Statistical Data Analysis Discussion notes – week 7

- Problem sheet 4
- Maximum Likelihood examples

#### Problem sheet 4

**Exercise 1(a)** [2 marks] Let  $x_N = \sum_{i=1}^N r_i$ , where the  $r_i$  are independent and uniformly distributed between 0 and 1. Find the mean  $\mu_N$  and standard deviation  $\sigma_N$  of  $x_N$  as a function of N.

1(a) [2 marks] The variable  $x_N$  is defined as

$$x_N = \sum_{i=1}^N r_i$$

where the  $r_i$  are independent and uniformly distributed on [0,1]. From the lectures (or rederive) we have the mean and variance  $E[r_i] = 1/2$  and  $V[r_i] = 1/12$ . The expectation value of  $x_N$  is therefore

$$E[x_N] \equiv \mu_N = E\left[\sum_{i=1}^N r_i\right] = \sum_{i=1}^N E[r_i] = \sum_{i=1}^N \frac{1}{2} = \frac{N}{2}$$

The variance of  $x_N$  is

$$V[x_N] = V\left[\sum_{i=1}^N r_i\right] = \sum_{i=1}^N V[r_i] = \sum_{i=1}^N \frac{1}{12} = \frac{N}{12}$$

and so its standard deviation is

$$\sigma[x_N] = \sqrt{\frac{N}{12}}$$

G. Cowan / RHUL Physics

1(b) [5 marks] Using the results from (a), construct the standardized variable

$$y_N = \frac{x_N - \mu_N}{\sigma_N} = \sqrt{\frac{12}{N}} \left( \sum_{i=1}^N r_i - \frac{N}{2} \right) \;.$$

Using the simpleMC program (either C++ or Python) from problem sheet 3 as a starting point, write a computer program to make histograms of 10000 values of  $y_N$  as defined above for N = 1, 2, 4, 12. Make sure to set the limits of the histogram such that the entire distribution is plotted.



Histograms of  $y_N$  for N = 1, 2, 4, 12. The means and standard deviations are indicated on the plots. As shown, they are close to  $\mu = 0$  and  $\sigma = 1$ , as must emerge by construction for the standardized variable  $y_N$ . **Exercise 2** Consider the pdf  $f(x) = 4x^3$ ,  $0 \le x \le 1$ .

**2(a)** [4 marks] Use the transformation method to find the function x(r) to generate random numbers according f(x). Implement the method in a short computer program and make a histogram with 10000 values.

2(a) [4 marks] Find the transformation x(r) to produce  $x \sim f(x)$ . First find the cumulative distribution



**2(b)** [4 marks] Write a program to generate random numbers according to f(x) using the acceptance-rejection technique. Plot a histogram of the results.



**Exercise 3** [5 marks] Suppose  $\vec{x} = (x_1, \ldots, x_n)$  follows an *n*-dimensional Gaussian distribution  $f(\vec{x}; \vec{\mu}, V)$  with  $\vec{\mu} = (\mu_1, \ldots, \mu_n)$  and covariance matrix  $V_{ij} = \operatorname{cov}[x_i, x_j]$ . (In the formulas below regard  $\vec{x}$  and  $\vec{\mu}$  to be column vectors.) Suppose we have two hypotheses for the vector of means,  $\vec{\mu}_0$  and  $\vec{\mu}_1$ , where for both one uses the same covariance matrix V, and consider the test statistic

$$t(\vec{x}) = \ln \frac{f(\vec{x}|\vec{\mu}_1)}{f(\vec{x}|\vec{\mu}_0)}$$

Show that this  $t(\vec{x})$  can be written in the form

$$t(\vec{x}) = w_0 + \sum_{i=1}^n w_i x_i \; ,$$

or equivalently  $t(\vec{x}) = w_0 + \vec{w}^T \vec{x}$ , where  $\vec{w}$  is a column vector of coefficients  $w_i, i = 1, ..., n$ .

$$i.l. f(\vec{x} | \vec{\mu}_{u}) = \frac{1}{(2\pi)^{n/2}} |V|^{\frac{1}{2}} \exp\left[-\frac{1}{2}(\vec{x} - \vec{\mu}_{u})^{T} V^{-1}(\vec{x} - \vec{\mu}_{u})\right]$$
  
Test statistic  $t(x) = \ln \frac{f(\vec{x} | \vec{\mu}_{u})}{f(\vec{x} | \vec{\mu}_{o})}$ 

G. Cowan / RHUL Physics

Ex. 3 (cont.)



Ex. 3 (cont.)



#### **Example of Maximum Likelihood**

Consider  $f(x) = (1+\theta)x^{\theta}$   $0 \le x \le 1$ with i.i.d. sample x1, ..., Xn  $L(\Theta) = \prod_{i=1}^{n} f(x_{i}; \Theta) = \prod_{i=1}^{n} (1+\Theta) x_{i}^{\Theta}$ iEI =)  $hL(0) = \sum_{i=1}^{n_i} \int h(1+0) + 0 h_i x_i^{n_i}$ To find MLE, set =  $\frac{\partial h L}{\partial \theta} = \frac{n}{1+\theta} + \sum_{i=1}^{\infty} h x_i$ h 25

G. Cowan / RHUL Physics

### MLE example (cont.)



Variance from asymptotic properties ~ 0 large n une 25 (1+ NI 0 F date X Inc 2 (1+0)2 lun L  $\frac{n}{(1+\theta)^2}$ E E ~ (1+0) 1+0 ->





## MLE for number of taxis

The number plate of taxis in every canton in Switzerland ends with a number N from 1 to  $N_{tot}$ , where  $N_{tot}$  is the total number of taxis.



Model the probability for observing plate number N with

$$P(N|N_{\text{tot}}) = \frac{1}{N_{\text{tot}}}, \quad 1 \le N \le N_{\text{tot}}$$

# MLE for $N_{\rm tot}$

Suppose you observe one taxi at random with plate number N.

The likelihood function is  $L(N_{\text{tot}}) = \frac{1}{N_{\text{tot}}}$ ,  $N_{\text{tot}} \ge N$ 

which is maximized for  $\widehat{N}_{tot} = N$ 

The expectation value and bias of the MLE are

$$E[\hat{N}_{\text{tot}}] = E[N] = \sum_{N=1}^{N_{\text{tot}}} \frac{N}{N_{\text{tot}}} = \frac{N_{\text{tot}} + 1}{2} \qquad b = \frac{1 - N_{\text{tot}}}{2}$$

For better estimators, see similar problem with tanks in WW2: https://en.wikipedia.org/wiki/German\_tank\_problem

E.g. the minimum-variance unbiased estimator is:  $\ \widehat{N}_{
m tot} = 2N-1$ 

## Cheap estimator for mass of W boson

The Particle Physics community has spent huge sums trying to estimate the mass of the W boson with the smallest possible statistical and systematic uncertainty.

Here is an estimator with zero statistical uncertainty. And it's free!

$$\widehat{M}_W = 80.4 \,\mathrm{GeV}$$

Here is its sampling distribution:

Does this violate the information inequality?



### Cheap estimator for mass of W boson (2)

This estimator's bias is  $b = E[\widehat{M}_W] - M_W = 80.4 \,\mathrm{GeV} - M_W$ 



Note best estimate of  $M_W$  is (in 2020) - 80.379±0.012 GeV, so the numerical  $M_W$  value of the bias may be fairly small.

But we have  $\frac{\partial b}{\partial M_W} = -1$  and so  $MVB = -\left(1 + \frac{\partial b}{\partial M_W}\right)^2 / E\left[\frac{\partial^2 \ln L}{\partial M_W^2}\right] = 0$ 

So the information inequality is still satisfied.

G. Cowan / RHUL Physics

### Extended ML example

Consider two types of events (e.g., signal and background) each of which predict a given pdf for the variable x:  $f_s(x)$  and  $f_b(x)$ .

We observe a mixture of the two event types, signal fraction = (, expected total number = v, observed total number = n.

Let  $\mu_{s} = \theta \nu$ ,  $\mu_{b} = (1 - \theta) \nu$ , goal is to estimate  $\int_{s}$ ,  $\int_{b}$ .

$$f(x; \mu_{s}, \mu_{b}) = \frac{\mu_{s}}{\mu_{s} + \mu_{b}} f_{s}(x) + \frac{\mu_{b}}{\mu_{s} + \mu_{b}} f_{b}(x)$$

$$P(n; \mu_{\rm S}, \mu_{\rm b}) = \frac{(\mu_{\rm S} + \mu_{\rm b})^n}{n!} e^{-(\mu_{\rm S} + \mu_{\rm b})}$$

$$\rightarrow \ln L(\mu_{s}, \mu_{b}) = -(\mu_{s} + \mu_{b}) + \sum_{i=1}^{n} \ln [(\mu_{s} + \mu_{b})f(x_{i}; \mu_{s}, \mu_{b})]$$

## Extended ML example (2)

Monte Carlo example with combination of exponential and Gaussian:

$$\mu_{\rm S} = 6$$
$$\mu_{\rm b} = 60$$

Maximize log-likelihood in terms of  $\mu_s$  and  $\mu_b$ :

$$\hat{\mu}_{s} = 8.7 \pm 5.5$$

 $\hat{\mu}_{b} = 54.3 \pm 8.8$ 



Here errors reflect total Poisson fluctuation as well as that in proportion of signal/background.