# S1: Principles of Data Science
# **Problem Sheet 3 Solutions**
## MPhil in Data Intensive Science

Matt Kenzie

mk652@cam.ac.uk

Michealmas Term 2023

## Problem Sheet 3

1. (a) Ensure p.d.f. normalises to unity:

$$\int_a^b N(1 + \alpha x + \beta x^2)dx = 1. \tag{1}$$

Therefore:

$$\frac{1}{N} = \int_a^b (1 + \alpha x + \beta x^2)dx \tag{2}$$

$$= \left[ x + \alpha \frac{x^2}{2} + \beta \frac{x^3}{3} \right]_a^b \tag{3}$$

$$= b + \alpha \frac{b^2}{2} + \beta \frac{b^3}{3} - a - \alpha \frac{a^2}{2} - \beta \frac{a^3}{3} \tag{4}$$

$$= (b - a) + \frac{\alpha}{2} \left( b^2 - a^2 \right) + \frac{\beta}{3} \left( b^3 - a^3 \right). \tag{5}$$

On reflection choosing $\alpha$ and $\beta$ as parameters, with $a$ and $b$ as limits is a poor choice that can easily lead to confusion. Alas, it is done.

(b) First two central moments are:

$$\mu_1 = x \tag{6}$$

$$\mu_2 = x^2. \tag{7}$$

Take the expectation values of these gives

$$E[\mu_1] = E[x] = \int_a^b Nx(1 + \alpha x + \beta x^2)dx \tag{8}$$

$$= N \int_a^b (x + \alpha x^2 + \beta x^3)dx \tag{9}$$

$$= N \left[ \frac{x^2}{2} + \alpha \frac{x^3}{3} + \beta \frac{x^4}{4} \right]_a^b \tag{10}$$

$$= N \left( \frac{1}{2}(b^2 - a^2) + \frac{\alpha}{3}(b^3 - a^3) + \frac{\beta}{4}(b^4 - a^4) \right) \tag{11}$$

$$= \frac{d_2 + \alpha d_3 + \beta d_4}{d_1 + \alpha d_2 + \beta d_3}. \tag{12}$$

$$E[\mu_2] = E[x^2] = \int_a^b Nx^2(1 + \alpha x + \beta x^2)dx \tag{13}$$

$$= N \int_a^b (x^2 + \alpha x^3 + \beta x^4)dx \tag{14}$$

$$= N \left[ \frac{x^3}{3} + \alpha \frac{x^4}{4} + \beta \frac{x^5}{5} \right]_a^b \tag{15}$$

$$= N \left( \frac{1}{3}(b^3 - a^3) + \frac{\alpha}{4}(b^4 - a^4) + \frac{\beta}{5}(b^5 - a^5) \right) \tag{16}$$

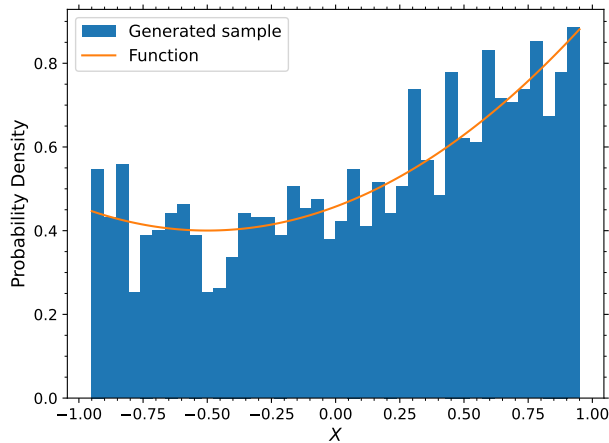$$= \frac{d_3 + \alpha d_4 + \beta d_5}{d_1 + \alpha d_2 + \beta d_3}. \tag{17}$$

$$\tag{18}$$

(c) To get estimations for $\hat\alpha$ and $\hat\beta$ we simply solve these equations for $\alpha$ and $\beta$ and plugin the sample estimates for the first and second central moments, $\hat\mu_1$ and $\hat\mu_2$. Solving the equations:

$$\hat\alpha = \frac{(\hat\mu_1 d_3 - d_4)(\hat\mu_2 d_1 - d_3) - (\hat\mu_1 d_1 - d_2)(\hat\mu_2 d_3 - d_5)}{(\hat\mu_1 d_2 - d_3)(\hat\mu_2 d_3 - d_5) - (\hat\mu_1 d_3 - d_4)(\hat\mu_2 d_2 - d_4)}, \tag{19}$$

$$\hat\beta = \frac{(\hat\mu_1 d_1 - d_2)(\hat\mu_2 d_2 - d_4) - (\hat\mu_1 d_2 - d_3)(\hat\mu_2 d_1 - d_3)}{(\hat\mu_1 d_2 - d_3)(\hat\mu_2 d_3 - d_5) - (\hat\mu_1 d_3 - d_4)(\hat\mu_2 d_2 - d_4)}. \tag{20}$$

(d) Accept-reject code we have already seen in a previous problem sheet. So just make a plot like this one

(e) Write some code similar to that which can be found in `method_of_moments.py` which computes MoM estimates and compares to MLE. I put this code below for reference

```python
## Solution to Problem Sheet 3 Question 1
## once again more optimal solutions will be available

import numpy as np
np.random.seed(210187)
from scipy.optimize import brute, minimize
from scipy.integrate import quad
from scipy.stats import moment
import matplotlib.pyplot as plt
plt.style.use('code/mphil.mplstyle')
from jacobi import propagate
from iminuit import cost, Minuit
import sys
sys.path.append('$PWD/code')
from accept_reject import accept_reject_1d, check_ok_plot
from tabulate import tabulate

if __name__=="__main__":

    xrange = [-0.95,0.95]
    alpha = 0.5
    beta = 0.5

    def dk(k):
        return (xrange[1]**k - xrange[0]**k)/k

    N = 1/( dk(1) + alpha*dk(2) + beta*dk(3) )

    func = lambda x: N*(1 + alpha*x + beta*x**2)

    # run the accept-reject
    dset = accept_reject_1d( func, xrange=xrange, size=2000 )
    check_ok_plot( func, xrange, dset, bins=40, save=f'figs/mom_gen.
        pdf' )
    dset = np.asarray(dset)
    np.save('../s1_principles_of_data_science/datasets/mom_data.npy'
        , dset)

    # MoM estimate
    m1_hat = moment(dset, moment=1, center=0)
    m2_hat = moment(dset, moment=2, center=0)

    def alpha(m1, m2):
        numerator = (m1*dk(3) - dk(4))*(m2*dk(1)-dk(3)) - (m1*dk(1)-
            dk(2))*(m2*dk(3)-dk(5))
        denominator = (m1*dk(2) - dk(3))*(m2*dk(3)-dk(5)) - (m1*dk(3
            )-dk(4))*(m2*dk(2)-dk(4))
        return numerator / denominator

    def beta(m1, m2):
        numerator = (m1*dk(1) - dk(2))*(m2*dk(2)-dk(4)) - (m1*dk(2)-
            dk(3))*(m2*dk(1)-dk(3))
        denominator = (m1*dk(2) - dk(3))*(m2*dk(3)-dk(5)) - (m1*dk(3
```

```
                )-dk(4))*(m2*dk(2)-dk(4))
49          return numerator / denominator
50
51      alpha_hat = alpha(m1_hat, m2_hat)
52      beta_hat = beta(m1_hat, m2_hat)
53
54      N = len(dset)
55      cov11 = 1/(N*(N-1)) * np.sum( (dset-m1_hat) * (dset-m1_hat) )
56      cov22 = 1/(N*(N-1)) * np.sum( (dset**2-m2_hat) * (dset**2-m2_hat
        ) )
57      cov12 = 1/(N*(N-1)) * np.sum( (dset-m1_hat) * (dset**2-m2_hat) )
58
59      cov_hat = np.array( [[cov11,cov12],[cov12,cov22]] )
60
61      # propagate to parameter estimates
62
63      transform = lambda p: np.array( [ alpha(*p), beta(*p) ] )
64
65      mom, mom_cov = propagate( transform, [m1_hat, m2_hat], cov_hat )
66
67      # Now some kind of MLE
68
69      def density(x, alpha, beta):
70          N = 1/( dk(1) + alpha*dk(2) + beta*dk(3) )
71          return N*(1 + alpha*x + beta*x**2)
72
73      nll = cost.UnbinnedNLL(dset, density)
74      mi = Minuit( nll, alpha=0.5, beta=0.5 )
75      mi.migrad()
76      mi.hesse()
77
78      print("---- METHOD OF MOMENTS ----")
79      print(tabulate(zip(['alpha','beta'], mom, np.diag(mom_cov)**0.5)
        , headers=['Parameter','Value','Error']))
80      print("---- MAXIMUM LIKELIHOOD ---")
81      print(tabulate(zip(['alpha','beta'], list(mi.values), list(mi.
        errors)), headers=['Parameter','Value','Error']))
82
83      plt.show()
```

2. There are 25 bins. A Gaussian has 2 shape parameters and there is one scale or "height" parameter for how much of the Gaussian there is. The background has no shape parameters (it is flat) but still has an overall height parameter. Thus the total is $25 - 3 - 1 = 21$.

3. The code for this solution is provided in `neutrinos.py`.

   (a) To compute the expectation of our Poisson then we can do

   $$\bar{\nu} = \frac{\sum w_i N_i}{\sum N_i} \tag{21}$$

   where $w_i$ is the number of events for an internal and $N_i$ is the number of intervals with that number of events.

   This results in a value of $\bar{\nu} = 0.777$.

4

I can now ask how many events I would expect to see from a Poisson distribution with this expectation and then multiply this by the total number of intervals I measured to give an extra column in my table:

| No. of events | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|---|
| Obs no. of intervals | 1042 | 860 | 307 | 78 | 15 | 3 | 0 | 0 | 0 | 1 |
| Pred no. of intervals | 1060.15 | 823.847 | 320.107 | 82.919 | 16.109 | 2.503 | 0.324 | 0.036 | 0.0035 | 0.0003 |

Then use this to calculate the $\chi^2$ and $p$-value. Including all data gives:

$$\chi^2 = 3313 \tag{22}$$

$$p\text{-value} = 0 \tag{23}$$

Clearly this is a dreadful fit.
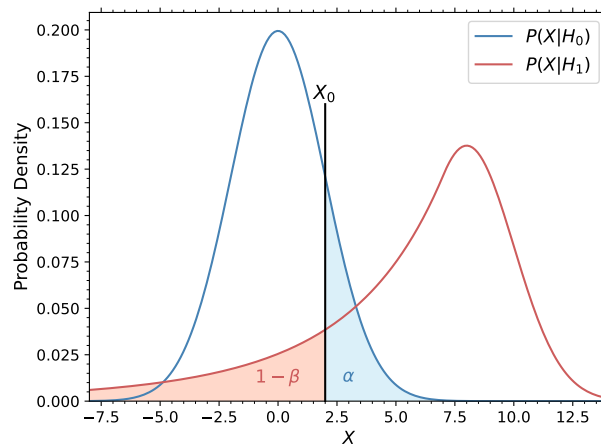
(b) Then omitting the data in the last bin gives

$$\chi^2 = 3.25 \tag{24}$$

$$p\text{-value} = 0.92 \tag{25}$$

Which is a very reasonable fit.

(c) The conclusion of course is that a simple Poisson only background cannot explain this occurence of an intervals countaining 9 neutrinos events. The conclusion of that is an observation has been made of some object which for a short period was emitting a high-neutrino flux. This was the supervnova S1987A.

4. I actually already do the proof of this in the lectures, but just to reproduce it here. We already saw this sort of plot in the lectures showing a test-statistic distribution under null and alternate hypotheses and then the regions $\alpha$ and $\beta$ which lie above the observation.

The regions $\alpha$ and $\beta$ are defined such that

$$\alpha = \int_{X_0}^{\infty} p(X; \theta_0) dX \tag{26}$$

$$\beta = \int_{X_0}^{\infty} p(X; \theta_1) dX \tag{27}$$

assuming that $\theta$ is some continuous parameter which differentiates between the null and alternate hypotheses (for example a number of signal). For a fixed $\alpha$ we need to find the test-statistic that will maximise $\beta$. Thus starting from the term of $\beta$ above:

$$\beta = \int_{X_0}^{\infty} p(X; \theta_1) dX \tag{28}$$

$$= \int_{X_0}^{\infty} \frac{p(X; \theta_1)}{p(X; \theta_0)} p(X; \theta_0) dX \qquad = E\left[\frac{p(X; \theta_1)}{p(X; \theta_0)}\right] \tag{29}$$

Thus we see that in order to maximise $\beta$ we have to maximise the expectation value on the right which is done with test-statistics of the form of the likelihood ratio.

5. A question inspired by measurements of the $CP$-violating phase $\gamma$ which appears in the CKM quark-mixing matrix

(a) A snippet of code which will perform the relevant stuff is in `polar.py` which I put below

```python
1  import numpy as np
2  from scipy.stats import multivariate_normal as mvn, chi2
3  from iminuit import Minuit
4  import matplotlib.pyplot as plt
5  plt.style.use('code/mphil.mplstyle')
6
7  names  = [ 'xp', 'yp', 'xm', 'ym' ]
8  titles = [ '$x^+$', '$y^+$', '$x^-$', '$y^-$' ]
9  values = np.array( [ -9.3, -1.3, 5.7, 6.5 ] )
10 errors = np.array( [  8.2,  8.4, 8.2, 8.3 ] )
11
12 corr = np.array([
13     [  1   , -0.1 , -0.05,   0.1 ],
14     [ -0.1 ,  1   ,  0.1 ,  -0.05],
15     [ -0.05,  0.1 ,  1   ,   0.1 ],
16     [  0.1 , -0.05,  0.1 ,   1   ] ] )
17
18 cova = np.ones_like( corr )
19 for i in range(4):
20     for j in range(4):
21         cova[i,j] = errors[i] * corr[i,j] * errors[j]
22
23 model  = mvn( mean=values, cov=cova )
24 ll_min = model.logpdf( values )
25
26 def draw_2d(mi):
27     fig, ax = plt.subplots()
28
29     ax.axvline( 0, c='k', lw=1)
30     ax.axhline( 0, c='k', lw=1)
```
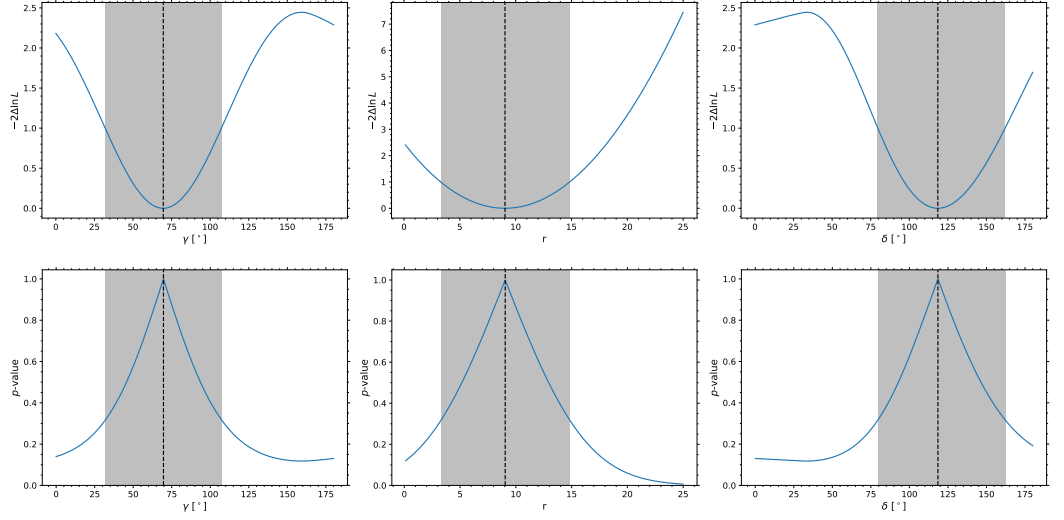
```python
31
32      mvnp = mvn( mean=values[:2], cov=cova[:2,:2] )
33      mvnm = mvn( mean=values[2:], cov=cova[2:,2:] )
34
35      x = np.linspace(-30,30,100)
36      y = np.linspace(-30,30,100)
37
38      X, Y = np.meshgrid(x,y)
39
40      pos = np.dstack([X,Y])
41
42      n2llp = -2 * ( mvnp.logpdf( pos ) - mvnp.logpdf( values[:2] ) )
43      n2llm = -2 * ( mvnm.logpdf( pos ) - mvnm.logpdf( values[2:] ) )
44
45      sigmas = np.array( [1, 2] )
46      levels = chi2.ppf( chi2.cdf( sigmas**2, 1), 2 )
47      ax.contour( X, Y, n2llp, colors='b', linestyles=['-','--'],
        levels=levels )
48      ax.contour( X, Y, n2llm, colors='r', linestyles=['-','--'],
        levels=levels )
49
50      ax.plot( *values[:2], 'bx', label='$(x^+, y^+)$' )
51      ax.plot( *values[2:], 'rx', label='$(x^-, y^-)$' )
52
53      ax.text(0.02, 0.02, 'Contours hold 68.3% and 95.4%', transform=
        ax.transAxes, fontsize=11, color='0.5' )
54
55      ax.legend(loc='lower right')
56
57      ax.set_xlabel('$x$')
58      ax.set_ylabel('$y$')
59
60      fig.savefig('figs/polar_plot.pdf')
61
62  def draw_1d_profile( mi, par, title, ax=None, convert=False ):
63
64      ax = ax or plt.gca()
65
66      x, y, valid = mi.mnprofile(par, subtract_min=True, bound=mi.
        limits[par], size=200)
67      if convert:
68          y = chi2.sf(y,1)
69      ax.plot( x[valid], y[valid] )
70
71      ax.axvspan( mi.values[par]+mi.merrors[par].lower, mi.values[par
        ]+mi.merrors[par].upper, fc='0.5', alpha=0.5 )
72      ax.axvline( mi.values[par], c='k', ls='--' )
73
74      ax.set_xlabel(title)
75      ax.set_ylabel('$-2\Delta\ln L$')
76      if convert:
77          ax.set_ylabel('$p$-value')
78
79      if convert:
80          ax.set_ylim(bottom=0)
81
```

```
82  def n2ll( gamma, r, delta ):
83      g = np.radians(gamma)
84      d = np.radians(delta)
85      xp = r*np.cos(d+g)
86      yp = r*np.sin(d+g)
87      xm = r*np.cos(d-g)
88      ym = r*np.sin(d-g)
89
90      ll_here = model.logpdf( [xp, yp, xm, ym] )
91
92      n2ll = -2*(ll_here - ll_min)
93
94      return n2ll
95
96  mi = Minuit( n2ll, gamma=100, r=1, delta=150 )
97  mi.limits['gamma'] = (0,180)
98  mi.limits['delta'] = (0,180)
99  mi.limits['r'] = (0,25)
100
101 mi.migrad()
102 mi.hesse()
103 mi.minos()
104 print(mi)
105
106 draw_2d(mi)
107
108 fig, axes = plt.subplots(1,3, figsize=(19.2, 4.8))
109 for par, title, ax in zip(['gamma','r','delta'], ['$\gamma$ [$^\
        circ$]', 'r', '$\delta$ [$^\circ$]'], axes.flatten()):
110     draw_1d_profile( mi, par, title, ax=ax )
111 fig.savefig('figs/polar_prof_nll.pdf')
112
113 fig, axes = plt.subplots(1,3, figsize=(19.2, 4.8))
114 for par, title, ax in zip(['gamma','r','delta'], ['$\gamma$ [$^\
        circ$]', 'r', '$\delta$ [$^\circ$]'], axes.flatten()):
115     draw_1d_profile( mi, par, title, ax=ax, convert=True )
116 fig.savefig('figs/polar_prof_pval.pdf')
117
118 plt.show()
```

My code will make the following plots:

(b) The profile likelihood looks fine for $r$ (depending on where you run the limits to). The problem is you hit a physical boundary at 0 and thus cannot quote an interval out to 95%.

(c) To compute the coverage you would through toys (pseduo-experiments) from the best-fit point. Refit those and compute what fraction of the time the $X\%$ interval contains the true value. This would tell you if you have a coverage issue. If you really did this for the 90% interval on $\gamma$ you would see that it would undercover.

(d) To set the appropriate upper limit on $r$ you should use the Feldman-Cousins method (as $r$ is near the physical boundary).

6. Most powerful is the log-likelihood ratio and each is Poisson distributed so

$$L_1/L_0 = \prod_i^N \text{Poiss}(d_i; s_i + b_i) / \prod_i^N \text{Poiss}(d_i; b_i) \tag{30}$$

$$= \prod_i^N \frac{e^{-(b_i+s_i)}(b_i + s_i)^{d_i}}{d_i!} / \prod_i^N \frac{e^{-b_i} b_i^{d_i}}{d_i!} \tag{31}$$

$$= \prod_i^N \frac{e^{-s_i}(s_i + b_i)^{d_i}}{b_i^{d_i}} \tag{32}$$

$$\ln(L_1/L_0) = \sum_i^N d_i \ln(s_i + b_i) - s_i - d_i \ln(b_i) \tag{33}$$