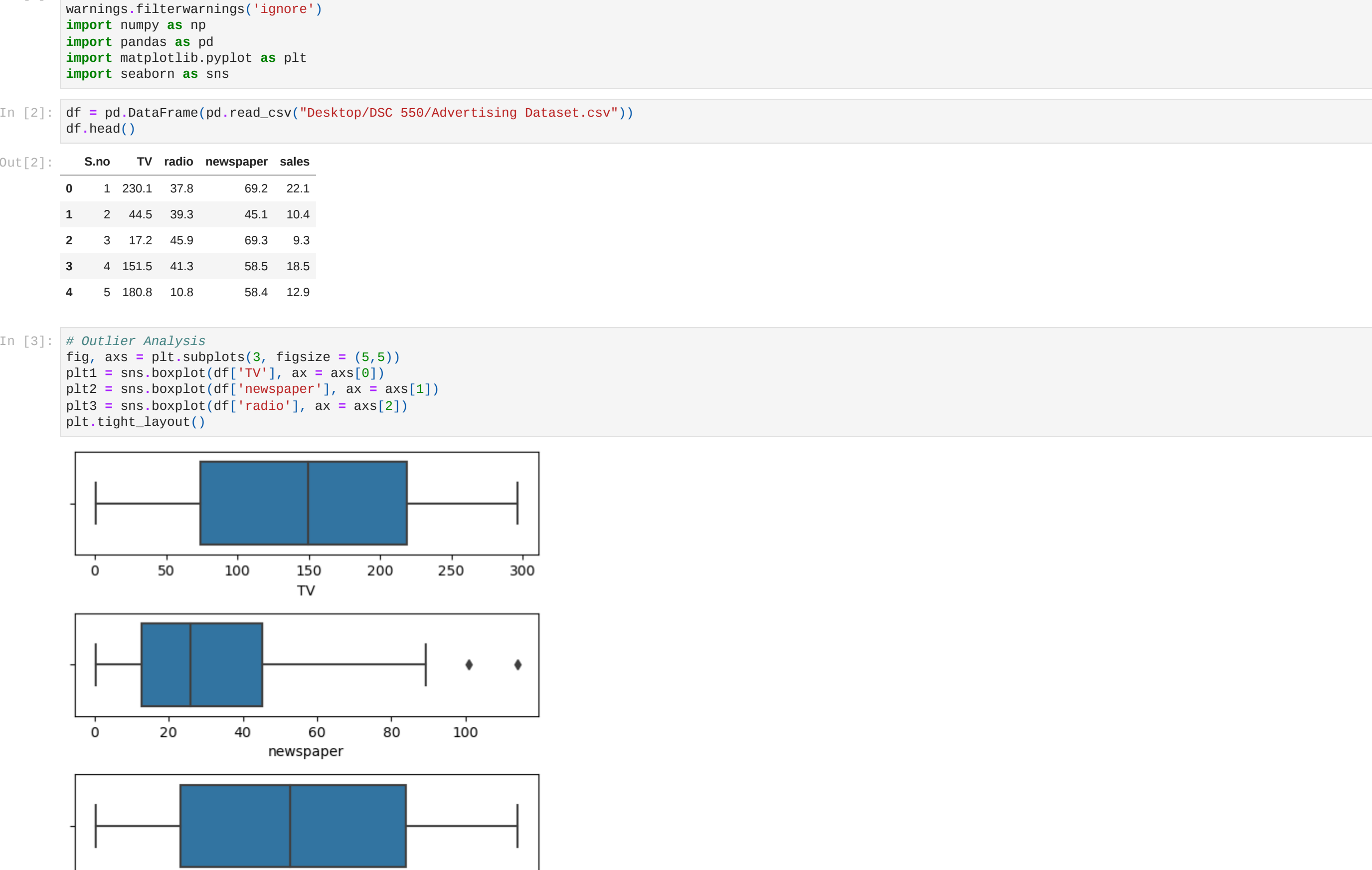
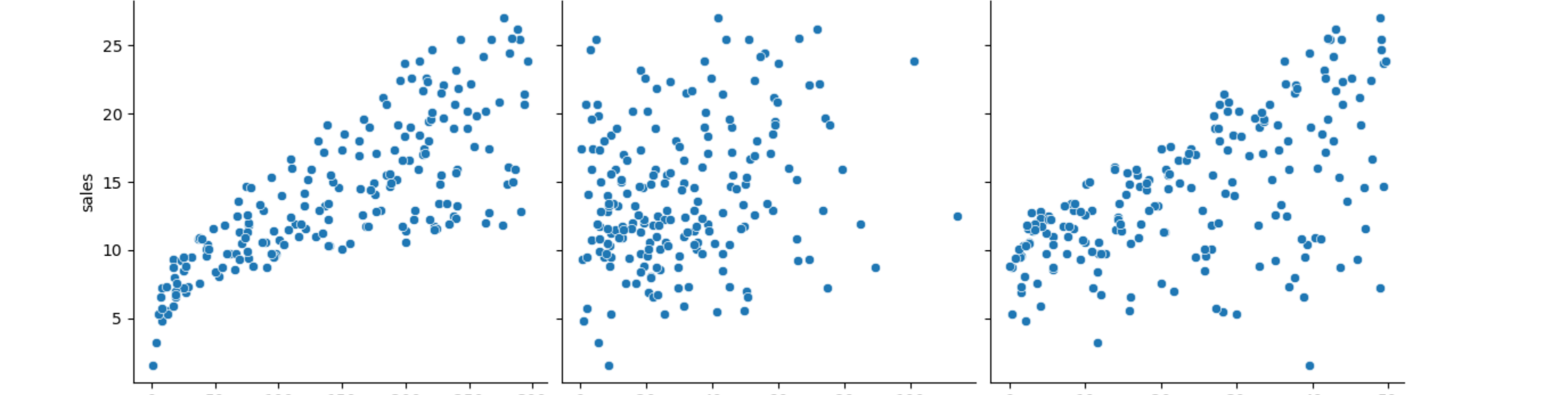
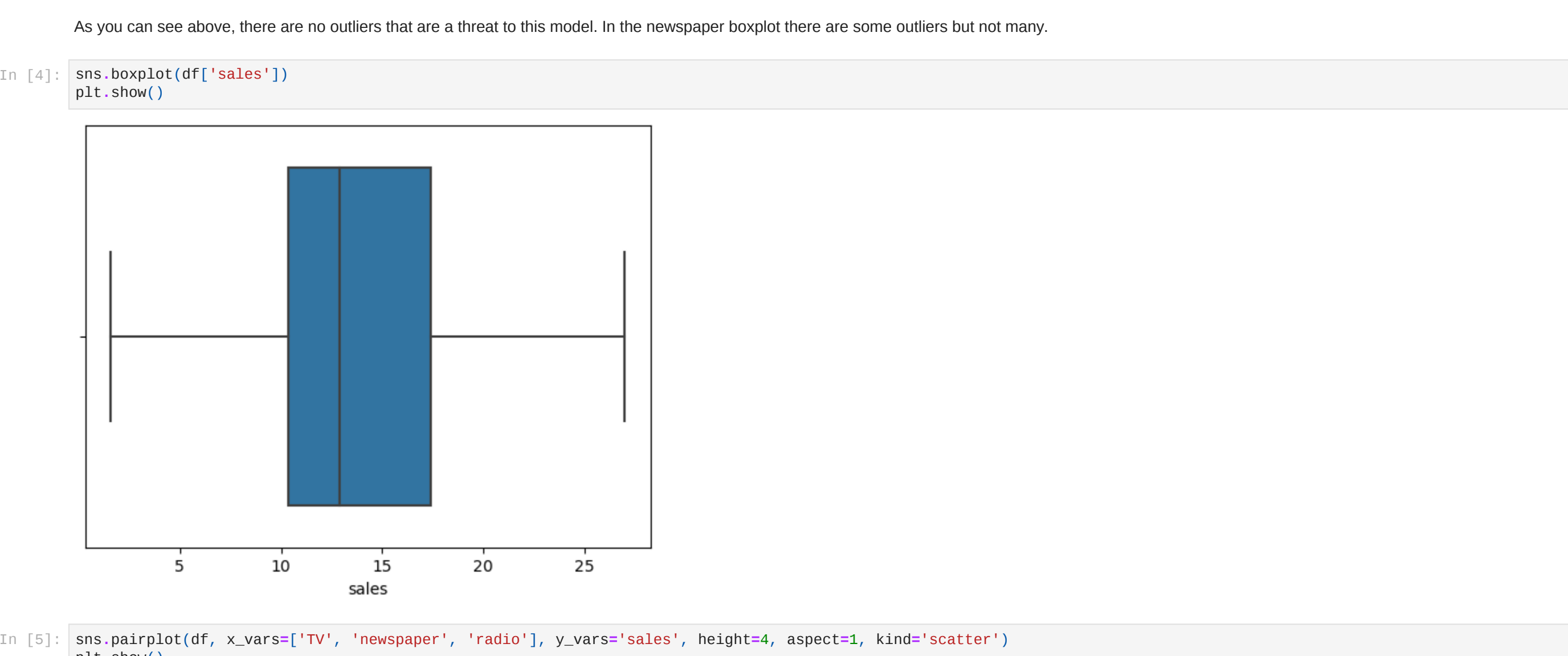


Project Milestone 1: Data Selection and EDA

Business Problem: The first step of this project was to identify a use for data mining in a business setting. I have settled on advertising expenditures vs sales. I have always enjoyed and appreciated regression models so I plan on building a regression model to predict the efficiency of different advertising channels. Every company contains some sort of advertising, and these different channels of advertising need to be looked at and utilized correctly to maximize sales. By looking at sales data, and advertising data one can build a regression model to predict how much sales each channel will generate based on how much was spent on advertising. The data I will be looking at is going to be advertising expenditures from three different channels. These different channels are, "Tv advertising", "Radio Advertising", and "Newspaper Advertising". Any company can use a multiple linear regression model to estimate a relationship between advertising spending and sales. The goal of this project is to show the potential of each respective channel of advertising. The model needs to be trained from historical data. By creating a model the company can then optimize the spending across the channels to maximize its revenue stream. They can also use the model to monitor the effectiveness of their advertising spending over time and adjust their strategies accordingly. Additionally, the regression model could be used to forecast revenue based on different advertising spending scenarios, enabling the company to make data-driven decisions about its advertising budget allocation. Analyzing advertising data is essential for any company because this can reveal underlying issues in a revenue stream. For example, newspaper and radio ads are not used as often as Facebook ads or e-commerce ads. So why do some companies continue advertising through newspapers or radio? By looking at the data we can see that newspapers and radio are still being watched by the consumer. We cannot make decisions without looking at data.



As you can see above, there are no outliers that are a threat to this model. In the newspaper boxplot there are some outliers but not many.

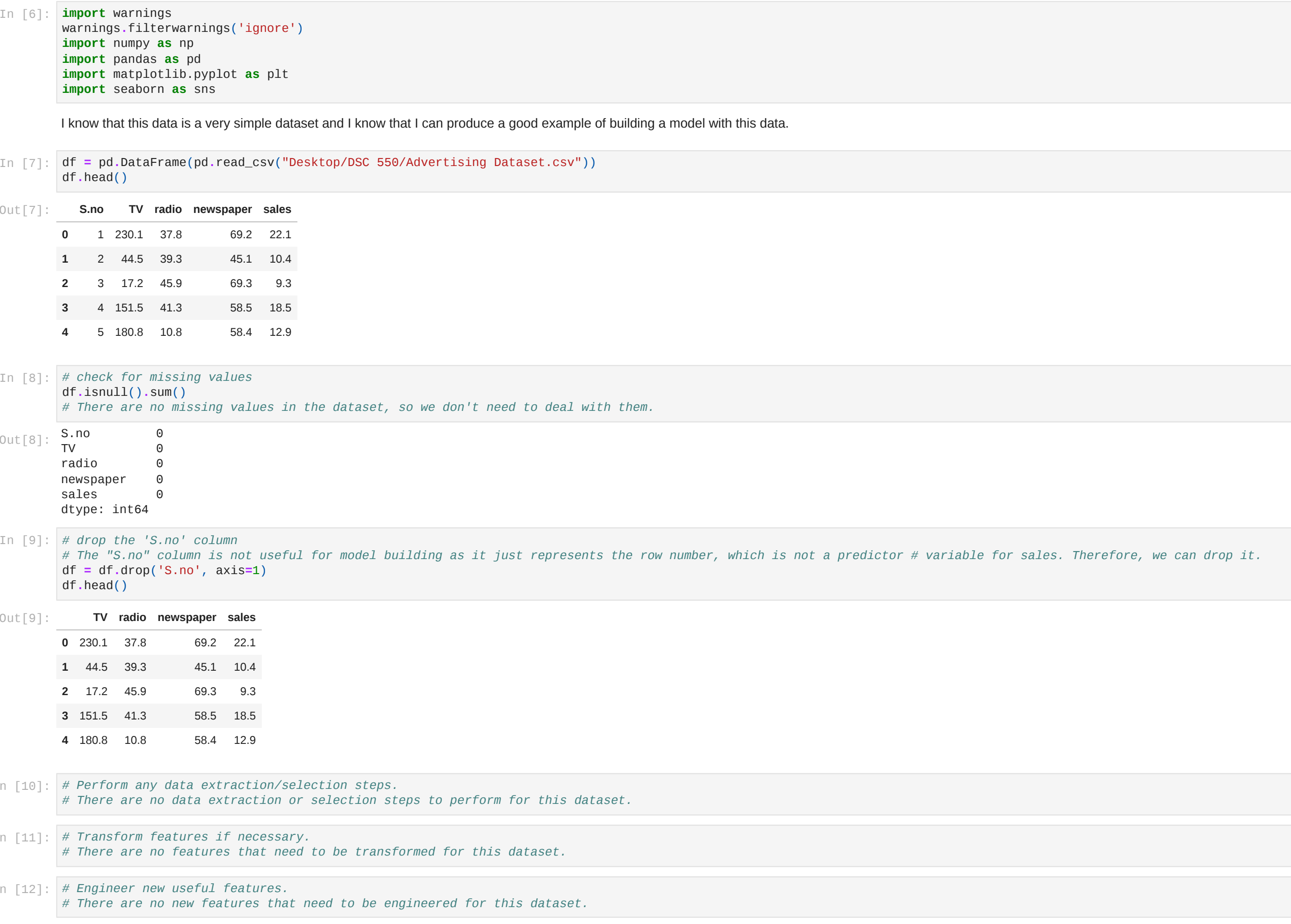


As you can see above, the most correlated variable with sales seems to be TV. The other variables have very low expenditures in advertising but seem to be vaguely correlated with sales.

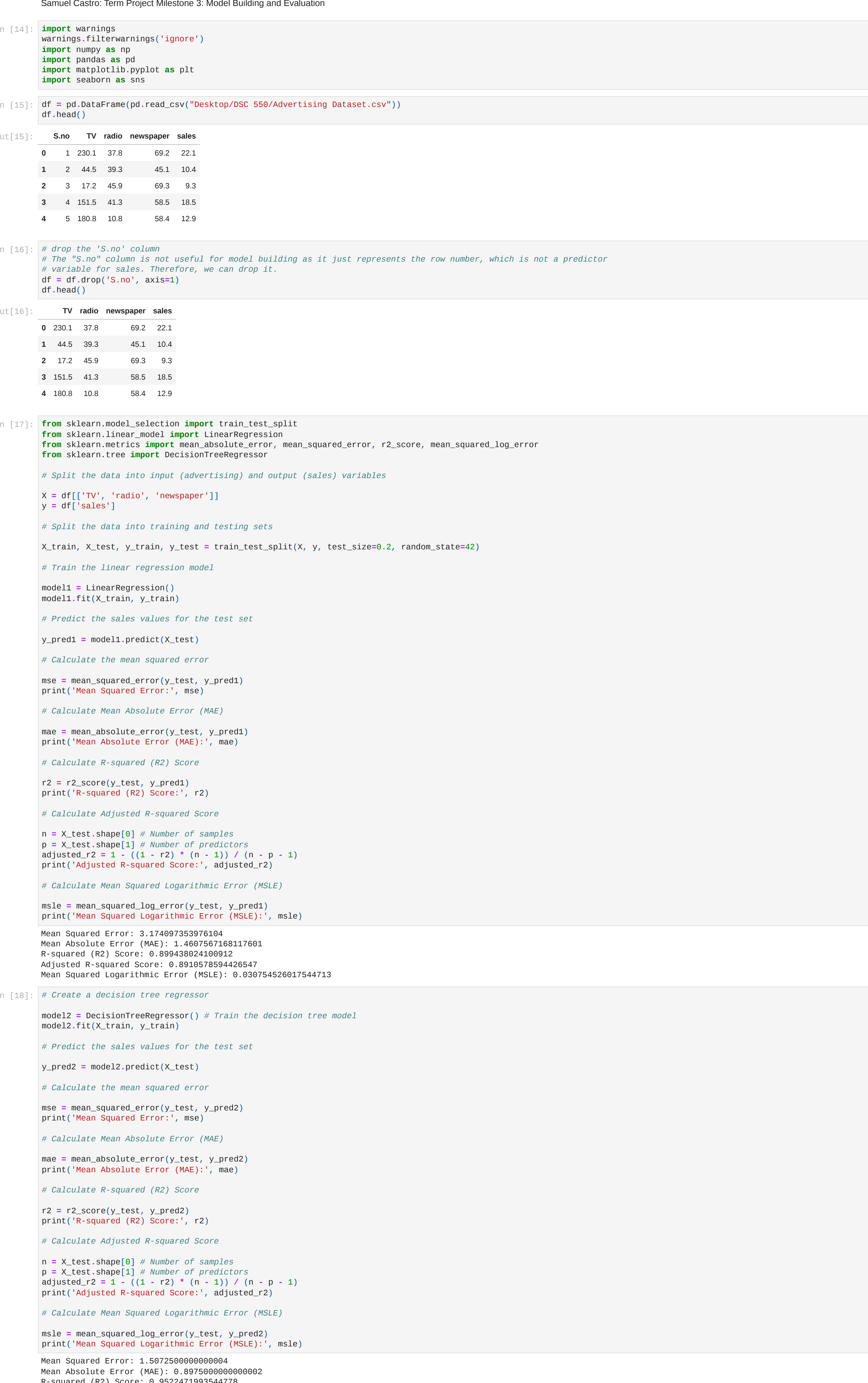
The conclusion on my very brief EDA is that each variable is correlated with sales in their respective ways. Tv seems to be the most correlated but our goal is to see which advertising stream generates the most revenue and which stream has the potential to create the most revenue.

Samuel Castro Project Mileston #2

Now that you have created your idea, located data, and have started your graphical analysis, you will move on to the data preparation process of your project. After completing Milestone 2, your data should be ready for the model building/evaluation phase.



Samuel Castro: Term Project Milestone 3: Model Building and Evaluation



I used both of these models because they have different strengths. Linear regression assumes a linear relationship between features and the target, provides interpretable coefficients, and is sensitive to outliers. Decision tree regressor makes no assumptions about linearity, can handle non-linear relationships, and may be less interpretable but can handle outliers and capture complex patterns. I chose these two because they are opposites when it comes to linear relationships. We have 2 extremes. One assuming a linear relationship and another making no assumptions.

As you can see with the results that the decision tree regressor is a better model because of the MSE results. The MSE measures the average squared difference between the predicted and actual values. It means that the predictions are off by the output. The lower the output of MSE the closer the model is to the actual sales values.

The MAE calculates the average absolute difference between the predicted and actual values. With a MAE of 0.955, it means that, on average, the absolute difference between the predicted and actual values of the target variable is 0.955.

The R-squared score measures the proportion of the variance in the target variable that is explained by the model. An R2 score of 0.9498 indicates that approximately 94.98% of the variance in the target variable can be explained by the linear regression model. The higher the r2 score is the better the model is. The r2 result for both is about the same.

The adjusted R-squared adjusts the R2 score by taking into account the number of predictors and the sample size. It penalizes the addition of unnecessary predictors that do not improve the model's performance. Higher adjusted R2 scores suggest a stronger fit of the model to the data. The R*2 was higher in the decision tree as well.

The MSLE measures the average logarithmic squared difference between the predicted and actual values. It is useful when the target variable has exponential growth patterns or spans a large range. The MSLE value of 0.0148 indicates that, on average, the logarithmic squared difference between the predicted and actual values is 0.0148.