

Statistische Methoden der Datenanalyse

Markus Schumacher

Übung VIII

Markus Warsinsky

19.12.2011

Anwesenheitsaufgaben

Aufgabe 48 *Kleinste Quadrat Anpassung an $e^+e^- \rightarrow \mu^+\mu^-$ Monte Carlo Daten*

Diese Übung baut auf der in Aufgabe 43 durchgeführten Maximum-Likelihood-Anpassung an gebinnte Daten auf.

Der differentielle Wirkungsquerschnitt für $e^+e^- \rightarrow \mu^+\mu^-$ folgt der theoretischen Verteilung

$$\frac{d\sigma}{d(\cos\theta)d\phi} \propto 1 + \alpha \cos\theta + \beta \cos^2\theta, \quad (1)$$

wobei θ und ϕ für die Winkel des μ^+ Teilchens in Kugelkoordinaten stehen. Der Winkel θ wird relativ zur z -Achse und der Winkel ϕ relativ zur x -Achse gemessen. In der Übung sollen Sie einen 'Binned' Kleinste-Quadrat-Fit durchführen, um die Parameter zu finden, welche zur Erstellung der Monte Carlo Ereignisse benutzt wurden. Ein solcher simulierter Datensatz befindet sich in `/home/warsinsk/sd_wise1112/ueb7/MuMuEvGen.root`.

- (i) Am besten starten Sie von Ihrer Lösung zu Aufgabe 43. Berechnen Sie zusätzlich zum Wert der log-Likelihood für jedes durchgescannte Wertepaar (α_i, β_j) den Wert des χ^2 , indem Sie folgende Gleichung benutzen:

$$\chi^2(\alpha_i, \beta_j) = \sum_{k=1}^N \frac{(n_k - \nu_k(\alpha_i, \beta_j))^2}{\sigma_k^2}$$

wobei σ_k der erwartete Fehler auf die Anzahl der Einträge im Bin k ist und sich ergibt zu $\sqrt{\nu_k}$. Die Theorievorhersage ν_k ist gegeben durch:

$$\nu_k(\alpha_i, \beta_j) = n_{tot} \int_{x_k^{min}}^{x_k^{max}} f(x; \alpha_i, \beta_j) dx$$

Beachten Sie, dass im Gegensatz zur log-Likelihoodfunktion die Gesamtanzahl der Einträge des Histogramms bekannt sein muss, um eine Anpassung durchzuführen. Diese können Sie mittels der Methode `TH1::Integral()` ermitteln.

- (ii) Füllen Sie für jedes α_i und β_j den χ^2 -Wert in einen zweidimensionalen Graph. Dies funktioniert analog zu einem eindimensionalen Graphen:

```
TGraph2D chi2Graph = TGraph2D(likepoints);
```

wobei `likepoints` die Gesamtanzahl an Kombinationen von α und β ist, für die die Likelihood ausgewertet werden soll.

Jeder Punkt des Graphen kann gesetzt werden durch die `TGraph2D` Funktion:

```
SetPoint(int pointNumber, float alpha, float beta, float chi2)
```

Denken Sie daran, sich `pointNumber` auszurechnen (dies ist nicht die Zählvariable in den Schleifen über die Werte von α und β !).

- (iii) Zeichnen Sie den Graphen (χ^2 gegen α und β) unter Benutzung von

```
chi2Graph.Draw("colz");
```

wobei die Option "colz" ROOT die Anweisung gibt, einen Surface-Plot zu erstellen. Weitere Optionen des Draw-Befehls können im 'ROOT User's Guide' in der TGraph2D Sektion gefunden werden.

- (iv) Ermitteln Sie aus den erzeugten Graph eine Abschätzung für $\hat{\alpha}$ und $\hat{\beta}$ und deren Standardabweichungen $\sigma_{\hat{\alpha}}$ und $\sigma_{\hat{\beta}}$. Dazu könnten folgende Befehle nützlich sein:

```
float minimum=chi2Graph.GetHistogram().GetMinimum();
```

gibt das Minimum des Graphen zurück.

```
chi2Graph.GetHistogram().SetMaximum(minimum+1);
```

setzt für die Darstellung die maximale Zeichenebene auf das Minimum des Graphen plus 1.0.

- (v) Was ändert sich, wenn anstatt der erwarteten Fehler für σ_k die Fehler auf die Histogrammeinträge benutzt werden? Diese können mittels der Methode `TH1::GetBinError(int k)` ermittelt werden.
- (vi) Vergleichen Sie die Werte Ihres Fits mit denen des ROOT Fitting Paketes, indem Sie die Histogramm-Anpassungsfunktion

```
hist.Fit("functionName", "I");
```

benutzen, wobei "functionName" der Name der WDF Funktion ist und die Option "I" eine Kleinste Quadrat Anpassung ausführt, indem die gegebene Funktion über jedes Bin des Histogramms integriert wird. Beachten Sie, dass die für die Maximum Likelihood Anpassung verwendete Funktion auf Eins normiert ist (was in diesem Fall egal war). Definieren Sie sich deshalb eine zweite Funktion, die auf die Anzahl der Histogrammeinträge geteilt durch die Binbreite normiert ist. Sie können dazu einen weiteren Funktionsparameter einführen und diesen mit der Methode `FixParameter` fest setzen, so dass er in der Anpassung nicht variiert wird.

Aufgabe 49 *Erweiterter Maximum Likelihood Fit für eine Signal- und Untergrundverteilung*

In dieser Übung werden wir einen Beispieldatensatz betrachten, welcher aus zwei verschiedenen Arten von Ereignissen besteht: Signalereignisse, welche in einer Gaussverteilung $f_s(x)$ vorliegen und Untergrundereignisse, welche nach einer Exponentialfunktion $f_b(x)$ verteilt sind. Aufgabe ist es, die erweiterte Maximum Likelihood Methode zu benutzen, um die Anzahl von Signal- und Untergrundereignissen aus dem Beispieldatensatz zu errechnen.

Betrachten Sie die WDF $f(x; \mu_s, \mu_b)$ für die Signal- und Untergrundverteilung:

$$f(x; \mu_s, \mu_b) = \frac{\mu_s}{\mu_s + \mu_b} f_s(x) + \frac{\mu_b}{\mu_s + \mu_b} f_b(x) \quad (2)$$

wobei entsprechend μ_s und μ_b für die Anzahl von Signal- und Untergrundereignissen stehen. Um das Beispiel zu vereinfachen wird angenommen, dass die Parameter der Signal- und Untergrund-WDF bekannt sind. Daher sind f_s und f_b Funktionen nur einer Variablen x .

Falls die Gesamtzahl an Ereignissen poissonverteilt ist, ergibt sich für die WDF:

$$P(n; \mu_s, \mu_b) = \frac{(\mu_s + \mu_b)^n}{n!} \exp(-(\mu_s + \mu_b)) \quad (3)$$

Daher ist die log-Likelihoodfunktion gegeben durch

$$\ln \mathcal{L} = -(\mu_s + \mu_b) + \sum_{i=1}^n \ln[(\mu_s + \mu_b) f(x_i; \mu_s, \mu_b)] \quad (4)$$

Diese Funktion sollen Sie minimieren, um die Parameter der Signal- und Untergrundnormalisierung zu finden.

Zuerst muss ein Datensatz mit der Signal- und der Untergrund-WDF erzeugt werden. Dies wurde bereits in `/home/warsinsk/sd_wise1112/ueb7/aufgabe44_i.C` vorbereitet. Kopieren Sie sich dieses Skript in ein beliebiges Verzeichnis und lassen es laufen. Es erledigt die folgenden Aufgaben:

- (i) Eine TF1 Funktion, welche die Summe einer Gauss- und einer Exponentialfunktion darstellt wird bereitgestellt:

$$f(x; \mu, \sigma, \tau, \mu_s, \mu_b) = \frac{\mu_s}{\mu_s + \mu_b} \frac{1}{\sqrt{2\pi}\sigma^2} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right) + \frac{\mu_b}{\mu_s + \mu_b} \frac{1}{\tau} \exp\left(-\frac{x}{\tau}\right) \quad (5)$$

in dem Intervall $0.0 < x < 100.0$. Die Parameter werden auf $\mu = 5.0$, $\sigma = 0.5$ und $\tau = 10.0$ gesetzt.

- (ii) Die TF1-Funktion `GetRandom` wird benutzt, um einen Satz von 1000 Werten der obigen Verteilungsfunktion zu erhalten und diese in ein Array gefüllt. Ein Array ist eine feste Aneinanderreihung von (in diesem Fall) Fliesskommazahlen. Um auf das i -te Element zuzugreifen, wird einfach `data[i]` benutzt.
- (iii) Die Daten werden weiterhin in ein Histogramm gefüllt und graphisch dargestellt.

Als nächstes führen Sie nun selbst einen erweiterten Maximum Likelihood Fit auf dem erstellten Datensatz durch. Ergänzen Sie hierzu das Skript um folgende Komponenten (fügen Sie die Komponenten am Ende des Skriptes hinzu):

- (iv) Definieren sie eine weitere TF1 Funktion, wie oben, der Form $f(x; \mu, \sigma, \tau, \mu_s, \mu_b)$. Benutzen Sie die TF1 Funktion `TF1::FixParameter(Int_t parNum, Double_t value)`, um die Funktionsparameter auf die Werte $\mu = 2.0$, $\sigma = 0.5$ und $\tau = 10.0$ zu setzen (man könnte diese Parameter variabel lassen, jedoch müsste man dann einen fünfdimensionalen Fit durchführen!).
- (v) Erstellen Sie, analog zu Aufgabe 43, eine Schleife über mögliche Werte des Parameters μ_s . Definieren Sie innerhalb dieser Schleife eine weitere Schleife über mögliche Werte von μ_b . Sinnvolle Bereiche für μ_s und μ_b sind zwischen 450 und 550 mit 50 äquidistanten Schritten.
- (vi) Setzen Sie innerhalb der Schleife die Parameter der neuen TF1-Funktion auf die gerade aktuellen Werte von μ_s und μ_b mittels der `SetParameter(int parnumber, float parwert)`-Funktion.
- (vii) Erstellen Sie für jeden Wert von μ_s und μ_b eine Schleife über das Array mit den generierten Daten und berechnen Sie den aufsummierten log-Likelihood Wert. Hier kann es nützlich sein, mit der TF1 Funktion `Eval(Double_t x)` den Wert der Funktion am Punkt x zu berechnen.
- (viii) Befüllen Sie ein `TGraph2D` Objekt mit den aufsummierten log-Likelihood Werten für jeden Wert von μ_s und μ_b .
- (ix) Zeichnen Sie den Graphen und schätzen Sie damit die Parameter $\hat{\mu}_s$ und $\hat{\mu}_b$ und deren Standardabweichungen $\sigma_{\hat{\mu}_s}$ und $\sigma_{\hat{\mu}_b}$ ab.
- (x) Vergleichen Sie die Werte Ihres Fits mit denen des ROOT fitting Paketes (benutzen Sie `dataHist.Fit("fitFunk2", "LI");`). Definieren Sie sich zuvor `fitFunk2` als separate Funktion und fixieren Sie wieder $\mu = 2.0$, $\sigma = 0.5$ und $\tau = 10.0$. Stellen Sie diese Funktion wie folgt bereit:

```
TF1 fitFunk2 = TF1("fitFunk2",
  "([0]*1.0/TMath::Sqrt(2*TMath::Pi()*[2]*[2])
  *exp(-0.5*((x-[1])/[2])**2)
  + [3]*(1.0/[4])*exp(-x/[4]))", 0.0, 100.0);
fitFunk2.FixParameter(1, 5.0);
fitFunk2.FixParameter(2, 0.5);
fitFunk2.FixParameter(4, 10.0);
fitFunk2.SetParameter(0, 1000.0);
fitFunk2.SetParameter(3, 1000.0);
```

Beachten Sie, dass die Fitparameter 0 und 3 mit einem Faktor 10 zu multiplizieren sind, da die Binbreite nicht 1 ist.

Hausaufgaben

Aufgabe 50 *Kleinste Quadrate und Abschätzung der Normierung*

4 Punkte

Betrachten Sie eine Anpassung mittels der Methode der kleinsten Quadrate an ein Histogramm mit Einträgen y_i und Bins $i = 1, \dots, N$ und vorhergesagten Werten

$$\lambda_i(\vec{\theta}) = n \int_{x_i^{\min}}^{x_i^{\max}} f(x; \vec{\theta}) dx,$$

wobei die WDF $f(x; \vec{\theta})$ von unbekanntem Parametern $\vec{\theta}$ abhängt. Nehmen Sie an, dass man die Gesamtanzahl n ersetzt durch einen Parameter ν und diesen simultan mit den anderen Parametern anpasst, indem man die Größe

$$\chi^2(\vec{\theta}, \nu) = \sum_{i=1}^N \frac{(y_i - \lambda_i(\vec{\theta}, \nu))^2}{\sigma_i^2}$$

minimiert.

(i) Zeigen Sie, dass für die Wahl $\sigma_i^2 = \lambda_i$ sich der Schätzer

$$\hat{\nu}_{\text{KQ}} = n + \frac{\chi_{\min}^2}{2}$$

ergibt.

(ii) Zeigen Sie, dass sich für $\sigma_i^2 = y_i$ (modifizierte kleinste Quadrate) der Schätzer

$$\hat{\nu}_{\text{MKQ}} = n - \chi_{\min}^2$$

ergibt.

Aufgabe 51 *Unsicherheiten in der Abschätzung von Parametern*

4 Punkte

Die Lösungen für die Schätzer einer Anpassung auf der Basis der Methode der kleinsten Quadrate $\hat{\vec{\theta}}$ können in Matrixform geschrieben werden als

$$\hat{\vec{\theta}} = (A^T V^{-1} A)^{-1} A^T V^{-1} \vec{y}.$$

Zeigen Sie durch Anwenden der allgemeinen Formel für Fehlerfortpflanzung, dass die Kovarianzmatrix geschrieben werden kann als

$$V(\hat{\vec{\theta}}) = (A^T V^{-1} A)^{-1}$$

mit $V = V(\vec{y})$.

Aufgabe 52 *Geradenanpassung mit Matrixmethoden*

6 Punkte

Betrachten Sie eine Stichprobe vom Umfang N von Messwerten $(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)$, die alle denselben Fehler σ auf die Messung von y_i haben sollen. Benutzen Sie die Matrixnotation der Methode der kleinsten Quadrate, um Schätzer $\hat{\theta}_0$ und $\hat{\theta}_1$ auf die Parameter einer anzupassenden Gerade der Form

$$\lambda_i = \theta_0 + \theta_1 x_i = \sum_{j=0}^n a_j(x_i) \theta_j = \sum_{j=0}^n A_{ij} \theta_j$$

mit $A_{ij} = a_j(x_i)$ anzugeben.

(i) Schreiben Sie zunächst die Matrix A und die Kovarianzmatrix V auf.

(ii) Berechnen Sie die Matrix für die Schätzer $\hat{\vec{\theta}}$, indem Sie die Matrizen $A^T V^{-1} A$ und $A^T V^{-1} \vec{y}$ berechnen.

- (iii) Schreiben Sie die Schätzer $\hat{\theta}_0$ und $\hat{\theta}_1$ in Abhängigkeit der Erwartungswerte und Varianzen von x , y , xy sowie der Kovarianz von x und y auf.
- (iv) Wie unterscheidet sich die Rechnung, wenn jeder Messwert einen nicht korrelierten Fehler σ_i auf y_i hat? Schreiben Sie die Matrizen für V , $A^T V^{-1} A$, $A^T V^{-1} y$ und somit $\hat{\theta}$ auf.

Aufgabe 53 *Methode der Kleinsten Quadrate mit Zwangsbedingungen: Winkelmessung im Dreieck* **6 Punkte**

In der Vorlesung wurde gezeigt, dass, um einen Satz von Messungen $\vec{y} = (y_1, y_2, \dots, y_N)$ durch eine lineare Funktion $A\vec{\theta}$ zu fitten, ein Fit mit der Methode der Linearen Kleinsten Quadrate durchgeführt werden kann. Dieser minimiert die Größe

$$\chi^2 = (\vec{y} - A\vec{\theta})^T V^{-1} (\vec{y} - A\vec{\theta}),$$

wobei V die Kovarianzmatrix der Messungen \vec{y} ist.

Weiterhin wurde gezeigt, dass unter Berücksichtigung eines Satzes von K Randbedingungen $\vec{b} = (b_1, b_2, \dots, b_K)$, welche die Gleichungen $B\vec{\theta} - \vec{b} = 0$ erfüllen, die Methode der kleinsten Quadrate durch Minimierung der Größe

$$\chi^2 = (\vec{y} - A\vec{\theta})^T V^{-1} (\vec{y} - A\vec{\theta}) + 2\vec{\lambda}^T (B\vec{\theta} - \vec{b})$$

verbessert werden kann, wobei $\vec{\lambda} = (\lambda_1, \lambda_2, \dots, \lambda_K)$ ein Vektor von Lagrange-Multiplikatoren ist. Folglich muss χ^2 minimiert werden in Bezug auf $\vec{\theta}$ und $\vec{\lambda}$. Die Lösung dieses Minimierungsproblems ergibt die Schätzer

$$\hat{\vec{\theta}} = C^{-1} \vec{c} - C^{-1} B^T V_B^{-1} (B C^{-1} \vec{c} - \vec{b}) = C^{-1} \vec{c} - C^{-1} B^T \hat{\vec{\lambda}}$$

wobei $C \equiv A^T V^{-1} A$, $\vec{c} \equiv A^T V^{-1} \vec{y}$ und $V_B \equiv B C^{-1} B^T$. Die Varianz errechnet sich daher zu

$$V[\hat{\vec{\theta}}] = C^{-1} - (B C^{-1})^T V_B^{-1} (B C^{-1}).$$

Die verbesserten Schätzer für die Messungen sind gegeben durch:

$$\hat{n} = A \hat{\vec{\theta}} = A \left[C^{-1} \vec{c} - C^{-1} B^T V_B^{-1} (B C^{-1} \vec{c} - \vec{b}) \right],$$

mit der Varianz

$$V[\hat{n}] = A V[\hat{\vec{\theta}}] A^T = A \left[C^{-1} - (B C^{-1})^T V_B^{-1} B C^{-1} \right] A^T.$$

- (i) Zeigen Sie, dass $E[\hat{\vec{\lambda}}] = 0$ und $E[\hat{\vec{\theta}}] = \vec{\theta}$. Daraus folgt, dass die Schätzungen der Parameter erwartungstreu sind.
- (ii) Betrachten Sie die Messung von drei Winkeln eines Dreiecks analog zu dem Beispiel aus der Vorlesung. Die drei Messungen sind gegeben durch $\vec{y} = (y_1, y_2, y_3)$ mit $\sigma_i = \sigma$ und sollen gefittet werden durch die Funktion $A\vec{\theta}$ mit $\vec{\theta} = (\theta_1, \theta_2, \theta_3)$, wobei A die Einheitsmatrix in drei Dimensionen ist.
- Was sind die Werte für B und \vec{b} ?
 - Berechnen Sie $\hat{\vec{\theta}}$ und $V[\hat{\vec{\theta}}]$ und somit \hat{n} und $V[\hat{n}]$.
 - Wie verbessert die Zwangsbedingung auf die Messungen die Schätzungen der gemessenen Werte der Dreieckswinkel?