

# STATS 202: Data Mining and Analysis

## Homework #2

Adam Kainikara

July 17, 2023

Problem 1 (5 points)

Chapter 4, Exercise 1 (p. 189).

$$4.2 \quad p(X) = \frac{e^{\beta_0 + \beta_1 X}}{1 + e^{\beta_0 + \beta_1 X}}$$

$$4.3 \quad \frac{p(X)}{1 - p(X)} = e^{\beta_0 + \beta_1 X}$$

Starting with 4.3

$$p(X) = (1 - p(X))e^{\beta_0 + \beta_1 X}$$

$$p(X) = e^{\beta_0 + \beta_1 X} - e^{\beta_0 + \beta_1 X} p(X)$$

$$p(X) + e^{\beta_0 + \beta_1 X} p(X) = e^{\beta_0 + \beta_1 X}$$

$$p(X)(1 + e^{\beta_0 + \beta_1 X}) = e^{\beta_0 + \beta_1 X}$$

$$p(X) = \frac{e^{\beta_0 + \beta_1 X}}{1 + e^{\beta_0 + \beta_1 X}}$$

This is 4.2

Therefore 4.2 is equivalent to 4.3

Problem 2 (5 points)

Chapter 4, Exercise 4 (p. 189).

a) Because we are using 10% of observations in the range of X. This is the same as  $[100x + 5)\%$ . To find the average, we will have to integrate it.

$$\int_{0.05}^{0.95} 10x dx + \int_0^{0.05} 10x + 5 dx + \int_{0.95}^1 105 - 100x dx = 9.75$$

b) X1 and X2 are independent. Under the same assumptions the fraction of available observations would be  $9.75\%^2 = \sim 0.95\%$

c) This is similar to (b) just with 100 times.  $9.75\%^{100} = \sim 0\%$

d) As the p increases, the fraction of amiable observations tends to 0.

e)  $p=100 \quad l = 0.1^{0.01}$

Problem 3 (5 points)

Chapter 4, Exercise 6 (p. 191).

a)  $\hat{\beta}_0 = -6, \hat{\beta}_1 = 0.05, \hat{\beta}_2 = 1$  are the coefficients for this question. These are coefficients for the constant, number of hours studied and undergrad GPA. Plugging in the values for this particular student gives  $-6 + (0.05 \times 40) + (1 \times 3.5) = -0.5$ . This value can be plugged into the logistic formula to calculate the probability.

$$\hat{p}(x) = \frac{e^x}{1 + e^x}$$

$$\hat{p}(x) = \frac{e^{-0.5}}{1 + e^{-0.5}} = 0.378$$

A student who studied for 40 hours and has an undergrad GPA of 3.5 has about a 37.8% chance of getting an A in the class.

b) How long does the student need to study for a 50% chance of an A.

Want  $\hat{p}(x) = 0.5$

$$0.5 = \frac{e^x}{1+e^x}$$

$$\frac{1}{2e^x} = \frac{1}{1+e^x}$$

$$2e^x = 1 + e^x$$

$$x = 0$$

$$0 = -6 + 0.05X_1 + 1X$$

Because its the same student and only the hours studied is changing  $X_2$ (the GPA) is still the same

$$0 = -6 + 0.05X_1 + 1(3.5)$$

$$X_1 = 50$$

If this student wants a 50% of an A then they should study for 50 hours.

Problem 4 (5 points)

Chapter 4, Exercise 8 (p. 191).

Using KNN with  $K = 1$ , the training error rate can be calculated to be 0%. This means that the test error rate has to be 36% in order for the average to be 18%. Because the logistic regression fit had a test error rate of 20%, we should use the logistic fit because it

has a lower test error rate.

Problem 5 (5 points)

Chapter 4, Exercise 13 parts a-h (p. 193)

a) See bar graph

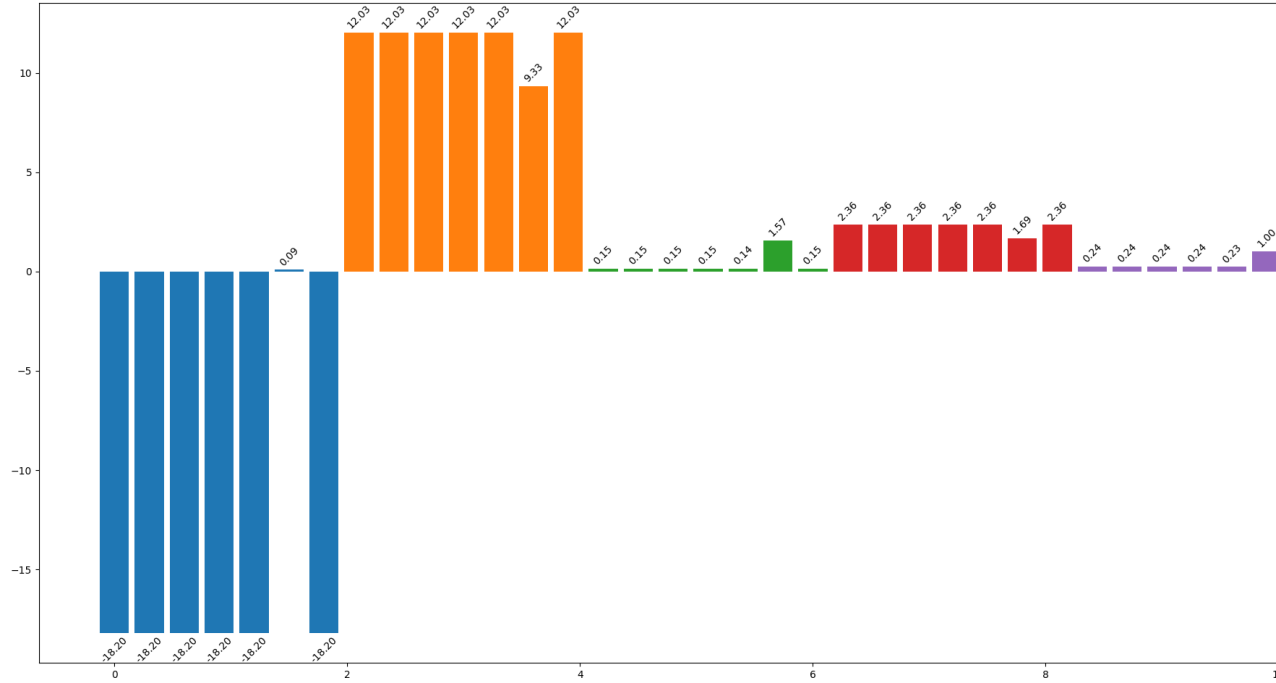


Figure 1: Min, Max, Mean Stand Dev and Median in the Colors Blue, Orange, Green, Red and Purple Respectively

The bar graph shows the min, max, mean stand dev and median in the colors blue, orange, green, red and purple respectively. These values are for Lag 1, 2, 3, 4, 5, Volume and Today. The bar graph shows similar results for each parameter. Only volume seems to change. The numerical summary for this data is below.

min: [-18.195 -18.195 -18.195 -18.195 -18.195 0.087465 -18.195 ]  
max: [12.026 12.026 12.026 12.026 12.026 9.328214 12.026 ]  
mean: [0.15058494 0.15107897 0.14720478 0.14581818 0.13989256 1.57461763  
0.14989899]  
stand dev: [2.35593008 2.35617168 2.35941794 2.35919491 2.36020027 1.68586184  
2.35584499]  
median: [0.241 0.241 0.241 0.238 0.234 1.00268 0.241 ]

b) The coefficients for the logistic regression are below.  
[[ 5.37367327e-02 5.34098644e-03 -1.50782005e-02 5.70986281e-02 9.18621394e-02 7.21885742e+00]]

The logistic regression was performed with direction as the response and the five lag variables plus volume as predictors.

Used stats models API to obtain p values. Lag 2 had a p value of 0.0296 which is lower than 0.05 so the null hypothesis can be rejected. Lag 2 does not influence direction.

c) I first did logistic regression on the whole data set. Then I computed the confusion matrix. The confusion matrix resulted the following:

```
[[482 2]
 [0 605]]
```

The model was very successful in predicting whether the data would go up or down. The model correctly predicted down every single time. When predicting up, there were only two instances the model predicted wrong of the 607 times. Of all 1089 data points, just two where predicted wrong.

d) Training data is 1990 to 2008 which is until row 986.

I fitted a logistic regression using the training data. I then computed the confusion matrix for the left out data (test data) which was the data from 2009 and 2010. The confusion matrix resulted the following:

Confusion Matrix of the test data with Lag2 as the only predictor

```
[[ 9 34]
 [ 5 56]]
```

Of the 104 data points from 2009 and 2010, 39 (37.5%) where incorrect. 37.5% is the test error rate. In this time period there were 61 ups and 43 downs. Of the 61 ups, 5 were incorrect (8.2%). Of the 43 downs, 34 where incorrect (79.1%).

e, f, g, h, i) Doing d again but with LDA, QDA, KNN.

Confusion Matrix of the training data with Lag2 as the only predictor but instead used LDA

```
[[ 22 419]
 [ 20 524]]
```

Confusion Matrix of the test data with Lag2 as the only predictor but instead used LDA

```
[[ 9 34]
 [ 5 56]]
```

Confusion Matrix of the training data with Lag2 as the only predictor but instead used QDA

```
[[ 0 441]
 [ 0 544]]
```

Confusion Matrix of the test data with Lag2 as the only predictor but instead used QDA

```
[[ 0 43]
 [ 0 61]]
```

Confusion Matrix of the training data with Lag2 as the only predictor but instead used NAIVE BAYES

```
[[ 0 441]
 [ 0 544]]
```

Confusion Matrix of the test data with Lag2 as the only predictor but instead used NAIVE BAYES

```
[[ 0 43]
```

[ 0 61]]

Confusion Matrix of the training data with Lag2 as the only predictor but instead used KNN

[[298 143]

[106 438]]

Confusion Matrix of the test data with Lag2 as the only predictor but instead used KNN

[[16 27]

[19 42]]

Error Rates (test data)	Lag 2 Logi	LDA	QDA	N
Test Error Rate	37.5%	37.5%	41.3%	
Incorrectly Predicted Down	79.1% (34/43 wrong)	79.1% (34/43 wrong)	100% (43/43 wrong)	100%
Incorrectly Predicted Up	8.2 % (5/61 wrong)	8.2 % (5/61 wrong)	0% (0/61 wrong)	0%

Lag 2 Logi only had the lowest test error rate. KNN was the best at predicting down. QDA and NAIVE BAYES was best at picking up.

j) I did some transformations. These where using Lag 2 and 3 and changing KNN to  $K = 7$ .

Confusion Matrix of the training data with Lag2 and Lag3 as the only predictors

[[ 23 418]

[ 23 521]]

Confusion Matrix of the test data with Lag2 and Lag 3 as the only predictor

[[ 8 35]

[ 4 57]]

Confusion Matrix of the training data with Lag2 and Lag3 as the only predictors but instead used LDA

[[ 22 419]

[ 22 522]]

Confusion Matrix of the test data with Lag2 and Lag3 as the only predictors but instead used LDA

[[ 8 35]

[ 4 57]]

Confusion Matrix of the training data with Lag2 and Lag 3 as the only predictors but instead used QDA

[[ 12 429]

[ 13 531]]

Confusion Matrix of the test data with Lag2 and Lag3 as the only predictors but instead used QDA

[[ 4 39]

[2 59]]

Confusion Matrix of the training data with Lag2 as the only predictor but instead used NAIVE BAYES

[[ 0 441]

[ 0 544]]

Confusion Matrix of the test data with Lag2 as the only predictor but instead used NAIVE BAYES

[[ 0 43]  
[ 0 61]]

Confusion Matrix of the training data with Lag2 and Lag 3 as the only predictor but instead used KNN

[[254 187]  
[120 424]]

Confusion Matrix of the test data with Lag2 and Lag 3 as the only predictor but instead used KNN [

[11 32]  
[17 44]]

Error Rates (test data)	Lag 2 and 3 Logi	LDA	QDA	N
Test Error Rate	37.5	37.5	39.4	
Incorrectly Predicted Down	81.4% (35/43 wrong)	81.4% (35/43 wrong)	90.7% (39/43 wrong)	100
Incorrectly Predicted Up	6.3 % (4/61 wrong)	6.3 % (4/61 wrong)	3.3% (2/61 wrong)	0

Problem 6 (5 points)

Chapter 5, Exercise 2 (p. 219).

a) Let  $n$  be the number of observations. The probability that the  $j$ th observation is in the bootstrap is  $\frac{1}{n}$  so the probability the  $j$ th observation is not in the bootstrap is  $1 - \frac{1}{n}$ .

b) Each bootstrap is independent. So the probability that the  $j$ th observation is not in the second is the same:  $1 - \frac{1}{n}$ .

c) Bootstrapping has sample with replacement and is independent. The probability that the  $j$ th observation is not in a bootstrap sample is  $1 - \frac{1}{n}$ . The probability that the  $j$ th observation is not in the bootstrap sample is the product of this. So it becomes  $(1 - \frac{1}{n}) \times (1 - \frac{1}{n}) \times \dots = (1 - \frac{1}{n})^n$

d)  $n = 5$ , finding if the  $j$ th observation is in the bootstrap.

$P(\text{observation in the bootstrap}) = 1 - P(\text{observation is not in the bootstrap})$

$$1 - (1 - \frac{1}{5})^5 = 0.67$$

e) Same as above but now  $n = 100$

$$1 - (1 - \frac{1}{100})^{100} = 0.63$$

f) Same as above but now  $n = 10000$

$$1 - (1 - \frac{1}{10000})^{10000} = 0.63$$

g) Note: For the graph I did  $n = 20,000$  because when I did  $n = 100,000$  my computer crashed and above 30,000 it was taking a very long time for the graph to load.

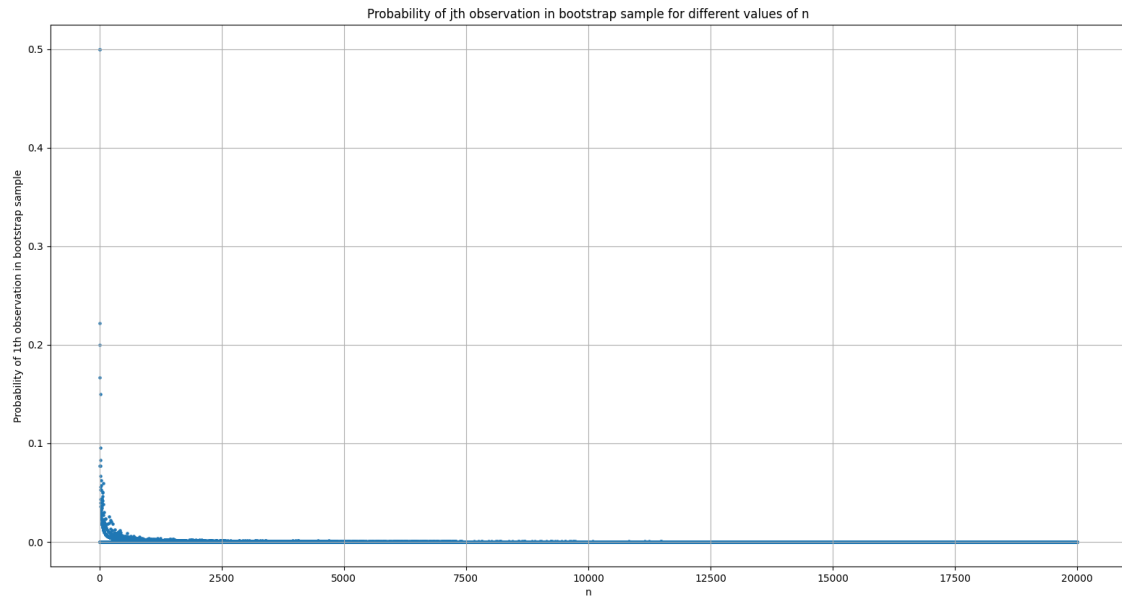


Figure 2: Probability of Jth Observation

The graph asymptotes around 0.63

h) As  $n \rightarrow \infty, p = 0.632$

Problem 7 (5 points)

Chapter 5, Exercise 5 (p. 220).

a) The coefficients for income and balance fit using logistic regression to predict default are:

[[5.64710797e-03,  
2.08091984e-05]]

b) 1) Did split into test and validation set. The data was split with 5000 observations in each set.

2) Fitted a multi logistic model. Coefficients are:

[[ 0.00041108  
-0.00012325]]

3) Classified using posterior probability

4) Computed the validation set error. The confusion matrix was:

[[4996 1]  
[ 3 0]]

The error for the validation set is 0.08%

c) 1) Did split into test and validation set. The data was split with 7500 observations in the training set and 2500 observations in the validation set.

2) Fitted a multi logistic model. Coefficients are:

```
[[5.79458398e-03
 2.30839912e-05]]
```

3) Classified using posterior probability

4) Computed the validation set error. The confusion matrix was:

```
[[7275 107]
 [ 112  6]]
```

The error for the validation set is 2.2%

5) Did split into test and validation set. The data was split with 2500 observations in the training set and 7500 observations in the validation set.

6) Fitted a multi logistic model. Coefficients are:

```
[[ 0.00043135
 -0.00012144]]
```

7) Classified using posterior probability

8) Computed the validation set error. The confusion matrix was:

```
[[2499  0]
 [  1  0]]
```

The error for the validation set is 0.04%

9) Did split into test and validation set. The data was split with 6000 observations in the training set and 4000 observations in the validation set.

10) Fitted a multi logistic model. Coefficients are:

```
[[5.95395704e-03
 3.26622332e-05]]
```

11) Classified using posterior probability

12) Computed the validation set error. The confusion matrix was:

```
[[3872  65]
 [  61  2]]
```

The error for the validation set is 3.2%

d) I used a list comprehension to switch the “yes” and “no” to 1 and 0. Then added it to the array. Then computed the logistic regression and made a confusion matrix of the validation set. The results are below:

```
[[4996  1]
 [  3  0]]
```

The error for the validation set is 0.08%

Adding a dummy variable, student, did not result in an increase or decrease in the error rate.

Problem 8 (5 points)

Chapter 5, Exercise 6 (p. 221).

a) The standard error for the coefficients balance, income and intercept are as follows.

<i>Name</i>	<i>Standard Error for the coefficients</i>
<i>balance</i>	0
<i>income</i>	$4.99e - 06$
<i>intercept</i>	0.435



b) See code  
c,d) The standard error was 3.804649256670676.  
Problem 9 (5 points)  
Chapter 5, Exercise 8 (p. 222).  
a) In this problem  $n = 100$  and  $p = 2$ . The model in this problem is  $Y = X - 2X^2 + \text{constant}$   
b) See scatter plot. The scatter plot has a upside down parabolic shape. Most of the points occur at the top of the curve.  
c, d ,e)

I first fitted these 4 models using linear algebra and coding it.

$$\begin{aligned} i & Y = \beta_0 + \beta_1 X + \epsilon \\ ii & Y = \beta_0 + \beta_1 X + \beta_2 X^2 + \epsilon \\ iii & Y = \beta_0 + \beta_1 X + \beta_2 X^2 + \beta_3 X^3 + \epsilon \\ iv & Y = \beta_0 + \beta_1 X + \beta_2 X^2 + \beta_3 X^3 + \beta_4 X^4 + \epsilon \end{aligned}$$

The seed was set to 1:

`rng = np.random.default_rng(1)`

The resulting coefficients for the models are as follows:

	$\beta_0$	$\beta_1$	$\beta_2$	$\beta_3$	$\beta_4$
i	-1.46496301	1.94936857			
ii	-0.07275529	0.96627276	-2.00470902		
iii	-0.05719669	1.1145842	-2.04709357	-0.06430033	
iv	0.10084766	0.90499786	-2.50592308	0.03376837	0.10421699

I then used LOOCV to help fit the models. The coefficients for each of the models that produced the lowest MSE are as follows:

	$\beta_0$	$\beta_1$	$\beta_2$	$\beta_3$	$\beta_4$	Error
i	-1.46518099	1.94913361				5.9226536056081205
ii	-0.0726598	0.96615516	-2.00472843			0.9834354527808048
iii	-0.05737607	1.11461298	-2.04701671	-0.06429704		0.9718403920130584
iv	0.10079418	0.90506572	-2.50587159	0.03374964	0.10420779	0.9201948751940691

Changed the seed to 100.

`rng = np.random.default_rng(100)`

The resulting coefficients for the models are as follows:

	$\beta_0$	$\beta_1$	$\beta_2$	$\beta_3$	$\beta_4$
i	-1.69029473	0.50898235			
ii	0.13096828	0.77172222	-1.94030595		
iii	0.15961224	0.46259444	0.46259444	0.13689711	
iv	0.07938098	0.41879764	-1.75905017	0.16392833	-0.06110937

I then used LOOCV to help fit the models. The coefficients for each of the models that produced the lowest MSE are as follows:

	$\beta_0$	$\beta_1$	$\beta_2$	$\beta_3$	$\beta_4$	Error
i	-1.69066004	0.50920579				5.790791169159457
ii	0.13115837	0.7717761	-1.94038458			1.0288896723158567
iii	0.15971326	0.46265244	-1.98187348	0.13688155		0.994999233705984
iv	0.07952053	0.41882023	-1.75921296	0.16391697	-0.0610765	0.9864312854873168

In both seeds model  $iv : Y = \beta_0 + \beta_1 X + \beta_2 X^2 + \beta_3 X^3 + \beta_4 X^4 + \epsilon$  had the lowest LOOCV. I had mostly expected it as I thought the error would decrease

but *iv* would over fit because we knew what the true function was. I had expected that *ii* would have the lower error because it has the same polynomial type as the true function but *iv* does have a lower error it just over fits the data.

f) The coefficients for  $\beta_0, \beta_1, \beta_2$  all have p values less than 0.05. This means that these coefficients significantly help predict y. This is what should happen because the original function is a quadratic.

Problem 10 (5 points)

Chapter 5, Exercise 9 (p. 223).

a) The population mean:  $\hat{\mu} = 22.53$

b) The standard error of the population mean 0.408. On average the mean of the population will be off from the population mean by 0.408.

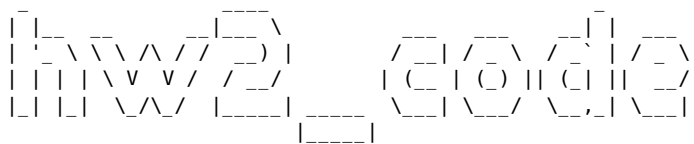
c) Using bootstrap I got a standard error of 0.4018559

d) Con Int: [21.752012851895014, 23.31359979632634

e)  $\mu_{med} = 21.2$

f) Standard Error of Median 0.36652890745478695

g,h) Tenth Percentile ( $\hat{\mu}_{0.1}$ ) of medv: 12.75 with an error of 0.49. This shows how much the 10 percentile median would be off by on average.



<a href="#">hw2v2/</a> .....	1
└── ch4p13.py .....	1
└── ch5p2.py .....	5
└── ch5p5.py .....	6
└── ch5p6.py .....	8
└── ch5p8.py .....	10
└── ch5p9.py .....	12



```

1 from numpy import *
2 import numpy as np
3 import matplotlib as mpl
4 import matplotlib.pyplot as plt
5 from scipy import stats
6
7 import sys
8 import matplotlib
9
10 import matplotlib.pyplot as plt
11 import numpy
12 from sklearn.datasets import load_iris
13 from sklearn.linear_model import LogisticRegression
14 from sklearn.metrics import confusion_matrix
15 #from sklearn.metrics import accuracy_score #works
16 from matplotlib.ticker import FormatStrFormatter
17 import statsmodels.api as sm
18
19 from sklearn.discriminant_analysis import LinearDiscriminantAnalysis as lda
20 from sklearn.discriminant_analysis import QuadraticDiscriminantAnalysis as qda
21 from sklearn.neighbors import KNeighborsClassifier
22 from sklearn.naive_bayes import GaussianNB
23
24 #Adam Kainikara
25 #This code is for
26 #CHAPTER 4 QUESTION 13
27 #THIS IS FOR PROBLEM 5
28 #OF HOMEWORK 2 FOR STANFORD SUMMER SESSION STATS 202
29
30
31 def data_loader(fname):
32     data_a = loadtxt(fname, skiprows=1, usecols=(1,2,3,4,5,6,7), delimiter=',')
33     direction_v = loadtxt(fname, skiprows=1, usecols=(8), delimiter=',', dtype=str)
34
35     return data_a, direction_v
36
37 def summaries0(x_a):
38     print("Numerical summaries")
39     for i in range(x_a.shape[1]):
40         column = x_a[:, i]
41         print(f'Column {i + 1}:')
42         print(f' - Min: {min(column)}')
43         print(f' - Max: {max(column)}')
44         print(f' - Mean: {mean(column)}')
45         print(f' - Median: {median(column)}')
46         print(f' - Stand Dev: {std(column)}')
47     print(x_a.shape[1])
48
49     summary = empty((8,5))
50     storage_l = []
51
52     for i in range(x_a.shape[1]):
53         column = x_a[:, i]
54         summary = [min(column), max(column), mean(column), median(column), std(column)]
55         print(summary)
56         storage_l.append(summary)
57     return np.array(storage_l)
58
59 def calc_summary(x_a) -> dict[str, ndarray]:
60     data_d = {
61         "min": x_a.min(axis = 0),
62         "max": x_a.max(axis = 0),
63         "mean": x_a.mean(axis = 0),
64         "stand dev": x_a.std(axis = 0),
65         "median": median(x_a, axis = 0)
66     }
67     return data_d
68
69 def print_summary(data_d: dict[str, ndarray]):
70     for k,v in data_d.items():
71         print(f'{k}: ', v)
72
73 def find_p_values(x_a, direction_v):
74     #y_v = direction_v
75     #log_reg = sm.Logit(y_v, ex x_a).fit()
76
77     # Extract the p-values
78     #print(log_reg.summary())
79     pass
80
81 def plot_summary(data_d):
82     data_d = {'min': data_d['min'], 'max': data_d['max'], 'mean': data_d['mean'], 'stan dev': data_d['stand dev'], 'median': data_d['median']}
83
84     nvar = len(data_d["min"])
85     print(data_d)
86     w = 0.25
87     stride = w + 0.05
88     initial_change = 0
89     labels = ("Min", "Max", "Mean", "Stan Dev", "Median")
90     x_v = arange(nvar) * stride + 0
91     print(x_v)
92
93     fig, ax = plt.subplots()
94
95     for attribute, measurement in data_d.items():
96         #shift = w * initial_change
97         rectangle = ax.bar(x_v, measurement, width=w, label = attribute)
98         x_v += nvar * stride
99         ax.bar_label(rectangle, fmt = '%.02f', rotation = 45, padding=2)
100         #initial change += 0.25
101         #ax.set_xticks(x_v, labels)

```

```

103 plt.show()
104
105
106 def summaries(x_a):
107     print("Numerical summaries")
108     #min_x_a = x_a.min(axis = 0)
109     print(f' Min:{x_a.min(axis = 0)}, Max{x_a.max(axis = 0)}, Mean{x_a.mean(axis = 0)}, Stand Dev {x_a.std(axis = 0)}, Median{median(x_a, axis = 0)}')
110
111 def logi_reg(x_a, direction_v):
112     predictors_a = x_a[:,1:7]
113     #print(predictors_a)
114     #print(x_a.shape)
115     response_v = direction_v
116
117     #response_v = x_a[:,8]
118     #print(response_v)
119     logreg = LogisticRegression()
120     logreg.fit(predictors_a, response_v)
121     coefficients = logreg.coef_
122
123     print(coefficients)
124
125 def make_prediction(x_a, direction_v):
126     clf = LogisticRegression()
127     clf.fit(x_a, direction_v)
128     ypredict_v = clf.predict(x_a)
129     return ypredict_v
130
131
132 def compute_confusion_mat(ypredict_v, direction_v):
133     truey_v = direction_v
134     confu_mat = confusion_matrix(truey_v, ypredict_v)
135     print(confu_mat)
136     return confu_mat
137
138 def lda_prediction(x_v, y_v):
139     #This only served to help me write the code, I did it in main otherwise
140     clf = lda()
141     clf.fit(x_v, y_v)
142     ypredict_ldatrain_v = clf.predict(x_v)
143     ypredict_ldatest_v = clf.predict(x_v)
144     #Rembr to do the confusion matrix after
145     return ypredict_ldatrain_v, ypredict_ldatest_v
146     #What I ended up doing
147     clf = lda()
148     clf.fit(xtrain_v, ytrain_v)
149
150     ypredict_ldatrain_v = clf.predict(xtrain_v)
151     print("Confusion Matrix of the training data with Lag2 as the only predictor but instead used LDA")
152     compute_confusion_mat(ypredict_ldatrain_v, ytrain_v)
153
154     ypredict_ldatest_v = clf.predict(xtest_v)
155     print("Confusion Matrix of the test data with Lag2 as the only predictor but instead used LDA")
156     compute_confusion_mat(ypredict_ldatest_v, ytest_v)
157
158 def qda_prediction(x_v, y_v):
159     #This only served to help me write the code, I did it in main otherwise
160     clf = qda()
161     clf.fit(x_v, y_v)
162     ypredict_qdatrain_v = clf.predict(x_v)
163     ypredict_qdatest_v = clf.predict(x_v)
164     #Rembr to do the confusion matrix after
165     return ypredict_qdatrain_v, ypredict_qdatest_v
166
167     #What I ended up doing
168     clf = qda()
169     clf.fit(xtrain_v, ytrain_v)
170     ypredict_qdatrain_v = clf.predict(xtrain_v)
171     print("Confusion Matrix of the training data with Lag2 as the only predictor but instead used QDA")
172     compute_confusion_mat(ypredict_qdatrain_v, ytrain_v)
173
174     ypredict_qdatest_v = clf.predict(xtest_v)
175     print("Confusion Matrix of the test data with Lag2 as the only predictor but instead used QDA")
176     compute_confusion_mat(ypredict_qdatest_v, ytest_v)
177
178 def naiv_prediction(x_v, y_v):
179     #This onl served to elp me write the code, I did it in main otherwise
180     clf = GaussianNB()
181     clf.fit(x_v, y_v)
182     ypredict_nbtrain_v = clf.predict(x_v)
183     ypredict_nbtrest_v = clf.predict(x_v)
184     return ypredict_nbtrain_v, ypredict_nbtrest_v
185
186 def knn_prediction(x_v, y_v):
187     neigh = KNeighborsClassifier(n_neighbors=3)
188     neigh.fit(x_v, y_v)
189     ypredict_knntrain_v = neigh.predict(x_v)
190     ypredict_knntest_v = neigh.predict(x_v)
191     return ypredict_knntrain_v, ypredict_knntest_v
192
193 def main():
194     x_a, y_v = data_loader("Weekly.csv")
195
196     xtrain_v = x_a[:,1:2][:985]
197     xtest_v = x_a[:,1:2][985:]
198
199     ytrain_v = y_v[:985]
200     ytest_v = y_v[985:]
201
202     '''
203     Above is the test and training data for x and y for the remainder of this problem. We will first train then do the test data.
204     Then do the confusion matrix
205     '''
206
207     '''
208     This first part (below) is fitting the training data and then predictng the y value based on the training data. I
209     t then computes the confusion matrix based on the training data
210     '''

```

```

210
211
212 clf = LogisticRegression()
213 clf.fit(xtrain_v, ytrain_v)
214
215
216 ytrain_pred_v = clf.predict(xtrain_v)
217 print("Confusion Matrix of the training data with Lag2 as the only predictor")
218 compute_confusion_mat(ytrain_pred_v, ytrain_v)
219
220
221 '''
222 This second part (below) is getting the predicted y value and computing
223 the confusion matrix based on the fit found earlier and the test data.
224 '''
225
226 ytest_pred_v = clf.predict(xtest_v)
227 print("Confusion Matrix of the test data with Lag2 as the only predictor")
228 compute_confusion_mat(ytest_pred_v, ytest_v)
229
230
231 =====
232 NOW DOING QDA
233 '''
234 clf = lda()
235 clf.fit(xtrain_v, ytrain_v)
236 ypredict_ldatrain_v = clf.predict(xtrain_v)
237 print("Confusion Matrix of the training data with Lag2 as the only predictor but instead used LDA")
238 compute_confusion_mat(ypredict_ldatrain_v, ytrain_v)
239
240 ypredict_ldatest_v = clf.predict(xtest_v)
241 print("Confusion Matrix of the test data with Lag2 as the only predictor but instead used LDA")
242 compute_confusion_mat(ypredict_ldatest_v, ytest_v)
243
244
245 =====
246 NOW DOING QDA
247 '''
248 clf = qda()
249 clf.fit(xtrain_v, ytrain_v)
250 ypredict_qdatrain_v = clf.predict(xtrain_v)
251 print("Confusion Matrix of the training data with Lag2 as the only predictor but instead used QDA")
252 compute_confusion_mat(ypredict_qdatrain_v, ytrain_v)
253
254 ypredict_qdatest_v = clf.predict(xtest_v)
255 print("Confusion Matrix of the test data with Lag2 as the only predictor but instead used QDA")
256 compute_confusion_mat(ypredict_qdatest_v, ytest_v)
257
258
259 =====
260 NOW DOING NAIVE BAEES
261 '''
262 clf = GaussianNB()
263 clf.fit(xtrain_v, ytrain_v)
264 ypredict_nbtrain_v = clf.predict(xtrain_v)
265 print("Confusion Matrix of the training data with Lag2 as the only predictor but instead used NAIVE BAYES")
266 compute_confusion_mat(ypredict_nbtrain_v, ytrain_v)
267
268 ypredict_nbtest_v = clf.predict(xtest_v)
269 print("Confusion Matrix of the test data with Lag2 as the only predictor but instead used NAIVE BAYES")
270 compute_confusion_mat(ypredict_nbtest_v, ytest_v)
271
272
273 =====
274 NOW DOING KNN
275 '''
276 neigh = KNeighborsClassifier(n_neighbors=3)
277 neigh.fit(xtrain_v, ytrain_v)
278 ypredict_kntrain_v = neigh.predict(xtrain_v)
279 print("Confusion Matrix of the training data with Lag2 as the only predictor but instead used KNN")
280 compute_confusion_mat(ypredict_kntrain_v, ytrain_v)
281
282 ypredict_kntest_v = neigh.predict(xtest_v)
283 print("Confusion Matrix of the test data with Lag2 as the only predictor but instead used KNN")
284 compute_confusion_mat(ypredict_kntest_v, ytest_v)
285
286
287 =====
288 NOW DOING LDA WITH A TWIST
289 '''
290 xtrain2_v = x_a[:,1:3][:985]
291 xtest2_v = x_a[:,1:3][985:]
292 ytrain_v = y_v[:985]
293 ytest_v = y_v[985:]
294 clf = lda()
295 clf.fit(xtrain2_v, ytrain_v)
296 ypredict_ldatrain2_v = clf.predict(xtrain2_v)
297 print("Confusion Matrix of the training data with Lag2 and Lag3 as the only predictors but instead used LDA")
298 compute_confusion_mat(ypredict_ldatrain2_v, ytrain_v)
299
300 ypredict_ldatest2_v = clf.predict(xtest2_v)
301 print("Confusion Matrix of the test data with Lag2 and Lag3 as the only predictors but instead used LDA")
302 compute_confusion_mat(ypredict_ldatest2_v, ytest_v)
303
304
305 =====
306 NOW DOING QDA WITH A TWIST
307 '''
308 clf = qda()
309 clf.fit(xtrain2_v, ytrain_v)
310 ypredict_qdatrain2_v = clf.predict(xtrain2_v)
311 print("Confusion Matrix of the training data with Lag2 and Lag 3 as the only predictors but instead used QDA")
312 compute_confusion_mat(ypredict_qdatrain2_v, ytrain_v)
313
314 ypredict_qdatest2_v = clf.predict(xtest2_v)
315 print("Confusion Matrix of the test data with Lag2 and Lag3 as the only predictors but instead used QDA")
316 compute_confusion_mat(ypredict_qdatest2_v, ytest_v)
317
318 =====

```

```

317 NOW DOING NAIVE BAYES WITH A TWIST
318 """
319 clf = GaussianNB()
320 clf.fit(xtrain2_v, ytrain_v)
321 ypredict_nbtrain2_v = clf.predict(xtrain2_v)
322 print("Confusion Matrix of the training data with Lag2 and Lag 3 as the only predictor but instead used NAIVE BAYES")
323 compute_confusion_mat(ypredict_nbtrain2_v, ytrain_v)
324
325 ypredict_nbtest2_v = clf.predict(xtest2_v)
326 print("Confusion Matrix of the test data with Lag2 and Lag 3 as the only predictor but instead used NAIVE BAYES")
327 compute_confusion_mat(ypredict_nbtest2_v, ytest_v)
328
329
330
331 =====
332 NOW DOING KNN WITH A TWIST
333 """
334 neigh = KNeighborsClassifier(n_neighbors=7)
335 neigh.fit(xtrain2_v, ytrain_v)
336 ypredict_kntrain2_v = neigh.predict(xtrain2_v)
337 print("Confusion Matrix of the training data with Lag2 and Lag 3 as the only predictor but instead used KNN")
338 compute_confusion_mat(ypredict_kntrain2_v, ytrain_v)
339
340 ypredict_kntest2_v = neigh.predict(xtest2_v)
341 print("Confusion Matrix of the test data with Lag2 and Lag 3 as the only predictor but instead used KNN")
342 compute_confusion_mat(ypredict_kntest2_v, ytest_v)
343
344 clf = LogisticRegression()
345 clf.fit(xtrain2_v, ytrain_v)
346
347
348 ytrain_pred2_v = clf.predict(xtrain2_v)
349 print("Confusion Matrix of the training data with Lag2 and Lag3 as the only predictors")
350 compute_confusion_mat(ytrain_pred2_v, ytrain_v)
351
352
353 ytest_pred2_v = clf.predict(xtest2_v)
354 print("Confusion Matrix of the test data with Lag2 and Lag 3 as the only predictor")
355 compute_confusion_mat(ytest_pred2_v, ytest_v)
356
357
358
359
360 #print(xtest_a)
361 #print(ytrain_v)
362
363 #ytest_v = y_v[:ntest]
364 #xtrain_v = x_a[ntest:]
365
366 #confusion_pract()
367 #print(x_a)
368 #logi_reg(x_a, y_v)
369 #find_p_values(x_a, direction_v)
370 #summaries(x_a)
371 summary = summaries(x_a)
372 #calc_summary(x_a)
373 data_dict = calc_summary(x_a)
374 print_summary(data_dict)
375 plot_summary(data_dict)
376
377
378 clf = LogisticRegression()
379 clf.fit(xtrain_v, ytrain_v)
380 ytrain_pred_v = clf.predict(xtrain_v)
381 compute_confusion_mat(ytrain_pred_v, ytrain_v)
382 ytest_pred_v = clf.predict(xtest_v)
383 compute_confusion_mat(ytest_pred_v, ytest_v)
384 print('Done----')
385
386 make_prediction(x_a, y_v)
387 ypredict_v = make_prediction(x_a, y_v)
388 print(ypredict_v)
389 compute_confusion_mat(ypredict_v, y_v)
390 #print(type(summary))
391 #print(summary.shape)
392 #plot_summaries(summary)
393
394
395
396 #Some Links I used
397 #https://scikit-learn.org/stable/auto_examples/classification/plot_lda_qda.html#sphx-glr-auto-examples-classification-plot-lda-qda-py
398 #https://stackoverflow.com/questions/46775155/importerror-no-module-named-sklearn-lda
399
400 if __name__ == '__main__':
401     main()

```



## # ch5p2.py

```
1 from numpy import *
2 import matplotlib.pyplot as plt
3
4 j = 1 # jth observation index
5 num_samples = 20000 # maximum value of n
6
7 n_values = arange(1, num_samples + 1)
8 probabilities = []
9
10 for n in n_values:
11     sample = random.choice(range(n), size=n, replace=True)
12     probability = mean(sample == j)
13     probabilities.append(probability)
14
15 # Creating a scatter plot
16 plt.scatter(n_values, probabilities, s=5)
17 plt.xlabel('n')
18 plt.ylabel(f'Probability of {j}th observation in bootstrap sample')
19 plt.title('Probability of jth observation in bootstrap sample for different values of n')
20 plt.grid(True)
21 plt.show()
```

## # ch5p5.py

```
1 from numpy import *
2 import numpy as np
3 import matplotlib as mpl
4 import matplotlib.pyplot as plt
5 from scipy import stats
6
7 import sys
8 import matplotlib
9
10 import matplotlib.pyplot as plt
11 import numpy
12 from sklearn.datasets import load_iris
13 from sklearn.linear_model import LogisticRegression
14 from sklearn.metrics import confusion_matrix
15 #from sklearn.metrics import accuracy_score #works
16 from matplotlib.ticker import FormatStrFormatter
17 import statsmodels.api as sm
18
19 from sklearn.discriminant_analysis import LinearDiscriminantAnalysis as lda
20 from sklearn.discriminant_analysis import QuadraticDiscriminantAnalysis as qda
21 from sklearn.neighbors import KNeighborsClassifier
22 from sklearn.naive_bayes import GaussianNB
23 from sklearn.metrics import accuracy_score
24
25
26 #Adam Kainikara
27 #This code is for
28 #CHAPTER 5 QUESTION 5
29 #THIS IS PROBLEM 7
30 #OF HOMEWORK 2 FOR STANFORD SUMMER SESSION STATS 202
31
32 def data_loader(fname):
33     num_data_a = loadtxt(fname, skiprows=1, usecols=(2,3), delimiter=',')
34     defa_student_a = loadtxt(fname, skiprows=1, usecols=(0,1), delimiter=',', dtype=str)
35
36     return num_data_a, defa_student_a
37
38 def logi_reg(x_a, y_a):
39     predictors_a = x_a
40     response_v = y_a[:,0]
41     #print(predictors_a, response_v)
42     logreg = LogisticRegression()
43     logreg.fit(predictors_a, response_v)
44     coefficients = logreg.coef_
45     print(coefficients)
46 def new_logi_reg(x_a, y_a):
47     pass
48
49 def compute_confusion_mat(ypredict_v, direction_v):
50     truey_v = direction_v
51     confu_mat = confusion_matrix(truey_v, ypredict_v)
52     print(confu_mat)
53     return confu_mat
54
55 def main():
56     x_a, y_a = data_loader("Default.csv")
57
58
59     ydefault_v = y_a[:,0]
60     print(y_a[:,1])
61     #ysudent = [1 if x == "yes" else 0 for x in x_a[:,1]]
62     ystudent = [1 if x == "Yes" else 0 for x in y_a[:,1]]
63     #print(ystudent)
64
65
66     #print(x_a)
67     student_a = array(ystudent)
68     #print(student_a.shape)
69     xall_a = array([x_a[:,0], x_a[:,1], student_a])
70     realx_a = transpose(xall_a)
71     print(realx_a.shape)
72     print(realx_a)
73
74     xalltrain_a = realx_a[:5000]
75     xallvalid_a = realx_a[5000:]
76     print(xalltrain_a, xallvalid_a)
77     ytrain_v = ydefault_v[:5000]
78     yvalid_v = ydefault_v[5000:]
79     clf = LogisticRegression()
```

```

80     clf.fit(xalltrain_a, ytrain_v)
81     ytrain_pred_v = clf.predict(xalltrain_a)
82     yvalid_pred_v = clf.predict(xallvalid_a)
83     coefficients = clf.coef_
84     print(coefficients)
85     compute_confusion_mat(ytrain_pred_v, yvalid_pred_v)
86     raise SystemExit
87
88     xall_a = ([x_a],[ysudent])
89     print(xall_a)
90     #logi_reg(x_a,y_a)
91
92
93     xtrain_a = x_a[:5000]
94     xvalid_a = x_a[5000:]
95     ytrain_v = ydefault_v[:5000]
96     yvalid_v = ydefault_v[5000:]
97
98
99
100    clf = LogisticRegression()
101    clf.fit(xtrain_a, ytrain_v)
102    ytrain_pred_v = clf.predict(xtrain_a)
103    yvalid_pred_v = clf.predict(xvalid_a)
104    coefficients = clf.coef_
105    print(coefficients)
106    #compute_confusion_mat(ytrain_pred_v, yvalid_pred_v)
107
108
109    clf = GaussianNB()
110    clf.fit(xtrain_a, ytrain_v)
111    posterior_probs = clf.predict_proba(xvalid_a)
112    predictions = (posterior_probs > 0.5)
113    print(predictions)
114    accuracy = accuracy_score(yvalid_v, predictions)
115    print("Accuracy:", accuracy)
116    #print(xtrain_a, ytrain_v)
117    raise SystemExit
118
119    ytrain_v = y_v[:985]
120    #ytest_v = y_v[985:]
121
122
123
124
125 if __name__ == '__main__':
126     main()

```

## # ch5p6.py

```
1 from numpy import *
2 import numpy as np
3 import matplotlib as mpl
4 import matplotlib.pyplot as plt
5 from scipy import stats
6
7 import sys
8 import matplotlib
9
10 import matplotlib.pyplot as plt
11 import numpy
12 from sklearn.datasets import load_iris
13 from sklearn.linear_model import LogisticRegression
14 from sklearn.metrics import confusion_matrix
15 #from sklearn.metrics import accuracy_score #works
16 from matplotlib.ticker import FormatStrFormatter
17 import statsmodels.api as sm
18
19 from sklearn.discriminant_analysis import LinearDiscriminantAnalysis as lda
20 from sklearn.discriminant_analysis import QuadraticDiscriminantAnalysis as qda
21 from sklearn.neighbors import KNeighborsClassifier
22 from sklearn.naive_bayes import GaussianNB
23 from sklearn.metrics import accuracy_score
24
25 import statsmodels.api as sm
26 rng = np.random.default_rng()
27 from scipy.stats import norm
28 #Adam Kainikara
29 #This code is for
30 #CHAPTER 5 QUESTION 6
31 #THIS IS PROBLEM 8
32 #OF HOMEWORK 2 FOR STANFORD SUMMER SESSION STATS 202
33
34 def data_loader(fname):
35     num_data_a = loadtxt(fname, skiprows=1, usecols=(2,3), delimiter=',')
36     defa_a = loadtxt(fname, skiprows=1, usecols=(0,1), delimiter=',', dtype=str)
37     ydefault = [1 if x == "Yes" else 0 for x in defa_a[:,0]]
38     ystudent = [1 if x == "Yes" else 0 for x in defa_a[:,1]]
39
40     default_a = transpose(array((ydefault, ystudent)))
41     #print(default_a)
42     return num_data_a, default_a
43
44 def use_sm(x_a, y_a):
45
46     b = ones((10000,1))
47     xareal_a = hstack((x_a,b))
48     print(y_a.dtype)
49
50     print(xareal_a)
51     logit_model = sm.Logit(y_a, xareal_a)
52     result = logit_model.fit()
53     print(result.summary())
54     predicted = result.predict(xareal_a)
55     return predicted, xareal_a
56
57
58 def boot_fn(x_a, y_a):
59     all_dataset = hstack((x_a,y_a))
60     n = all_dataset.shape[0]
61     index = arange(n)
62     #print(index.shape)
63     index_and_const = empty((n,2))
64     index_and_const[:,0] = index
65     index_and_const[:,1] = 1
66     #print(index_and_const, index_and_const.shape)
67
68     data_and_index = hstack((all_dataset, index_and_const))
69     #print(data_and_index, data_and_index.shape)
70     y_default = y_a[:,0]
71     clf = GaussianNB()
72     clf.fit(data_and_index, y_default)
73     probs = clf.predict_proba(data_and_index)
74     print(probs)
75     #predicted, xareal_a = use_sm(x_a, y_a)
76     #boot_fn(x_a)
77     return data_and_index, y_default, probs
78
79 def main():
```

```

80     x_a, y_a = data_loader("Default.csv")
81     #use_sm(x_a, y_a)
82     dist = norm(loc=2, scale=4)
83
84     data = dist.rvs(size=10, random_state=rng)
85     std_true = dist.std()
86
87     print(std_true)
88     std_sample = np.std(data)
89     # https://docs.scipy.org/doc/scipy/reference/generated/scipy.stats.bootstrap.html
90     print(std_sample)
91     raise SystemExit
92     data_and_index, y_default, probs = boot_fn(x_a, y_a)
93     #logit_model = sm.Logit(probs, y_default)
94     #print(probs.shape, data_and_index.shape)
95     #result = logit_model.fit()
96     #print(result.summary())
97     glmmodel = sm.GLM(probs, data_and_index)
98     result = glmmodel.fit
99     print(result.summary())
100
101 if __name__ == '__main__':
102     main()

```

```

1 from numpy import *
2 import numpy as np
3 import matplotlib as mpl
4 import matplotlib.pyplot as plt
5 from scipy import stats
6 from scipy.interpolate import CubicSpline
7 from scipy.interpolate import splrep, BSpline
8 import statsmodels.api as sm
9 from sklearn.model_selection import LeaveOneOut
10
11 def data_loader():
12     #random.seed(1)
13
14     rng = np.random.default_rng(100)
15     x_v = rng.normal(size = 100)
16     y_v = (x_v) - (2 * x_v**2) + (rng.normal(size = 100))
17     print(x_v.shape, y_v.shape)
18
19     x_a = empty((100,5), dtype=float64)
20     x_a[:,0] = 1
21     x_a[:,1] = x_v
22     x_a[:,2] = (x_v)**2
23     x_a[:,3] = (x_v)**3
24     x_a[:,4] = (x_v)**4
25     #print("x_v", x_v)
26     #print("x_a", x_a)
27     return x_a, x_v, y_v
28
29 def data_scatterplot(x_v,y_v):
30     plt.scatter(x_v,y_v)
31     plt.show()
32
33 def line_lin_fit(x_a, y_v):
34     #y_v = X@B
35     b_v = linalg.pinv(x_a[:,0:2])@y_v
36     print(b_v)
37
38 def line_loocv_fit(x_a, y_v):
39     loo = LeaveOneOut()
40     loo.get_n_splits(x_a)
41     degree_v = arange(1,5)
42
43     result_l = []
44
45     for degree in degree_v:
46
47         for train_i_v, test_i_v in loo.split(x_a):
48             xtrain_a = x_a[train_i_v,:degree+1]
49             ytrain_v = y_v[train_i_v]
50             b_v = linalg.pinv(xtrain_a)@ytrain_v
51             #print(b_v)
52
53             yfit_v = x_a[:,degree+1:] @ b_v
54             mse = ((y_v - yfit_v)**2).mean()
55             result_l.append((b_v, mse))
56             #print(mse)
57     return result_l
58
59
60 def quad_lin_fit(x_a, y_v):
61     #y_v = X@B
62     b_v = linalg.pinv(x_a[:,0:3])@y_v
63     print(b_v)
64 def cubic_lin_fit(x_a, y_v):
65     #y_v = X@B
66     b_v = linalg.pinv(x_a[:,0:4])@y_v
67     print(b_v)
68 def xtofour_lin_fit(x_a, y_v):
69     #y_v = X@B
70     b_v = linalg.pinv(x_a[:,0:5])@y_v
71     print(b_v)
72 def pvalue(x_a,y_v,name=''):
73     if name:
74         print(f'\n\n----- {name} -----')
75         model = sm.OLS(y_v, x_a)
76         results = model.fit()
77         print(results.summary())
78     #Using stats models to get p value even though I did my own regression
79 def main():
80     data_loader()
81     x_a, x_v, y_v = data_loader()
82     #data_scatterplot(x_a, y_v)
83     line_lin_fit(x_a, y_v)
84     quad_lin_fit(x_a, y_v)
85     cubic_lin_fit(x_a, y_v)
86     xtofour_lin_fit(x_a, y_v)
87     print("This separates normal and LOOCV")
88     mod_mse = line_loocv_fit(x_a, y_v)
89     #smse_l = sorted(mod_mse, key = lambda x_t: x_t[1])
90     #That sorted all the models, and found the one and its coefficients that produced the lowest mean squared error value
91     #smse_l = sorted(mod_mse, key = lambda x_t: x_t[1]+x_t[0].shape[0]*100000)
92     #This one aimed at finding the linear model with the lowest mean squared error value. This was done by using the kornicer delta.
93     #If the shape of the first term (where we had all the coefficients) was not 2 (!= is not equal ) (intercept and slope) it would increase the MSE
94     #By 100,000 which means it wouldnt show up cause it is sorted by decreasing MSE
95
96     for degree in range(1,5):
97         #now changed it a bit so that it loops and prints what i need for all of them
98         smse_l = sorted(mod_mse, key = lambda x_t: x_t[1]+(x_t[0].shape[0]-degree+1)*100000)
99

```

```
100     print('Degree', degree, smse_l[0])
101     #print(mod_mse_l[:10])
102     #print(mod_mse_l[100:110])
103
104
105 if __name__ == '__main__':
106     main()
```

## # ch5p9.py

```
1 from numpy import *
2 import numpy as np
3 import matplotlib as mpl
4 import matplotlib.pyplot as plt
5 import numpy as np
6 from sklearn.utils import resample
7 from scipy import stats
8 from scipy.stats import bootstrap
9 def data_loader(fname):
10     data_a = loadtxt(fname, skiprows=1, usecols=(0,1,2,3,4,5,6,7,8,9,10,11,12,13), delimiter=',')
11
12     return data_a
13
14 def pop_mean(data_a):
15     #Want population mean of medv
16     medv_v = data_a[:,13]
17     muhat = medv_v.mean()
18     #print(muhat)
19     #print(medv_v)
20     return medv_v, muhat
21
22 def stand_error_muhat(medv_v, muhat):
23     #Hint: We can compute the standard error of the sample mean by dividing the sample standard deviation by the square root of the
24     # number of observations.
25     stdmuhat = medv_v.std()
26     n = medv_v.shape[0]
27
28     stand_err_muhat = stdmuhat/(n**0.5)
29     print(stand_err_muhat)
30     return stand_err_muhat
31
32 def newbootstrapstderror(medv_v):
33     #I am not sure if i coded a method of bootstrap correctly. I referenced the following websites
34     #https://docs.scipy.org/doc/scipy/reference/generated/scipy.stats.bootstrap.html
35     #https://www.statology.org/bootstrapping-in-python/
36     #https://medium.com/swlh/bootstrap-sampling-using-pythons-numpy-85822d868977
37
38     nbootstrap = 1000
39     bootstrapmeans = []
40     for _ in range(nbootstrap):
41         bootstrap_sample = random.choice(medv_v, size=len(medv_v), replace=True)
42         bootstrap_sample_mean = mean(bootstrap_sample)
43         bootstrapmeans.append(bootstrap_sample_mean)
44
45     standard_error = np.std(bootstrapmeans)
46
47     print("Standard Error of  $\mu$  using Bootstrap:", standard_error)
48     return standard_error
49
50 def newbootstrapstderror_median(medv_v):
51     #I am not sure if i coded a method of bootstrap correctly. I referenced the following websites
52     #https://docs.scipy.org/doc/scipy/reference/generated/scipy.stats.bootstrap.html
53     #https://www.statology.org/bootstrapping-in-python/
54     #https://medium.com/swlh/bootstrap-sampling-using-pythons-numpy-85822d868977
55
56     nbootstrap = 1000
57     bootstrapmedian = []
58     for _ in range(nbootstrap):
59         bootstrap_sample = random.choice(medv_v, size=len(medv_v), replace=True)
60         bootstrap_sample_median = median(bootstrap_sample)
61         bootstrapmedian.append(bootstrap_sample_median)
62
63     standard_error = std(bootstrapmedian)
64
65     print("Standard Error of  $\mu$  hat median using Bootstrap:", standard_error)
66     return standard_error
67
68 def newbootstrapstderror_tenpercen(medv_v):
69     #I am not sure if i coded a method of bootstrap correctly. I referenced the following websites
70     #https://docs.scipy.org/doc/scipy/reference/generated/scipy.stats.bootstrap.html
71     #https://www.statology.org/bootstrapping-in-python/
72     #https://medium.com/swlh/bootstrap-sampling-using-pythons-numpy-85822d868977
73
74     nbootstrap = 1000
75     bootstrappercen = []
76     for _ in range(nbootstrap):
77         bootstrap_sample = random.choice(medv_v, size=len(medv_v), replace=True)
78         bootstrap_sample_percen = percentile(bootstrap_sample, 10)
79         bootstrappercen.append(bootstrap_sample_percen)
80
81     standard_error = std(bootstrappercen)
82
83     print("Standard Error of  $\mu$  hat 0.1 using Bootstrap:", standard_error)
84
85 def muhat_median(data_a):
86     medv_v = data_a[:,13]
87     muhatmed = median(medv_v)
88     return muhatmed
89
90 def main():
91     data_a = data_loader("Boston.csv")
92     medv_v, muhat = pop_mean(data_a)
```



```

90     stand_error_muhat(medv_v, muhat)
91     standard_error = newbootstrapstderror(medv_v)
92     standard_error_median = newbootstrapstderror_median(medv_v)
93     print("standard error of median", standard_error_median)
94     print(f'Con Int: [{muhat - 2*standard_error}, {muhat + 2*standard_error}']')
95     muhatmed = muhat_median(data_a)
96     print(muhatmed)
97     tenth_percen = percentile(medv_v, 10)
98     print("Tenth Percentile ( $\mu_{0.1}$ ) of medv:", tenth_percen)
99     standard_error_percentile = newbootstrapstderror_tenpercen(medv_v)
100    print(standard_error_percentile)
101 if __name__ == '__main__':
102     main()

```