

# Data Science Analysis - Assignment 3

Varenya Upadhyaya

January 30, 2023

1. The following code generates the bootstrap samples from a standard normal distributions and plots the histogram along with the gaussian fit:

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 from scipy import stats
4 from astroML.resample import bootstrap
5 from astroML.stats import median_sigmaG
6
7 N = 1000
8 N_boot = 10000
9 np.random.seed(10)
10 dist = stats.norm(0, 1).rvs(N)
11 sample_boot, sigmaG = bootstrap(dist, N_boot, median_sigmaG, kwargs = dict(axis=1))
12 x = np.linspace(-1,1,1000)
13 sigma = np.sqrt(np.pi/(2*N))
14 pdf = stats.norm(np.mean(sample_boot), sigma).pdf(x)
15
16 #plotting the histogram/distribution
17 plt.figure(figsize=(7,5))
18 plt.hist(sample_boot, bins=20, density=True, label='Bootstrap samples histogram',color
19          = '#48b5c4')
20 plt.plot(x, pdf, label='Gaussian fit',color='#115f9a' )
21 #minor
22 plt.xlim(-0.3,0.3)
23 plt.title('Bootstrap Samples')
24 plt.legend()
25 plt.savefig('1.png')
```

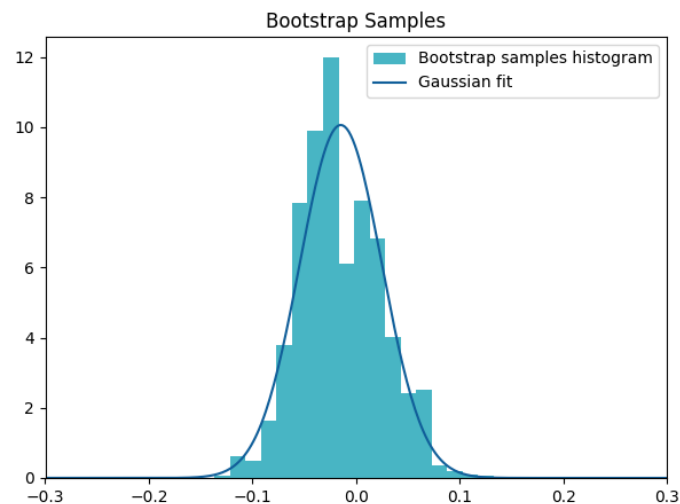


Figure 1: Data with the best fit

2. The following code plots Fig. (2) containing the data and the best fit line computed using `scipy.optimize.curve_fit`

```

1 from scipy.optimize import curve_fit
2 import numpy as np
3 import matplotlib.pyplot as plt
4
5 def line(x, m, c):
6     return c+m*x
7
8 ID,x,y,sigma_y,sigma_x,rho_xy = np.loadtxt('data.txt', unpack=True)
9
10 #use curve_fit to perform chi^2 minimization
11 param, param_cov = curve_fit(line, x, y, sigma= sigma_y,absolute_sigma=True)
12 perr = np.sqrt(np.diag(param_cov))
13
14 #printing outputs
15 m = str(round(param[0],2))
16 c = str(round(param[1],1))
17 err_m = str(round(perr[0],2))
18 err_c = str(round(perr[1],1))
19 print("m = {}, err_m = {}".format(m,err_m))
20 print('c = {}, err_c = {} '.format(c,err_c))
21
22 #plotting the data/line
23 x_axis = np.linspace(0,300)
24 plt.figure(figsize=(9,7))
25 plt.text(125,150,"$y=({}\pm{})x+({}\pm{})$".format(m,err_m,c,err_c),fontsize = 16,
26         fontweight='bold',color='black')
27 plt.errorbar(x,y,sigma_y, fmt='h', ms=6, color='#115f9a', mfc='#115f9a', mew=1, ecolor
28             = '#115f9a', alpha=0.75, capsize=2.0, zorder=0, label='Data');
29 plt.plot(x_axis, param[0]*x_axis+param[1], '-', color='#a6d75b',label = 'Best-fit line')
30
31 #minor
32 plt.title('Data vs. Best Fit Line',fontweight='bold',fontsize=16)
33 plt.xlim(0,300)
34 plt.ylim(0,700)
35 plt.legend(fontsize=16)
36 plt.savefig('2.png')

```

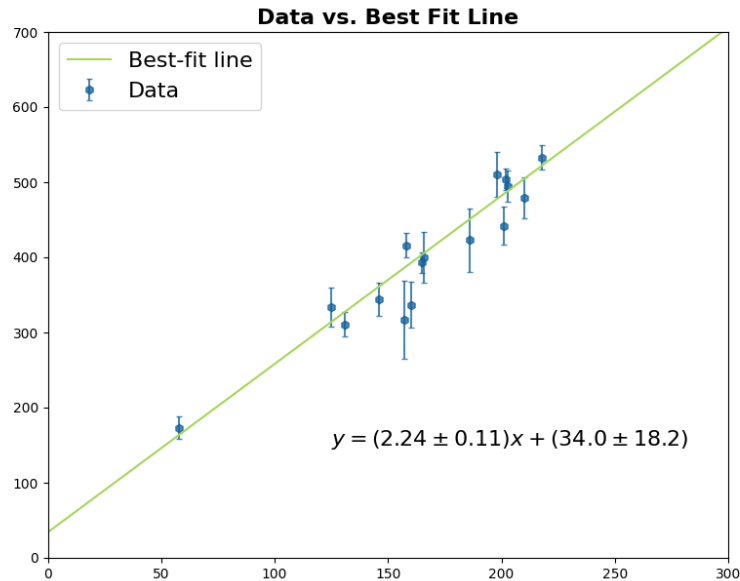


Figure 2: Data with the best fit

Code output:

```
m = 2.24, err_m = 0.11
c = 34.0, err_c = 18.2
```

3. The p-values of the four  $\chi^2$  are calculated using the following code:

```
1 import numpy as np
2 from scipy import stats
3
4 N = 50 #from the source code
5 chi2 = np.array([0.96,0.24,3.84,2.85])*(N-1)
6 p = 1 - stats.chi2.cdf(chi2,N-1)
7 print('p-values =',p)
```

Code output:

```
p-values = [5.52926434e-01 9.99999992e-01 0.00000000e+00 1.21072929e-10]
```

All the codes and figures used in this assignment can be found [in this repository](#)