4	19.4	19.5	19.5
3	10.1	9.55	9.28	9.12	9.01	8.94	8.69	8.85	8.81	8.79	8.76	8.70	8.66	8.58	8.54
4	7.71	6.94	6.59	6.39	6.26	6.16	6.09	6.04	6.00	5.96	5.94	5.86	5.80	5.70	5.66
5	6.61	5.79	5.41	5.19	5.05	4.95	4.88	4.82	4.77	4.74	4.70	4.62	4.56	4.44	4.40
6	5.99	5.14	4.76	4.53	4.39	4.28	4.21	4.15	4.10	4.06	4.03	3.94	3.87	3.75	3.71
7	5.59	4.74	4.35	4.12	3.97	3.87	3.79	3.73	3.68	3.64	3.60	3.51	3.44	3.32	3.27
8	5.32	4.46	4.07	3.84	3.69	3.58	3.50	3.44	3.39	3.35	3.31	3.22	3.15	3.02	2.97
9	5.12	4.26	3.86	3.63	3.48	3.37	3.29	3.23	3.18	3.14	3.10	3.01	2.94	2.80	2.76
10	4.96	4.10	3.71	3.48	3.33	3.22	3.14	3.07	3.02	2.98	2.94	2.84	2.77	2.64	2.59
11	4.84	3.98	3.59	3.36	3.20	3.09	3.01	2.95	2.90	2.85	2.82	2.72	2.65	2.51	2.46
15	4.54	3.68	3.29	3.06	2.90	2.79	2.71	2.64	2.59	2.54	2.51	2.40	2.33	2.18	2.12
20	4.35	3.49	3.10	2.87	2.71	2.60	2.51	2.45	2.39	2.35	2.31	2.20	2.12	1.97	1.91
50	4.03	3.18	2.79	2.56	2.40	2.29	2.20	2.13	2.07	2.03	1.99	1.87	1.78	1.60	1.52
100	3.94	3.09	2.70	2.46	2.31	2.19	2.10	2.03	1.97	1.93	1.89	1.77	1.68	1.48	1.39

Abbildung 2: Kritische Werte der F-Verteilung an 5% Signifikanz

Aufgabe 56 *Zählexperiment für eine Signal- und Untergrundmessung***11 Punkte**

Betrachtet wird ein Experiment, dessen Ziel die Entdeckung eines neuen Teilchens oberhalb eines von momentanen Theorien vorhergesagten Untergrundes ist. Dabei könnte es sich beispielsweise um das Higgs-Boson oder auch um supersymmetrische Teilchen handeln, wobei die Untergrundvorhersage durch das Standardmodell der Teilchenphysik erfolgt.

Die zu betrachtenden Hypothesen sind also die Nullhypothese H_0 , dass nur Ereignisse aus Untergrundprozessen gemessen wurden, sowie die Alternativhypothese H_1 , dass sowohl Signal- als auch Untergründereignisse beobachtet wurden.

Im Experiment wurde eine Gesamtanzahl x von Ereignissen aufgezeichnet. Weiterhin wurde eine andere, signalfreie kinematische Region definiert, aus der man die Normierung des Untergrundes bestimmen kann. In dieser Region wurden y Ereignisse gefunden. Das Verhältnis der Untergründereignisse in der signalfreien Region zu denjenigen in der Signalregion sei τ . Die mittlere Anzahl der Untergründereignisse in der Signalregion sei b und die mittlere Anzahl der Signaleereignisse im Falle von Hypothese H_1 sei s .

- (i) Stellen Sie die Likelihoodfunktionen für die Hypothesen H_0 und H_1 auf. Nehmen Sie dabei an, dass die Gesamtanzahlen von Ereignissen in Signal- und Kontrollregion jeweils Poissonverteilt sind.
- (ii) Betrachten Sie jetzt die Schätzer für s und b unter der Hypothese H_1 (\hat{s} bzw. \hat{b}) sowie den Schätzer für b unter der Nullhypothese, $\hat{\hat{b}}$.
- a) Stellen Sie die Profile Likelihood λ auf.
- b) Bestimmen Sie Ausdrücke für \hat{s} , \hat{b} und $\hat{\hat{b}}$. Betrachten Sie dazu

$$\left. \frac{\partial L(H_0)}{\partial b} \right|_{\hat{\hat{b}}}, \quad (6)$$

sowie die die beiden gleichzeitigen Einschränkungen

$$\left. \frac{\partial L(H_1)}{\partial s} \right|_{\hat{s}} \quad \text{und} \quad \left. \frac{\partial L(H_1)}{\partial b} \right|_{\hat{b}}. \quad (7)$$

- c) Berechnen Sie $q = -2 \ln \lambda$ und zeigen Sie dadurch, dass die in der Vorlesung angegebene Gleichung

$$q = 2 \left[x \ln x + y \ln y - (x + y) \ln \left(\frac{x + y}{1 + \tau} \right) - y \ln \tau \right]$$

korrekt ist.

TABLE 7.2
 CRITICAL VALUES OF t
 For various values of confidence levels and n

Confidence (2 tailed)	60%	80%	90%	95%	98%	99%
(1 tailed)	80%	90%	95%	97.5%	99%	99.5%
$n = 1$	1.376	3.078	6.314	12.706	31.820	63.651
2	1.061	1.886	2.920	4.303	6.965	9.925
3	0.978	1.638	2.353	3.182	4.541	5.841
4	0.941	1.533	2.132	2.776	3.747	4.604
5	0.920	1.476	2.015	2.571	3.365	4.032
6	0.906	1.440	1.943	2.447	3.143	3.707
7	0.896	1.415	1.895	2.365	2.998	3.499
8	0.889	1.397	1.860	2.306	2.896	3.355
9	0.883	1.383	1.833	2.262	2.821	3.250
10	0.879	1.372	1.812	2.228	2.764	3.169
11	0.876	1.363	1.796	2.201	2.718	3.106
12	0.873	1.356	1.782	2.179	2.681	3.055
13	0.870	1.350	1.771	2.160	2.650	3.012
14	0.868	1.345	1.761	2.145	2.624	2.977
15	0.866	1.341	1.753	2.131	2.602	2.947
16	0.865	1.337	1.746	2.120	2.583	2.921
17	0.863	1.333	1.740	2.110	2.567	2.898
18	0.682	1.330	1.734	2.101	2.552	2.878
19	0.861	1.328	1.729	2.093	2.539	2.861
20	0.860	1.325	1.725	2.086	2.528	2.845
21	0.859	1.323	1.721	2.080	2.518	2.831
22	0.858	1.321	1.717	2.074	2.508	2.819
23	0.858	1.319	1.714	2.069	2.500	2.807
24	0.857	1.318	1.711	2.064	2.492	2.797
∞	0.842	1.282	1.645	1.960	2.326	2.576

Abbildung 3: Kritische Werte für t .