

Lab 2 (4 points)

Linear Regression

This lab consists of two parts, which use linear regression for prediction and statistical inference.

Part I

This question involves the use of multiple linear regression on the Auto data set.

- Produce a scatterplot matrix, which includes all of the variables in the data set.
- Compute the matrix of correlations between the variables using the function `cor()`. You will need to exclude the name variable, which is qualitative.
- Use the `lm()` function to perform a multiple linear regression with `mpg` as the response and all other variables except `name` as the predictors. Use the `summary()` function to print the results. Comment on the output. For instance:
 - Is there a relationship between the predictors and the response?
 - Which predictors appear to have a statistically significant relationship to the response?
 - What does the coefficient for the year variable suggest?
- Use the `plot()` function to produce diagnostic plots of the linear regression fit. Comment on any problems you see with the fit. Do the residual plots suggest any unusually large outliers? Does the leverage plot identify any observations with unusually high leverage?
- Use the `*` and `:` symbols to fit linear regression models with interaction effects. Do any interactions appear to be statistically significant?
- Try a few different transformations of the variables, such as $\log(X)$, $X^{0.5}$, X^2 . Comment on your findings.

Part II

This question should be answered using the Carseats data set, which is part of the ISLR package.

- Fit a multiple regression model to predict Sales using Price, Urban, and US.
- Provide an interpretation of each coefficient in the model. Be careful—some of the variables in the model are qualitative!
- Write out the model in equation form, being careful to handle the qualitative variables properly.
- For which of the predictors can you reject the null hypothesis $H_0: \beta_j = 0$?
- On the basis of your response to the previous question, fit a smaller model that only uses the predictors for which there is evidence of association with the outcome.
- How well do the models in (a) and (e) fit the data?
- Using the model from (e), obtain 95% confidence intervals for the coefficient(s).