

Activity_Perform multiple linear regression

June 5, 2024

1 Activity: Perform multiple linear regression

1.1 Introduction

As you have learned, multiple linear regression helps you estimate the linear relationship between one continuous dependent variable and two or more independent variables. For data science professionals, this is a useful skill because it allows you to compare more than one variable to the variable you're measuring against. This provides the opportunity for much more thorough and flexible analysis.

For this activity, you will be analyzing a small business' historical marketing promotion data. Each row corresponds to an independent marketing promotion where their business uses TV, social media, radio, and influencer promotions to increase sales. They previously had you work on finding a single variable that predicts sales, and now they are hoping to expand this analysis to include other variables that can help them target their marketing efforts.

To address the business' request, you will conduct a multiple linear regression analysis to estimate sales from a combination of independent variables. This will include:

- Exploring and cleaning data
- Using plots and descriptive statistics to select the independent variables
- Creating a fitting multiple linear regression model
- Checking model assumptions
- Interpreting model outputs and communicating the results to non-technical stakeholders

1.2 Step 1: Imports

1.2.1 Import packages

Import relevant Python libraries and modules.

```
[ ]: # Import libraries and modules.  
  
### YOUR CODE HERE ###
```

1.2.2 Load dataset

Pandas was used to load the dataset `marketing_sales_data.csv` as `data`, now display the first five rows. The variables in the dataset have been adjusted to suit the objectives of this lab. As shown in this cell, the dataset has been automatically loaded in for you. You do not need to download the .csv file, or provide more code, in order to access the dataset and proceed with this lab. Please continue with this activity by completing the following instructions.

```
[ ]: # RUN THIS CELL TO IMPORT YOUR DATA.  
  
### YOUR CODE HERE ###  
data = pd.read_csv('marketing_sales_data.csv')  
  
# Display the first five rows.  
  
### YOUR CODE HERE ###
```

1.3 Step 2: Data exploration

1.3.1 Familiarize yourself with the data's features

Start with an exploratory data analysis to familiarize yourself with the data and prepare it for modeling.

The features in the data are:

- TV promotional budget (in “Low,” “Medium,” and “High” categories)
- Social media promotional budget (in millions of dollars)
- Radio promotional budget (in millions of dollars)
- Sales (in millions of dollars)
- Influencer size (in “Mega,” “Macro,” “Micro,” and “Nano” categories)

Question: What are some purposes of EDA before constructing a multiple linear regression model?

[Write your response here. Double-click (or enter) to edit.]

1.3.2 Create a pairplot of the data

Create a pairplot to visualize the relationship between the continuous variables in `data`.

```
[ ]: # Create a pairplot of the data.  
  
### YOUR CODE HERE ###
```

Hint 1

Refer to [the content](#) where creating a pairplot is demonstrated.

Hint 2

Use the function in the **seaborn** library that allows you to create a pairplot showing the relationships between variables in the data.

Hint 3

Use the `pairplot()` function from the **seaborn** library and pass in the entire DataFrame.

Question: Which variables have a linear relationship with **Sales**? Why are some variables in the data excluded from the preceding plot?

[Write your response here. Double-click (or enter) to edit.]

1.3.3 Calculate the mean sales for each categorical variable

There are two categorical variables: **TV** and **Influencer**. To characterize the relationship between the categorical variables and **Sales**, find the mean **Sales** for each category in **TV** and the mean **Sales** for each category in **Influencer**.

```
[ ]: # Calculate the mean sales for each TV category.

### YOUR CODE HERE ###

# Calculate the mean sales for each Influencer category.

### YOUR CODE HERE ###
```

Hint 1

Find the mean **Sales** when the **TV** promotion is **High**, **Medium**, or **Low**.

Find the mean **Sales** when the **Influencer** promotion is **Macro**, **Mega**, **Micro**, or **Nano**.

Hint 2

Use the `groupby` operation in **pandas** to split an object (e.g., `data`) into groups and apply a calculation to each group.

Hint 3

To calculate the mean **Sales** for each **TV** category, group by **TV**, select the **Sales** column, and then calculate the mean.

Apply the same process to calculate the mean **Sales** for each **Influencer** category.

Question: What do you notice about the categorical variables? Could they be useful predictors of **Sales**?

[Write your response here. Double-click (or enter) to edit.]

1.3.4 Remove missing data

This dataset contains rows with missing values. To correct this, drop all rows that contain missing data.

```
[ ]: # Drop rows that contain missing data and update the DataFrame.  
  
    ### YOUR CODE HERE ###
```

Hint 1

Use the **pandas** function that removes missing values.

Hint 2

The **dropna()** function removes missing values from an object (e.g., **DataFrame**).

Hint 3

Use **data.dropna(axis=0)** to drop all rows with missing values in **data**. Be sure to properly update the **DataFrame**.

1.3.5 Clean column names

The **ols()** function doesn't run when variable names contain a space. Check that the column names in **data** do not contain spaces and fix them, if needed.

```
[ ]: # Rename all columns in data that contain a space.  
  
    ### YOUR CODE HERE ###
```

Hint 1

There is one column name that contains a space. Search for it in **data**.

Hint 2

The **Social Media** column name in **data** contains a space. This is not allowed in the **ols()** function.

Hint 3

Use the **rename()** function in **pandas** and use the **columns** argument to provide a new name for **Social Media**.

1.4 Step 3: Model building

1.4.1 Fit a multiple linear regression model that predicts sales

Using the independent variables of your choice, fit a multiple linear regression model that predicts **Sales** using two or more independent variables from **data**.

```
[ ]: # Define the OLS formula.

### YOUR CODE HERE ###

# Create an OLS model.

### YOUR CODE HERE ###

# Fit the model.

### YOUR CODE HERE ###

# Save the results summary.

### YOUR CODE HERE ###

# Display the model results.

### YOUR CODE HERE ###
```

Hint 1

Refer to the content that discusses [model building](#) for linear regression.

Hint 2

Use the `ols()` function imported earlier—which creates a model from a formula and `DataFrame`—to create an OLS model.

Hint 3

You previously learned how to specify in `ols()` that a feature is categorical.

Be sure the string names for the independent variables match the column names in `data` exactly.

Question: Which independent variables did you choose for the model, and why?

[Write your response here. Double-click (or enter) to edit.]

1.4.2 Check model assumptions

For multiple linear regression, there is an additional assumption added to the four simple linear regression assumptions: **multicollinearity**.

Check that all five multiple linear regression assumptions are upheld for your model.

1.4.3 Model assumption: Linearity

Create scatterplots comparing the continuous independent variable(s) you selected previously with **Sales** to check the linearity assumption. Use the pairplot you created earlier to verify the linearity assumption or create new scatterplots comparing the variables of interest.

```
[ ]: # Create a scatterplot for each independent variable and the dependent variable.  
  
### YOUR CODE HERE ###
```

Hint 1

Use the function in the **seaborn** library that allows you to create a scatterplot to display the values for two variables.

Hint 2

Use the **scatterplot()** function in **seaborn**.

Hint 3

Pass the independent and dependent variables in your model as the arguments for **x** and **y**, respectively, in the **scatterplot()** function. Do this for each continuous independent variable in your model.

Question: Is the linearity assumption met?

[Write your response here. Double-click (or enter) to edit.]

1.4.4 Model assumption: Independence

The **independent observation assumption** states that each observation in the dataset is independent. As each marketing promotion (i.e., row) is independent from one another, the independence assumption is not violated.

1.4.5 Model assumption: Normality

Create the following plots to check the **normality assumption**:

- **Plot 1:** Histogram of the residuals
- **Plot 2:** Q-Q plot of the residuals

```
[ ]: # Calculate the residuals.  
  
### YOUR CODE HERE ###  
  
# Create a histogram with the residuals.  
  
### YOUR CODE HERE ###
```

```
# Create a Q-Q plot of the residuals.
```

```
### YOUR CODE HERE ###
```

Hint 1

Access the residuals from the fit model object.

Hint 2

Use `model.resid` to get the residuals from a fit model called `model`.

Hint 3

For the histogram, pass the residuals as the first argument in the `seaborn histplot()` function.

For the Q-Q plot, pass the residuals as the first argument in the `statsmodels qqplot()` function.

Question: Is the normality assumption met?

[Write your response here. Double-click (or enter) to edit.]

1.4.6 Model assumption: Constant variance

Check that the **constant variance assumption** is not violated by creating a scatterplot with the fitted values and residuals. Add a line at $y = 0$ to visualize the variance of residuals above and below $y = 0$.

```
[ ]: # Create a scatterplot with the fitted values from the model and the residuals.
```

```
### YOUR CODE HERE ###
```

```
# Add a line at y = 0 to visualize the variance of residuals above and below 0.
```

```
### YOUR CODE HERE ###
```

Hint 1

Access the fitted values from the model object fit earlier.

Hint 2

Use `model.fittedvalues` to get the fitted values from a fit model called `model`.

Hint 3

Call the `scatterplot()` function from the `seaborn` library and pass in the fitted values and residuals.

Add a line to a figure using the `axline()` function.

Question: Is the constant variance assumption met?

[Write your response here. Double-click (or enter) to edit.]

1.4.7 Model assumption: No multicollinearity

The **no multicollinearity assumption** states that no two independent variables (X_i and X_j) can be highly correlated with each other.

Two common ways to check for multicollinearity are to:

- Create scatterplots to show the relationship between pairs of independent variables
- Use the variance inflation factor to detect multicollinearity

Use one of these two methods to check your model's no multicollinearity assumption.

```
[ ]: # Create a pairplot of the data.
```

```
### YOUR CODE HERE ###
```

```
[ ]: # Calculate the variance inflation factor (optional).
```

```
### YOUR CODE HERE ###
```

Hint 1

Confirm that you previously created plots that could check the no multicollinearity assumption.

Hint 2

The `pairplot()` function applied earlier to `data` plots the relationship between all continuous variables (e.g., between `Radio` and `Social Media`).

Hint 3

The `statsmodels` library has a function to calculate the variance inflation factor called `variance_inflation_factor()`.

When using this function, subset the data to only include the continuous independent variables (e.g., `Radio` and `Social Media`). Refer to external tutorials on how to apply the variance inflation factor function mentioned previously.

Question 8: Is the no multicollinearity assumption met?

1.5 Step 4: Results and evaluation

1.5.1 Display the OLS regression results

If the model assumptions are met, you can interpret the model results accurately.

First, display the OLS regression results.

```
[ ]: # Display the model results summary.
```



```
### YOUR CODE HERE ###
```

Question: What is your interpretation of the model's R-squared?

[Write your response here. Double-click (or enter) to edit.]

1.5.2 Interpret model coefficients

With the model fit evaluated, you can look at the coefficient estimates and the uncertainty of these estimates.

Again, display the OLS regression results.

```
[ ]: # Display the model results summary.  
  
### YOUR CODE HERE ###
```

Question: What are the model coefficients?

[Write your response here. Double-click (or enter) to edit.]

Question: How would you write the relationship between **Sales** and the independent variables as a linear equation?

[Write your response here. Double-click (or enter) to edit.]

Question: What is your interpretation of the coefficient estimates? Are the coefficients statistically significant?

[Write your response here. Double-click (or enter) to edit.]

Question: Why is it important to interpret the beta coefficients?

[Write your response here. Double-click (or enter) to edit.]

Question: What are you interested in exploring based on your model?

[Write your response here. Double-click (or enter) to edit.]

Question: Do you think your model could be improved? Why or why not? How?

[Write your response here. Double-click (or enter) to edit.]

1.6 Conclusion

What are the key takeaways from this lab?

[Write your response here. Double-click (or enter) to edit.]

What results can be presented from this lab?

[Write your response here. Double-click (or enter) to edit.]

How would you frame your findings to external stakeholders?

[Write your response here. Double-click (or enter) to edit.]

References Saragih, H.S. (2020). *Dummy Marketing and Sales Data*.

Congratulations! You've completed this lab. However, you may not notice a green check mark next to this item on Coursera's platform. Please continue your progress regardless of the check mark. Just click on the "save" icon at the top of this notebook to ensure your work has been logged.