

## STATS 202: DATA MINING AND ANALYSIS HOMEWORK #3

INSTRUCTOR: LINH TRAN, HOMEWORK #3, DUE DATE: AUGUST 2, 2023,  
STANFORD UNIVERSITY, AND STUDENT: ADAM KAINIKARA

Problem 1 (7 points)

Chapter 6, Exercise 3 (p. 283).

$\sum_{i=1}^n (y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij})^2$  subject to  $\sum_{j=1}^p |\beta_j| \leq s$

The left square part is RSS

a) (iv) As we increase  $s$  from 0 the training RSS will steadily increase. Increasing  $s$  makes us restrict the  $\beta$  (coefficients) less and less. By restricting coefficients less, the model will become more flexible. When the model becomes more flexible the training RSS will decrease.

b) (ii) As we increase  $s$  from 0 the test RSS would at first decrease but then slowly start increasing and would form a  $U$  shape. Increasing  $s$  makes us restrict the  $\beta$  (coefficients) less and less. By restricting coefficients less, the model will become more flexible. However at some point on test data, the RSS will once again increase.

c) (iii) As we increase  $s$  from 0, variance would steadily increase. Increasing  $s$  makes us restrict the  $\beta$  (coefficients) less and less. By restricting coefficients less, the model becomes more flexible. As flexibility increases, variance increases. This effect is like the bias variance trade off graph.

d) (iv) As we increase  $s$  from 0, squared bias would steadily decrease. Increasing  $s$  makes us restrict the  $\beta$  (coefficients) less and less. By restricting coefficients less, the model becomes more flexible. As flexibility increases, squared bias decreases. This effect is like the bias variance trade off graph.

e) (v) As we increase  $s$  from 0, the irreducible error would stay the same. Irreducible error is always there and does not come from the fitted model. So changing the flexibility of the model will have no impact.

Problem 2 (7 points)

Chapter 6, Exercise 4 (p. 284).

$\sum_{i=1}^n (y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij})^2 + \lambda \sum_{j=1}^p \beta_j^2$

Ridge?

a) As we increase  $\lambda$  from 0 the training RSS will increase. Increasing  $\lambda$  increases the penalty term. With the penalty increasing and becoming more significant, the coefficients decrease. This leads to more error.

b) Initial, increasing  $\lambda$  shrinks the coefficients. So as it initially increases RSS will decrease since over fitting is reduced. However, if we continue to increase  $\lambda$ , the model may become too simple (due to smaller coefficients) and start to under fit, leading to an increase in the test RSS.

c) As we increase  $\lambda$  from 0 the variance will increase. With the coefficients changing as  $\lambda$  increases, the variability between how well the models can fit increases which increases the variance.

d) As we increase  $\lambda$  the bias will decrease. As  $\lambda$  increases, and the coefficients decrease, the model may become more simple. This leads to a decrease in bias.

e) Remain constant. The irreducible error represents the inherent noise in the data that cannot be reduced through modeling. As  $\lambda$  increases, irreducible error is always there.

Problem 3 (7 points)

Chapter 6, Exercise 9 (p. 286). Don't do parts (e), (f), and (g).

a) Split the data into a 75% training and 25% test.

b) Fit a linear model. Got a MSE for linear model of 1503017.4360986822 and a r squared value of 0.9134874545115684

c) Fit a ridge model. Got a MSE of 1502973.7121541149 and a r squared value of 0.9134921934779143.

d) Fit a lasso model. Got a MSE of 35152969.33526452 and 9 non zero coefficients

Coding question

Problem 4 (7 points)

Chapter 7, Exercise 1 (p. 321).

a) We are given that  $f(x) = \beta_0 + \beta_1x + \beta_2x^2 + \beta_3x^3 + \beta^4(x - \xi)^3 +$  where  $(x - \xi)^3$  is its normal polynomial self if  $x \geq \xi$  and is 0 otherwise

We are also given that  $f_1(x) = a_1 + b_1x + c_1x^2 + d_1x^3$

We want to find the polynomial and coefficients such that  $f(x) = f_1(x)$

In part 'a' it is given that  $x \leq \xi$  and because of this  $f(x) = \beta_0 + \beta_1x + \beta_2x^2 + \beta_3x^3$

So in order for  $f(x) = f_1(x)$  to be true the coefficients are:

$$a_1 = \beta_0, b_1 = \beta_1, c_1 = \beta_2, d_1 = \beta_3$$

b) Now given  $f_2(x) = a_2 + b_2x + c_2x^2 + d_2x^3$

In part 'b' it is given that  $x \geq \xi$  so  $f(x) = \beta_0 + \beta_1x + \beta_2x^2 + \beta_3x^3 + \beta_4(x - \xi)^3$   
 $= \beta_0 + \beta_1x + \beta_2x^2 + \beta_3x^3 + \beta_4(x^3 - 3x^2\xi + 3x\xi^2 - \xi^3)$

Multiply out  $\beta_4$  and rearrange the equation to follow the form of  $f_2(x)$  ie :  $x^0, x^1, x^2 \dots$

$$(\beta_0 - \beta_4\xi^3) + (\beta_1 + 3\beta_4\xi^2)x + (\beta_2 - 3\beta_4\xi)x^2 + (\beta_3 + \beta_4)x^3$$

In order for  $f(x) = f_2(x)$

$$a_2 = \beta_0 - \beta_4\xi^3, b_2 = \beta_1 + 3\beta_4\xi^2, c_2 = \beta_2 - 3\beta_4\xi, d_2 = \beta_3 + \beta_4$$

c) Show  $f_1(\xi) = f_2(\xi)$

$$f_1(\xi) = \beta_0 + \beta_1\xi + \beta_2\xi^2 + \beta_3\xi^3$$

$$f_2(\xi) = (\beta_0 - \beta_4\xi^3) + (\beta_1 + 3\beta_4\xi^2)\xi + (\beta_2 - 3\beta_4\xi)\xi^2 + (\beta_3 + \beta_4)\xi^3$$

$$\beta_0 + \beta_1\xi + \beta_2\xi^2 + \beta_3\xi^3 = (\beta_0 - \beta_4\xi^3) + (\beta_1 + 3\beta_4\xi^2)\xi + (\beta_2 - 3\beta_4\xi)\xi^2 + (\beta_3 + \beta_4)\xi^3$$

$$= \beta_0 + \beta_1\xi + \beta_2\xi^2 + \beta_3\xi^3 = \beta_0 - \beta_4\xi^3 + \beta_1\xi + 3\beta_4\xi^3 + \beta_2\xi^2 - 3\beta_4\xi^3 + \beta_3\xi^3 + \beta_4\xi^3$$

$$\beta_0 + \beta_1\xi + \beta_2\xi^2 + \beta_3\xi^3 = \beta_0 + \beta_1\xi + \beta_2\xi^2 + \beta_3\xi^3$$

d) Show :  $f'_1(\xi) = f'_2(\xi)$

$$f'_1(\xi) = \beta_1 + 2\beta_2\xi + 3\beta_3\xi^2$$

$$f'_2(\xi) = -3\beta_4\xi^3 + \beta_1 + 9\beta_4\xi^2 + 2\beta_2\xi - 9\beta_4\xi^2 + 3\beta_3\xi^2 + 3\beta_4\xi^2$$

$$f'_2(\xi) = \beta_1 + 2\beta_2\xi + 3\beta_3\xi^2$$

Therefore  $f'_1(\xi) = f'_2(\xi)$

e) Show :  $f''_1(\xi) = f''_2(\xi)$

$$f''_1(\xi) = 2\beta_2 + 6\beta_3\xi$$

$$f''_2(\xi) = 2\beta_2 + 6\beta_3\xi$$

Therefore  $f''_1(\xi) = f''_2(\xi)$

Problem 5 (7 points)

Chapter 7, Exercise 8 (p. 324). Find at least one non-linear estimate which does better than linear regression, and justify this using a t-test or by showing an improvement in the

cross-validation error with respect to a linear model. You must also produce a plot of the predictor  $X$  vs. the non-linear estimate  $\hat{f}(X)$ .

Warning: Ignoring XDG\_SESSION\_TYPE=wayland on Gnome. Use QT\_QPA\_PLATFORM=wayland to run on Wayland anyway. T-test statistic: 0.16492863090989973 P-value: 0.8690858275585926  
Linear CV error: 6.502276260524977 Polynomial CV error: 6.502276260524972

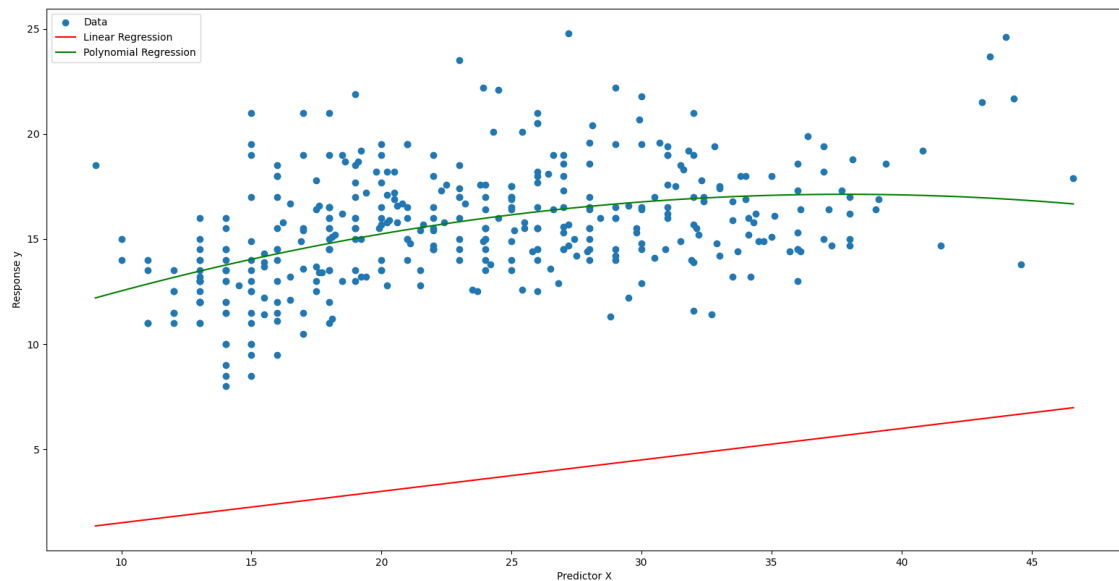


FIGURE 0.1. Poly and line

Problem 6 (7 points)  
Chapter 9, Exercise 1 (p. 398).  
Drawing a hyper plane

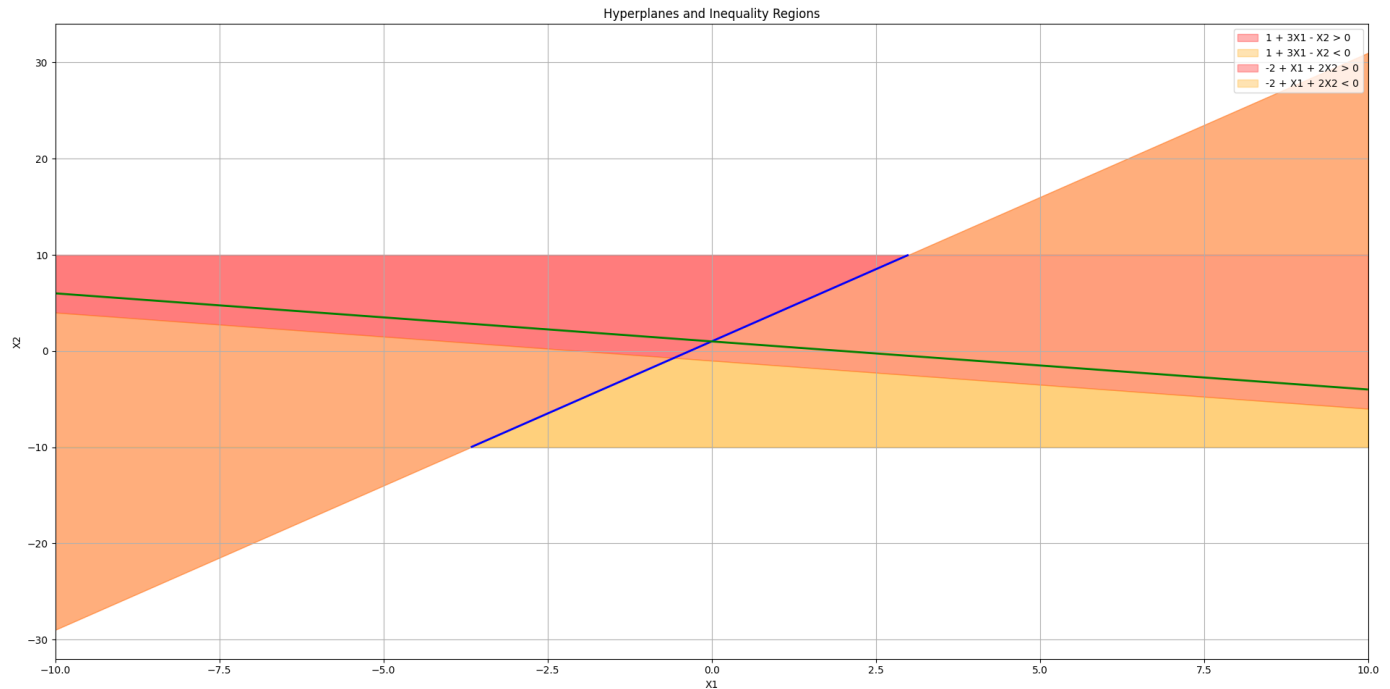


FIGURE 0.2. Hyper planes

Problem 7 (8 points)

Chapter 9, Exercise 8 (p. 401).

Note: Likely got some part of it wrong as I got the same accuracy scores for all and the best  $c$  value was 0.01.

Fitted a support vector classifier to the training data using  $C = 0.01$ , with Purchase as the response and the other variables as predictors. There were 612 support points. Training accuracy score of 0.3824999999999995 and test accuracy score of 0.4111111111111111 0.1 Fitted a support vector classifier to the training data using the best  $C = 0.1$ , and got Training accuracy score of 0.3824999999999995 and test accuracy score of 0.4111111111111111 Fitted a support vector classifier to the training data using  $C = 0.01$ , with Purchase as the response and the other variables as predictors. Used radial. There were 612 support points. Training accuracy score using radial is 0.3824999999999995 and test accuracy score of 0.4111111111111111 Fitted a support vector classifier to the training data using the best  $C = 0.01$  with radial, and got Training accuracy score of 0.3824999999999995 and test accuracy score of 0.4111111111111111 Fitted a support vector classifier to the training data using  $C = 0.01$ , with Purchase as the response and the other variables as predictors. Used poly. There were 612 support points. Training accuracy score using radial is 0.3824999999999995 and test accuracy score of 0.4111111111111111 Fitted a support vector classifier to the training data using the best  $C = 0.01$  with poly, and got Training accuracy score of 0.3824999999999995 and test accuracy score of 0.4111111111111111

# hw3

hw3/	1
└─ ch6p9.py	1
└─ ch7p8.py	3
└─ ch9p1.py	5
└─ ch9p8.py	6



## # ch6p9.py

```
1 from wsgiref.headers import tspecials
2 from numpy import *
3 import numpy as np
4 from sklearn.linear_model import LinearRegression
5 from sklearn.metrics import mean_squared_error
6 from sklearn.linear_model import RidgeCV
7 from sklearn.linear_model import LassoCV
8 from sklearn.linear_model import LassoCV
9
10
11 # THIS IS CH6 P9 FOR QUESTION 3
12 def data_loader(fname):
13     data_a = loadtxt(fname, skiprows=1, usecols=(2,3,4,5,6,7,8,9,10,11,12,13,14,15,16,17,18), delimiter=',')
14
15     return data_a
16
17 def lin_model(x_a, y_v):
18     #y_v = X@B
19     b_v = linalg.pinv(x_a)@y_v
20     return b_v
21
22 def lin_fit(x_a, b_v):
23     ypred_v = x_a@b_v
24     return ypred_v
25
26 def rid_reg(xtrain, ytrain, xtest, ytest):
27     alphas = logspace(-2, 2, 5)
28
29
30     rid_model = RidgeCV(alphas=alphas, store_cv_values=True)
31
32     rid_model.fit(xtrain, ytrain)
33     ypred = rid_model.predict(xtest)
34     rsq = 1 - (var(ytest-ypred))/(var(ytest))
35     test_error = mean((ytest-ypred)**2)
36
37     return rsq, test_error
38
39 def lass_reg(xtrain, ytrain, xtest, ytest):
40     alphas = logspace(-2, 2, 5)
41     lass_model = LassoCV(alphas=alphas)
42     lass_model.fit(xtrain, ytrain)
43     ypred = lass_model.predict(xtest)
44     rsq = 1 - (var(ytest-ypred))/(var(ytest))
45     test_error = mean((ytest-ypred)**2)
46     non_zero = sum(lass_model.coef_ != 0)
47     return rsq, test_error, non_zero
48
49
50 def main():
51     data_a = data_loader("College.csv")
52     #print(data_a)
53     #print(len(data_a))
54     '''Getting data set up'''
55     n = int(0.75*len(data_a))
56
57     xtrain = data_a[:n, 1:]
58     xtrainreal = hstack((ones((xtrain.shape[0], 1)), xtrain))
59     ytrain = data_a[:n, :1]
60
61     xtest = data_a[n:, 1:]
62     xtestreal = hstack((ones((xtest.shape[0], 1)), xtest))
63     ytest = data_a[n:, :1]
64
65     print(xtrainreal.shape)
66     print(ytrain.shape)
67     print(xtestreal.shape)
68     print(ytest.shape)
69
70     '''Doing fit'''
71     coef = lin_model(xtrainreal, ytrain)
72
73     print(coef)
74
75     pred_apps = lin_fit(xtestreal, coef)
76     print(pred_apps)
77     test_error = mean((ytest-pred_apps)**2)
78
79     rsq = 1 - (var(ytest-pred_apps))/(var(ytest))
```

```

80 #print(rsq)
81 print(f'The mse for lin reg is {test_error} and the r squared value for lin is {rsq}')
82
83 rsqrid, mserid = rid_reg(xtrainreal, ytrain, xtestreal, ytest)
84 print(f'The r squared value for ridge is {rsqrid} and the mse for ridge regression is {mserid}')
85
86 rsqlass, mselass, nonzero = lass_reg(xtrainreal, ytrain, xtestreal, ytest)
87 print(f'The r squared value for lasso is {rsqlass} and the mse for lasso is {mselass} and {nonzero} coefficents')
88
89 if __name__ == '__main__':
90     main()

```



```

1 from numpy import *
2 import matplotlib.pyplot as plt
3 from sklearn.model_selection import cross_val_score
4 from sklearn.linear_model import LinearRegression
5 from sklearn.preprocessing import PolynomialFeatures
6 from sklearn.linear_model import LinearRegression
7 from sklearn.metrics import mean_squared_error
8 from sklearn.linear_model import RidgeCV
9 from sklearn.linear_model import LassoCV
10 from sklearn.linear_model import LassoCV
11 from scipy.stats import ttest_ind
12 from scipy.stats import t
13 from numpy import *
14 import numpy as np
15 import matplotlib as mpl
16 import matplotlib.pyplot as plt
17 from scipy import stats
18 from scipy.interpolate import CubicSpline
19 from scipy.interpolate import splrep, BSpline
20 import statsmodels.api as sm
21
22 #THIS IS CH7 P8 FOR QUESTION 5
23
24 '''
25
26 To whoever is grading this problem: This file was a huge hot mess and still is a bit of a hot mess I'm sorry. I tried to clean it up a bit.
27 '''
28
29 def data_loader(fname):
30     data_a = loadtxt(fname, skiprows=1, usecols=(1,2,3,4,5,6,7), delimiter=',')
31     name_v = loadtxt(fname, skiprows=1, usecols=(0), delimiter=',')
32     return data_a, name_v
33
34
35 def lin_regression(x_a, y_v, name=''):
36     if name:
37         print(f'\n\n----- {name} -----')
38         model = sm.OLS(y_v, x_a)
39         results = model.fit()
40         print(results.summary())
41         #Using stats models to get p value even though I did my own regression
42         # y_v = X@B
43         b_v = linalg.pinv(x_a)@y_v
44         #print(b)
45         return b_v
46
47 def fitted_func(x_a, b_v):
48     yfit_v = x_a@b_v
49     # with np.printoptions(precision=2):
50     #     print(f'predicted mpg of cars {yfit_v=}')
51     return yfit_v
52
53 def r_square(y_v, yfit_v):
54     # This function is to find the r squared value
55     # This will be calculated by doing 1 - variance of (actual - predicted)/variance of actual
56     rsq = 1 - (var(y_v - yfit_v)) / (var(y_v))
57
58     return rsq
59
60 def scatter_matrix(data_a, name_l):
61     n, p = data_a.shape
62     fig, axs = plt.subplots(4, 7)
63     ax_l = list(axs.flat)
64     mpl.rcParams['figure.autolayout'] = True
65     font = {'family': 'normal',
66            'weight': 'bold',
67            'size': 10}
68
69     #mpl.rc('font', **font)
70     for i in range(p):
71         for j in range(i+1, p):
72             print(i, j, name_l[i], name_l[j])
73             x_v = data_a[:, i]
74             y_v = data_a[:, j]
75             ax = ax_l.pop(0)
76             ax.scatter(x_v, y_v, s=2**2)
77             title = f'{name_l[i]} vs {name_l[j]}'
78             ax.set_title(title[:25])
79     plt.tight_layout()
80     plt.show()
81
82 def fit_polynomial_regression(data_a, name_v, degree):
83     poly = PolynomialFeatures(degree=degree)
84     X_poly = poly.fit_transform(data_a.reshape(-1, 1))
85     model = LinearRegression()
86     model.fit(X_poly, name_v)
87     return model
88
89 def fit_polynomial_regression2(data_a, name_v, degree):
90     poly = PolynomialFeatures(degree=degree)
91     X_poly = poly.fit_transform(data_a.reshape(-1, 1))
92     model = LinearRegression()
93     model.fit(X_poly, name_v)
94     return model
95
96 def lass_reg(xtrain, ytrain, xtest, ytest):
97     alphas = logspace(-2, 2, 5)
98     lass_model = LassoCV(alphas=alphas)
99     lass_model.fit(xtrain, ytrain)
100     ypred = lass_model.predict(xtest)

```

```

105     rsq = 1 - (var(ytest-ypred))/(var(ytest))
106     test_error = mean((ytest-ypred)**2)
107     non_zero = sum(lass_model.coef_ != 0)
108     return rsq, test_error, non_zero
109
110
111 def fit_linear_regression(data_a, name_v):
112     X_linear = column_stack((ones(len(data_a)), data_a))
113     model = LinearRegression()
114     model.fit(X_linear, name_v)
115     return model
116
117 def fit_quadratic_regression(data_a, name_v):
118     X_quad = column_stack((ones(len(data_a)), data_a, data_a**2))
119     model = LinearRegression()
120     model.fit(X_quad, name_v)
121     return model
122
123 def perform_t_test(linear_model, quadratic_model, data_a, name_v):
124     linear_residuals = name_v - linear_model.predict(column_stack((ones(len(data_a)), data_a)))
125     quadratic_residuals = name_v - quadratic_model.predict(column_stack((ones(len(data_a)), data_a, data_a**2)))
126     t_stat = (linear_residuals.T @ linear_residuals - quadratic_residuals.T @ quadratic_residuals) / len(name_v)
127     return t_stat
128
129 def compute_cross_val_error(model, X, name_v):
130     cv_error = mean(cross_val_score(model, X, name_v, scoring='neg_mean_squared_error', cv=5))
131     return -cv_error
132
133 def plot_results(data_a, name_v, X_pred, linear_estimate, quadratic_estimate):
134     plt.scatter(data_a, name_v, label='Data')
135     plt.plot(X_pred, linear_estimate, label='Linear Regression', color='r')
136     plt.plot(X_pred, quadratic_estimate, label='Quadratic Polynomial Regression', color='g')
137     plt.xlabel('Predictor X')
138     plt.ylabel('Response y')
139     plt.legend()
140     plt.show()
141
142 def main():
143     data_a, name_v = data_loader('Auto.csv')
144
145     name_v = name_v.astype(float)
146
147
148     linear_model = fit_linear_regression(data_a, name_v)
149
150     quadratic_model = fit_quadratic_regression(data_a, name_v)
151
152     t_test_stat = perform_t_test(linear_model, quadratic_model, data_a, name_v)
153
154     X = column_stack((data_a, data_a**2)) # Combine linear and quadratic features
155     linear_cv_error = compute_cross_val_error(linear_model, X, name_v)
156     quadratic_cv_error = compute_cross_val_error(quadratic_model, X, name_v)
157
158     X_pred = linspace(data_a.min(), data_a.max(), 100)
159     X_pred_reshaped = column_stack((ones(100), X_pred)) # Add a column of ones for linear regression
160     linear_estimate = linear_model.predict(X_pred_reshaped)
161     quadratic_estimate = quadratic_model.predict(column_stack((ones(100), X_pred, X_pred**2)))
162     plot_results(data_a, name_v, X_pred, linear_estimate, quadratic_estimate)
163
164     print("T-test statistic:", t_test_stat)
165     print("Linear CV error:", linear_cv_error)
166     print("Quadratic CV error:", quadratic_cv_error)
167     polynomial_models = []
168     for degree in range(1, 6):
169         polynomial_model = fit_polynomial_regression(data_a, name_v, degree)
170         polynomial_models.append(polynomial_model)
171
172     linear_residuals = name_v - linear_model.predict(column_stack((ones(len(data_a)), data_a)))
173     t_test_stat, p_value = ttest_ind(linear_residuals, zeros(len(linear_residuals))) # Null hypothesis: linear model has no effect
174
175     polynomial_t_stats = []
176     polynomial_p_values = []
177     for polynomial_model in polynomial_models:
178         polynomial_residuals = name_v - polynomial_model.predict(PolynomialFeatures(degree=polynomial_model.degree).fit_transform(data_a.reshape(-1, 1)))
179         t_stat, p_value = ttest_ind(polynomial_residuals, zeros(len(polynomial_residuals))) # Null hypothesis: polynomial model has no effect
180         polynomial_t_stats.append(t_stat)
181         polynomial_p_values.append(p_value)
182
183     #x_a, y_v, data_a, name_l = data_loader('Auto.csv')
184     #cor_a = corrccoef(data_a, rowvar=False)
185     #print(cor_a.shape)
186     #with np.printoptions(precision=4):
187     #    print(cor_a)
188
189     #print(x_a,y_v)
190     #scatter_matrix(data_a, name_l)
191
192     #print(x_a,y_v)
193     #b_v = lin_regression(x_a,y_v, name='Main Regression')
194     #yfit_v = fitted_func(x_a, b_v)
195     #i_v = abs(b_v).argsort()[::-1]
196     #print(f'Coefficients {b_v=}')
197
198
199 if __name__ == "__main__":
200     main()

```

## # ch9p1.py

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3
4 # Hyperplane 1:  $1 + 3X_1 - X_2 = 0$ 
5 # Hyperplane 2:  $-2 + X_1 + 2X_2 = 0$ 
6
7 #THIS IS CHAPTER 9 P 1 FOR QUESTION 6
8 x1 = np.linspace(-10, 10, 100)
9 x2 = np.linspace(-10, 10, 100)
10 X1, X2 = np.meshgrid(x1, x2)
11 hyperplane1 = 1 + 3 * X1 - X2
12 hyperplane2 = -2 + X1 + 2 * X2
13 plt.figure(figsize=(8, 6))
14 plt.contour(X1, X2, hyperplane1, levels=[0], colors='blue', linewidths=2)
15 plt.contour(X1, X2, hyperplane2, levels=[0], colors='green', linewidths=2)
16 plt.fill_between(x1, 1 + 3 * x1, 10, color='red', alpha=0.3, label='1 + 3X1 - X2 > 0')
17 plt.fill_between(x1, -10, 1 + 3 * x1, color='orange', alpha=0.3, label='1 + 3X1 - X2 < 0')
18 plt.fill_between(x1, (-2 - x1) / 2, 10, color='red', alpha=0.3, label='-2 + X1 + 2X2 > 0')
19 plt.fill_between(x1, -10, (-2 - x1) / 2, color='orange', alpha=0.3, label='-2 + X1 + 2X2 < 0')
20 plt.xlabel('X1')
21 plt.ylabel('X2')
22 plt.legend()
23 plt.title('Hyperplanes and Inequality Regions')
24 plt.grid(True)
25 plt.show()
```

# ch9p8.py

```
1 from wsgiref.headers import tspecials
2 from numpy import *
3 import numpy as np
4 from sklearn.linear_model import LinearRegression
5 from sklearn.metrics import mean_squared_error
6 from sklearn.linear_model import RidgeCV
7 from sklearn.linear_model import LassoCV
8 from sklearn.linear_model import LassoCV
9 from sklearn.model_selection import train_test_split
10 from sklearn.svm import SVC
11 from sklearn.metrics import accuracy_score
12 from sklearn.model_selection import GridSearchCV
13
14 #THIS IS CH9 P 8 FOR QUESTION 7
15 def data_loader(fname):
16     # 1WeekofPurchase 2StoreID 3PriceCH 4PriceMM 5DiscCH 6DiscMM 7SpecialCH 8SpecialMM 9LoyalCH 10SalePriceMM 11SalePriceCH 12PriceDiff
17     # 14PctDiscMM 15PctDiscCH 16ListPriceDiff 17STORE
18     data_a = loadtxt(fname, skiprows=1, usecols=(1,2,3,4,5,6,7,8,9,10,11,12,13,14,15,16,17), delimiter=',')
19     #0 Purchase 13Store7
20     #For store 7 i replaced no with 0 and yes with 1
21     purchase_v = loadtxt(fname, skiprows=1, usecols=(0), delimiter=',', dtype=str)
22
23     return data_a, purchase_v
24
25 def sup_vec_class(xtrain, ytrain):
26     #https://scikit-learn.org/stable/modules/generated/sklearn.svm.SVC.html#sklearn.svm.SVC
27     #https://www.datacamp.com/tutorial/svm-classification-scikit-learn-python
28     #https://pythonprogramming.net/linear-svc-example-scikit-learn-svm-python/
29     sup_vec_classifier = SVC(C=0.01)
30     sup_vec_classifier.fit(xtrain, ytrain)
31     support = len(sup_vec_classifier.support_vectors_)
32     return sup_vec_classifier, support
33
34 def sup_vec_classv2(xtrain, ytrain, bestc):
35     #https://scikit-learn.org/stable/modules/generated/sklearn.svm.SVC.html#sklearn.svm.SVC
36     #https://www.datacamp.com/tutorial/svm-classification-scikit-learn-python
37     #https://pythonprogramming.net/linear-svc-example-scikit-learn-svm-python/
38     sup_vec_classifier = SVC(C=bestc)
39     sup_vec_classifier.fit(xtrain, ytrain)
40     #support = len(sup_vec_classifier.support_vectors_)
41     return sup_vec_classifier
42
43 def sup_vec_class_rbf(xtrain, ytrain):
44     sup_vec_class_radial = SVC(C=0.01, kernel='rbf')
45     sup_vec_class_radial.fit(xtrain, ytrain)
46     support = len(sup_vec_class_radial.support_vectors_)
47     return sup_vec_class_radial, support
48 def sup_vec_class_rbfv2(xtrain, ytrain, bestc):
49     #https://scikit-learn.org/stable/modules/generated/sklearn.svm.SVC.html#sklearn.svm.SVC
50     #https://www.datacamp.com/tutorial/svm-classification-scikit-learn-python
51     #https://pythonprogramming.net/linear-svc-example-scikit-learn-svm-python/
52     sup_vec_class_best_c_radial = SVC(C=bestc, kernel='rbf')
53     sup_vec_class_best_c_radial.fit(xtrain, ytrain)
54     #support = len(sup_vec_classifier.support_vectors_)
55     return sup_vec_class_best_c_radial
56
57 def sup_vec_class_poly(xtrain, ytrain):
58     sup_vec_class_poly = SVC(C=0.01, kernel='poly', degree=2)
59     sup_vec_class_poly.fit(xtrain, ytrain)
60     support = len(sup_vec_class_poly.support_vectors_)
61     return sup_vec_class_poly, support
62
63 def sup_vec_class_polyv2(xtrain, ytrain, bestpoly):
64     sup_vec_class_best_c_poly = SVC(C=bestpoly, kernel='rbf')
65     sup_vec_class_best_c_poly.fit(xtrain, ytrain)
66     #support = len(sup_vec_classifier.support_vectors_)
67     return sup_vec_class_best_c_poly
68
69 def acc_score(svmclass, xtrain, ytrain, xtest, ytest):
70     ytrainpred = svmclass.predict(xtrain)
71     trainscoretrain = 1 - accuracy_score(ytrain, ytrainpred)
72
73     ytestpred = svmclass.predict(xtest)
74     testcoretest = 1 - accuracy_score(ytest, ytestpred)
75     return trainscoretrain, testcoretest
76
77 def find_best_c(svmclass, xtrain, ytrain):
78     # this part was heavily influenced by
79     #https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.GridSearchCV.html
80     # https://scikit-learn.org/stable/modules/grid_search.html
81     #https://www.vebuso.com/2020/03/svm-hyperparameter-tuning-using-gridsearchcv/
82     #https://stats.stackexchange.com/questions/305201/optimal-grid-search-for-c-in-svm
83     #https://www.baeldung.com/cs/ml-svm-c-parameter
84     #finding c = logspace(-2, 1, 5) #5 values from 0.01 to 10
85     param_grid = {'C': [0.1, 10]}
86     grid_search = GridSearchCV(svmclass, param_grid, cv=100)
87     grid_search.fit(xtrain, ytrain)
88     best_c = grid_search.best_params_['C']
89     return best_c
90 def find_best_c_radial(svmclassradial, xtrain, ytrain):
91     #same as above but now with rbf
92     param_grid_radial = {'C': [0.01, 10]}
93     svm_cv_radial = SVC(kernel='rbf')
94     grid_search_radial = GridSearchCV(svm_cv_radial, param_grid_radial, cv=100)
95     grid_search_radial.fit(xtrain, ytrain)
96     best_c_radial = grid_search_radial.best_params_['C']
97     return best_c_radial
98
99 def find_best_c_poly(svmclasspoly, xtrain, ytrain):
100     param_grid_poly = {'C': [0.01, 10]}
101     svm_cv_poly = SVC(kernel='poly', degree=2)
102     grid_search_poly = GridSearchCV(svm_cv_poly, param_grid_poly, cv=100)
103     grid_search_poly.fit(xtrain, ytrain)
104     best_c_poly = grid_search_poly.best_params_['C']
105     return best_c_poly
106
107 def main():
108     data_a, purchase_v = data_loader('0J.csv')
109
110
111     xtrain, xtest, ytrain, ytest = train_test_split(data_a, purchase_v, train_size=800, random_state=42)
112     #print(f'xtrain{xtrain} ytrain{ytrain} xtest{xtest} ytest{ytest}')
113
114     '''part b'''
115     svmclass, support = sup_vec_class(xtrain, ytrain)
116
117     print(f'Fitted a support vector classifier to the training data using C = 0.01, with Purchase as the response and the other variables as predictors. There were {support} support points.')
118
119     '''part c'''
120     trainscore, testcore = acc_score(svmclass, xtrain, ytrain, xtest, ytest)
121     print(f'Training accuracy score of {trainscore} and test accuracy score of {testcore}')
122     '''part d'''
123     c = find_best_c(svmclass, xtrain, ytrain)
124     print(c)
125     '''part e'''
126     svmclassv2 = sup_vec_classv2(xtrain, ytrain, c)
127     newtrainscore, newtestscore = acc_score(svmclassv2, xtrain, ytrain, xtest, ytest)
128     print(f'Fitted a support vector classifier to the training data using the best C = {c}, and got Training accuracy score of {newtrainscore} and test accuracy score of {newtestscore}')
```

```

130
131 '''now doing with radial'''
132 svmclassradial, supportradial = sup_vec_class_rbf(xtrain, ytrain)
133 print(f'Fitted a support vector classifier to the training data using C = 0.01, with Purchase as the response and the other variables as predictors. Used radial. There were {supportradial} support vectors.')
134 trainscoreradial, testscoreradial = acc_score(svmclassradial, xtrain, ytrain, xtest, ytest)
135 print(f'Training accuracy score using radial is {trainscoreradial} and test accuracy score of {testscoreradial}.')
136 bestc = find_best_c_radial(svmclassradial, xtrain, ytrain)
137 svmclassradialv2 = sup_vec_class_rbfv2(xtrain, ytrain, bestc)
138 newtrainscorerad, newtestscorerad = acc_score(svmclassradialv2, xtrain, ytrain, xtest, ytest)
139 print(f'Fitted a support vector classifier to the training data using the best C = {bestc} with radial, and got Training accuracy score of {newtrainscorerad} and test accuracy score of {newtestscorerad}.')
140
141 '''now doing with poly'''
142 svmclasspoly, supportpoly = sup_vec_class_poly(xtrain, ytrain)
143 print(f'Fitted a support vector classifier to the training data using C = 0.01, with Purchase as the response and the other variables as predictors. Used poly. There were {supportpoly} support vectors.')
144 traincorepoly, testcorepoly = acc_score(svmclasspoly, xtrain, ytrain, xtest, ytest)
145 print(f'Training accuracy score using radial is {traincorepoly} and test accuracy score of {testcorepoly}.')
146 bestcpoly = find_best_c_poly(svmclasspoly, xtrain, ytrain)
147 svmclasspolyv2 = sup_vec_class_polyv2(xtrain, ytrain, bestcpoly)
148 newtraincorepoly, newtestcorepoly = acc_score(svmclasspolyv2, xtrain, ytrain, xtest, ytest)
149 print(f'Fitted a support vector classifier to the training data using the best C = {bestcpoly} with poly, and got Training accuracy score of {newtraincorepoly} and test accuracy score of {newtestcorepoly}.')
150
151 if __name__ == '__main__':
152     main()

```