

# Exam 1 A - MATH 4322

Instructor - Cathy Poliak

Fall 2022

Name: \_\_\_\_\_

PSID: \_\_\_\_\_

## Instructions

- Allow one sheet of notes front and back to be turned in for extra credit.
- Allow calculator.
- Total possible points 100.
- For multiple choice circle your answer on this test paper.
- For short answer questions answer fully on this test paper, partial credit will be given.
- Once completed turn in to TA or instructor.
- Data sets are coming from

[UCI Machine Learning Repository](#)

## Problem 1

(36 possible points) We want to understand how the input variables relate to miles per gallon, `mpg`. The input variables are:

- `cylinders` - as qualitative 4, 6 or 8
- `displacement` - cubic inches
- `horsepower` - gross horsepower
- `weight` - per 1000 pounds

- a. Is this a inference or prediction statistical learning problem?
- b. Is this a regression or classification problem?
- c. Give the model formula for our problem. Use the variable names in the formula.
- d. Give the R code to get the model for predicting the `mpg` based on the 4 input variables.

e. The following is the output from the data. Write out the equation with the estimates.

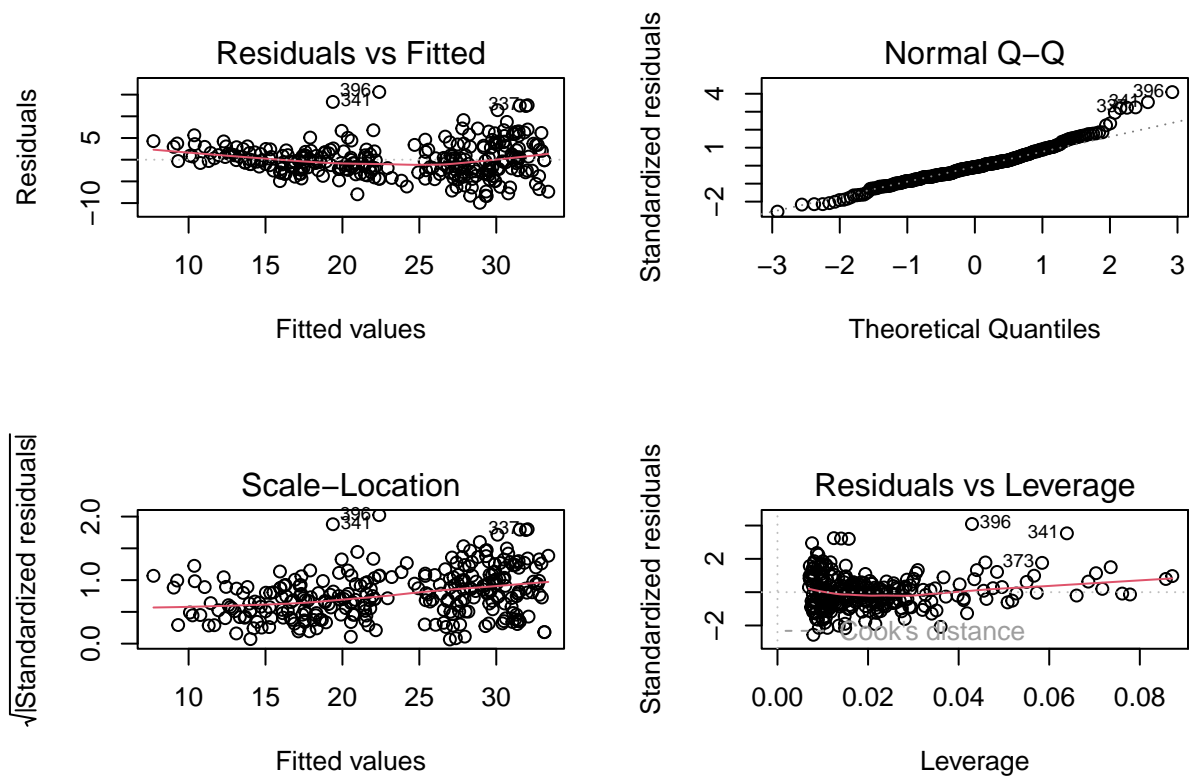
	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	44.5313	1.5297	29.1111	0.0000
cylinders6	-4.1951	1.0709	-3.9174	0.0001
cylinders8	-2.6522	1.8066	-1.4681	0.1432
displacement	0.0090	0.0104	0.8617	0.3896
horsepower	-0.0576	0.0152	-3.7828	0.0002
weight	-5.1079	0.8346	-6.1202	0.0000

f. Give the interpretation of the coefficient for the variable **horsepower**.

g. Are there any variables that are not needed in this model? Justify your answer.

h. What are the assumptions of this model?

- i. The plot below are the diagnostics plots. Are any of the assumptions violated with this model?



## Problem 2

(32 possible points) We want to predict whether a person will donate blood or not. The variables are:

- **Monetary** - total blood donated in c.c per 1000.
- **Recency** - months since last donation.
- **Donate** - a binary variable representing whether he/she donated blood (1 stand for donating blood; 0 stands for not donating blood).

a. Is this a inference or prediction statistical learning problem?

b. Is this a regression or classification problem?

c. Give the model formula for our problem. Use the variable names in the formula.

d. Give the R code to get the model for predicting the probability of donating blood based on the 2 input variables.

e. The following is the output from the data. Write out the equation with the estimates.

	Estimate	Std. Error	z value	Pr(> z )
(Intercept)	-0.5643	0.1988	-2.84	0.0045
Recency	-0.1228	0.0190	-6.47	0.0000
Monetary	0.2616	0.0769	3.40	0.0007

f. Give the predicted probability of donating blood for a donor that has donated 1400 c.c. of blood and last donation was 4 months ago.

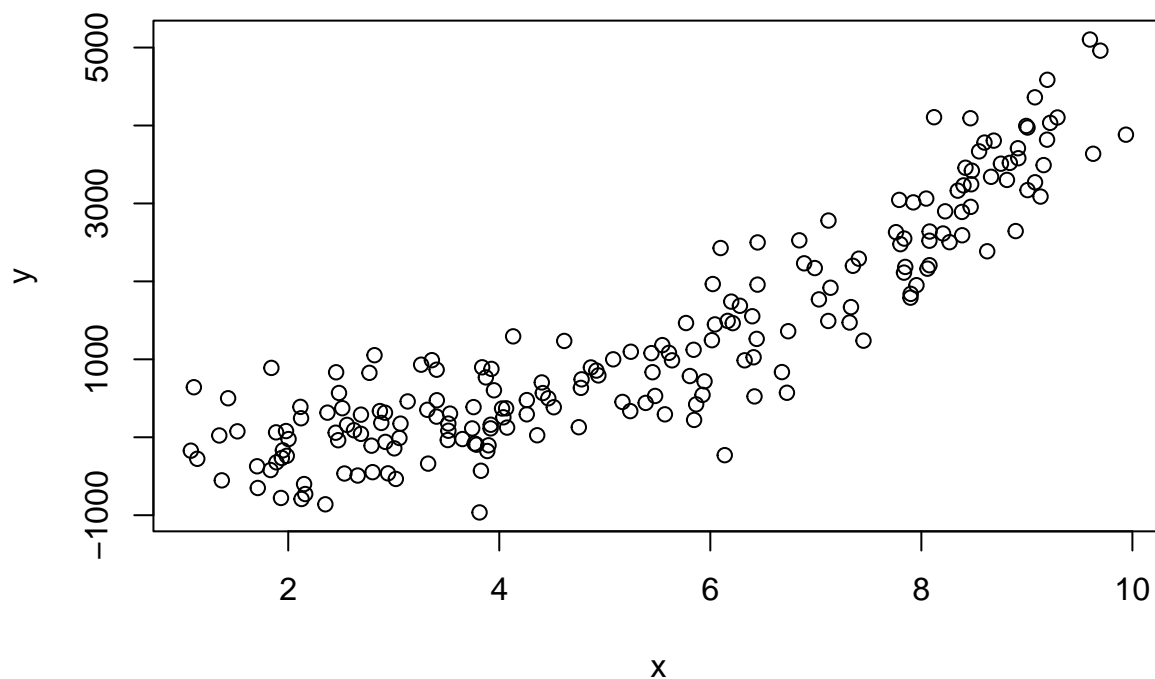
g. The following is the output from R. Determine  $R^2$  and give an interpretation.

Null deviance:	625.95	on	560	degrees of freedom
Residual deviance:	550.16	on	558	degrees of freedom

### Problem 3

(8 possible points)

- a. Using the following plot below do we have a linear relationship?



- b. The following is an output for a regression model with degree 1, 2, 3 and 4 respectively, based on the data represented from the plot above. According to these statistics, write out the formula for the best model.

	Adj.R2	Cp	BIC
Degree 1	0.81	114.39	-322.54
Degree 2	0.88	2.57	-408.60
Degree 3	0.88	3.21	-404.69
Degree 4	0.88	5.00	-399.61

### Problem 4

(8 points) A graduate program is making decisions to admit students into the program with the variables **GPA**, and the score on the **GRE**. The response variable is **Decision**, there are three decisions that are made; *yes*, *no*, and *conditional*.

- a. Circle the best model to use for this example.
- i. Simple Linear Regression
  - ii. Logistic Regression
  - iii. Multiple Linear Regression
  - iv. Linear Discriminat Analysis (LDA)
  - v. Polynomial Regression
- b. The following is the confusion matrix based on the model. What is the error rate?

	Yes	No	Conditional
Yes	21	0	3
No	0	21	2
Conditional	1	0	20

- i. 0.0882
- ii. 0.875
- iii. 0.913
- iv. 0.9118
- v. 0.9524

### Problem 5

(4 points) The following is the ANOVA table from problem 1, where  $n = 288$  and the MSE for the full model from problem 1 is 15.2. What is the  $C_p$  statistic?

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
horsepower	1	10782.97	10782.97	656.78	0.0000
weight	1	1868.22	1868.22	113.79	0.0000
Residuals	285	4679.13	16.42		

- a. 425.38
- b. -161.09
- c. 4
- d. 27.83
- e. 23.83



### Problem 6

(4 points) Suppose we have  $p = 4$  predictors. How many possible additive models contain subsets of the 4 predictors?

- a. 4
- b. 8
- c. 16
- d. 32
- e. 100

### Problem 7

(4 points) Which stepwise selection begins with a model containing no predictors, and then adds predictors to the model, one-at-a-time, until all of the significant predictors are in the model.

- a. forward
- b. backward
- c. best subset
- d. none of these

### Problem 8

(4 points) The following is a 95% prediction interval for the `mpg` from problem 1, with only `weight` as the predictor. We wanted to predict where `weight` is 2845 pounds. Which statement is correct?

	fit	lwr	upr
1	24.50	16.41	32.58

- a. For one automobile that weighs 2845, we predict the `mpg` to be between 16.41 and 32.58 with 95% confidence.
- b. On average for all automobiles that weigh 2845, we we predict the `mpg` to be between 16.41 and 32.58 with 95% confidence.
- c. For one automobile that regardless of the weight, we predict the `mpg` to be between 16.41 and 32.58 with 95% confidence.
- d. On average for all automobiles regardless of the weight, we we predict the `mpg` to be between 16.41 and 32.58 with 95% confidence.
- e. For an automobile that weights 2845 pounds, the `mpg` will be 24.5.