

## Homework # 5

In this Homework, you will learn practice simple & multiple linear regression analysis for prediction. You will also practice polynomial regression for prediction. Feature selection using PCA & random forest technique will also be discussed in this assignment. The numbers in graphs in this LAB may be only for the display and may not be accurate.

The dataset that is selected for this assignment is '[Carseat.csv](#)' which can be found on the course GitHub. All numbers must be displayed with 3-digit decimal precision. Add x-label, y-label, grid, title, and legend [if applicable] to all submitted plots.

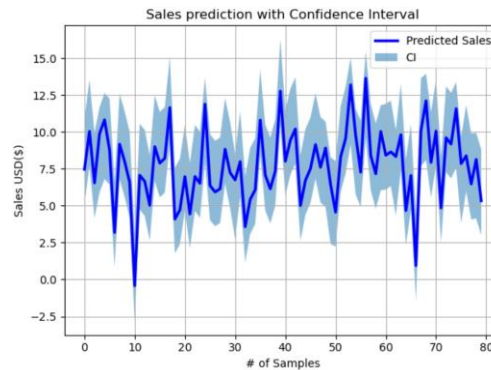
The '[Carseat.csv](#)' dataset contains several features (qualitative and quantitative) and the dependent variable 'Sales' (child car seat) in 400 locations. The goal is to predict the sales using the regression model. Use random\_state = 5805 as needed.

1. **EDA**: Develop a python program that: [15 points]
  - a. Read the dataset and plot the group horizontal bar plot that shows the 'ShelveLoc' versus Sales with respect to the 'US' (legend). Which Shelve location has the highest sales and whether the best sales are inside the US or outside? Hint: You need to aggregate using summation method.
  - b. Using the one-hot-encoding converts the qualitative features to quantitative features and makes it ready for regression analysis. Avoid the dummy trap. Display the first 5 rows of the converted features on the console.
  - c. Split the dataset into train-test 80-20 then standardize the dataset. Turn the shuffle ON and use the random\_state = 5805. Display the first 5 rows of the train and test set. Hint: The encoded features should not be standardized. Use the following for the split. Set the shuffle = True.

```
from sklearn.model_selection import train_test_split
```

2. **Feature selection & Prediction**: Develop a python program that performs a backward stepwise regression analysis and eliminate features using p-value of t-test (threshold 0.01). Hint: Drop one feature at a time while monitoring Adjusted R-squared [20 points]
  - a. Create a table that shows the process and justifies the one-by-one elimination. You need to display the AIC, BIC and Adjusted  $R^2$  as a predictive accuracy for each elimination inside the table. Display the eliminated and final selected features on the console. Display the p-value (of a feature to be eliminated) of the feature inside the table.

- b. Drop the insignificant features and perform a regression analysis OLS on the selected features. Take a screen shot of the OLS summary and place it here after each elimination.
  - c. After removing insignificant features, write the final regression model. Make a prediction on sales and compare it with the test set and plot the original test set (sales without transformation) versus the predicted values (not a scatter plot). Hint: You need to perform a reverse transformation.
  - d. Display the Mean Squared Error (MSE) on this prediction on the console.
- 3. **Feature selection & Prediction**: Develop a python program that performs a Principal Component Analysis (PCA). [15 points]
  - a. How many features are needed that explain more than 95% of the dependent variance.
  - b. Plot the cumulative explained variance versus the number of features. Hint: The number of features lie on the x-axis and must start off at 1.
  - c. Draw a vertical line and horizontal line showing the exact 95% threshold and the corresponding number of features.
  - d. What does the PCA say? Do the PCA tell you which feature(s) need to remove? Explain your answer.
- 4. **Feature selection & Prediction**: Develop a python program that performs a Random Forest Analysis and eliminates features based on the importance. [25 points]
  - a. Plot the horizontal bar graph (once random forest analysis is complete) that displays the features importance in descending order. Hint: Make sure that the y-axis is labeled according to the feature labels.
  - b. Display the eliminated and final selected features on the console. Do the selected features base on random forest and stepwise regression method identical? Explain your answer.
  - c. Drop the insignificant features (pick a threshold) and perform a regression analysis OLS on the selected features. Take a screen shot of the OLS summary and place it here.
  - d. Make a prediction using the test set and plot the original test set (sales without transformation) versus the predicted values (not a scatter plot). Hint: You need to perform a reverse transformation.
  - e. Display the Mean Squared Error (MSE) on this prediction on the console.
- 5. **Create a table that compares** the R-squared, Adjusted R-squared, AIC, BIC and MSE of the prediction in step 2 and 4. Which method of feature selection will you recommend for this problem and what features do you recommend for elimination? Explain your answer. [10 points]
- 6. **Prediction Interval**: Based on the answer in the previous question and the selected features using stepwise regression, find the prediction intervals for the regression model. Plot the predicted sales value (need to perform reverse transformation), the lower and upper prediction intervals (95%) in one graph. Hint: You can use `.get_prediction()` function. [10 points]

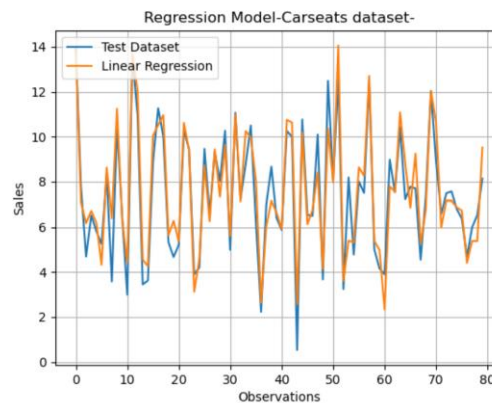


7. **Polynomial regression and Grid Search:** Let suppose we want to find a polynomial regression model that regresses Sales (dependent variable) versus the “price” feature only. Hint: You may use the following packages: [25 points]

```
from sklearn.model_selection import GridSearchCV
from sklearn.linear_model import LinearRegression
from sklearn.preprocessing import PolynomialFeatures
from sklearn.pipeline import make_pipeline
```

$$y = \beta_0 + \beta_1 x + \beta_2 x^2 + \beta_3 x^3 + \dots + \beta_n x^n + \varepsilon.$$

- Perform a grid search that minimized RMSE.
- What is the optimum order for n?
- Plot the RMSE versus the n order. The search intervals [1,15].
- Split the dataset into train-test 80-20 random\_state =5805. Train the regression model using OLS and the optimum nth order derived from the previous section. Plot the test set (sales) versus the predicted sales values.
- What is the MSE of this nth other polynomial regression prediction?



8. In a simple linear regression with n observations. Hint: You need to construct the sum of squared of error function and minimize the function with respect to beta\_0 and beta\_1. This requires a derivative. [10 points]

$$y_i \approx \hat{\beta}_0 + \hat{\beta}_1 x_i$$

Prove the following:

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

where  $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$  and  $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$  are the sample mean.

Upload the solution as a **report (as a single pdf)** plus **the .py file** through Canvas by the due date.