# Predicting restaurant tips using predictive analytics on Excel.

## In this Project:

1. I selected all the data and filtered if there is any blank cell or not

2. Checked and removed duplicates

3. Found Independent and dependent variables

   Independent Variables = Sex, Smoker, Day, Time, size, total bill

   Depended variable = Tips

4. Used multiple regression analysis to obtain data model as shown in screenshots and excel files

5. Converted categorical values into numerical using IF condition

6. Then using predictive analysis multiple regression formula, found the predictive tips alongside   actual ones

7. Then found errors in next column

8. Finally calculated RMSE


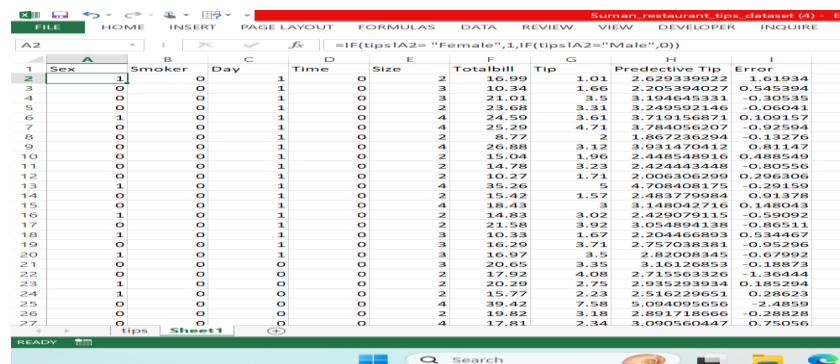Encode the categorical variables to numeric values using IF conditions.

### Problem 1:

Differentiate Gender by giving them the different values

### Process:

Converting categorical variables to numeric values cab be done by adding dummy values to each column.
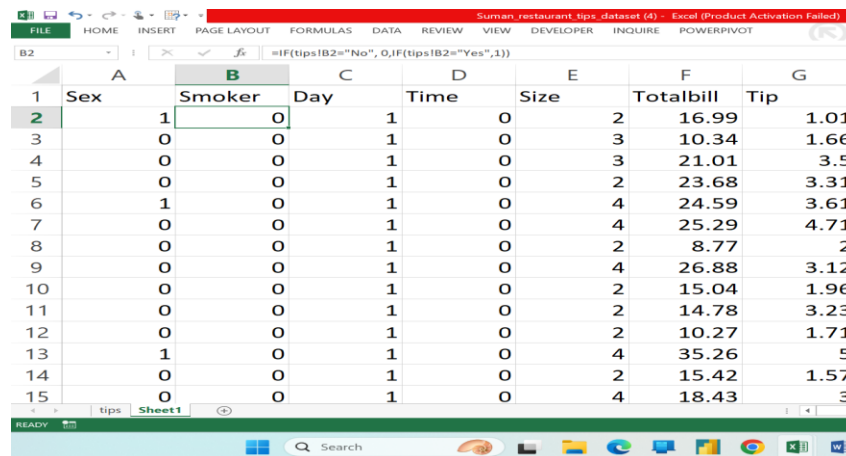

### Output:

## Problem 2:

Know whether the smoker is convertible or not? If convertible change them.

## Process:

They are convertible because there is no varchar data type in the values of the column smoker.
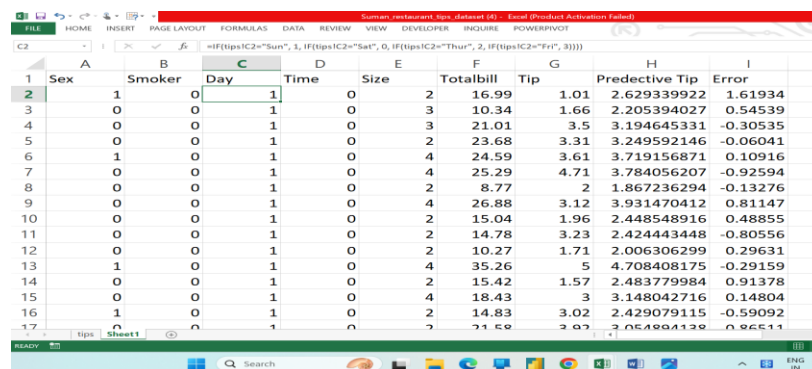
## Output:



## Problem 3:

Analyses the days in the week based that are crowded.

## Process:

According to the data the 3 week days are Thursday, Friday, Sunday and Saturday.

The columns are created and also variables are created based on the DAY column.

## Output:
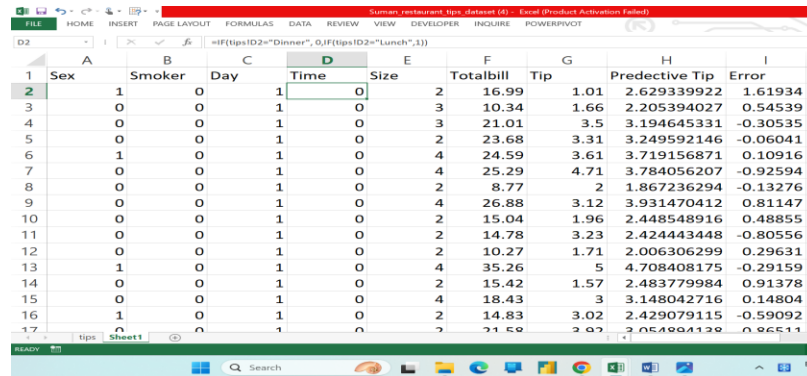
## Problem 4:

Calculate the dinner column by converting categorical variables in to numerical to know the bill size

### Process:

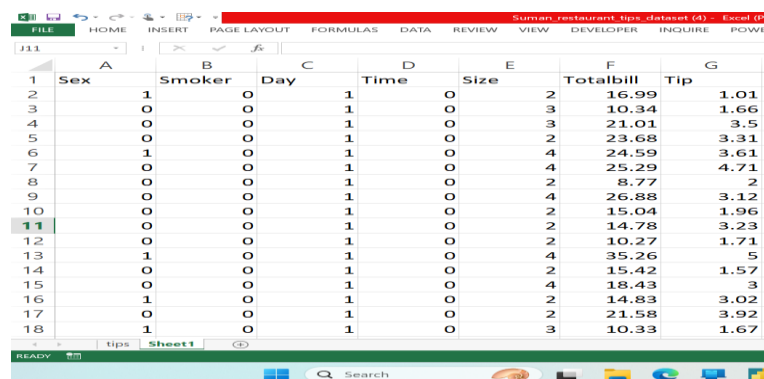As predicted tips was calculated by dinner size, the column was created.

### Output:



## Problem 5:

Identify an appropriate model with the dataset.

### Process:

To identify the right model the complete dataset should be cleaned and being formatted.

### Output:



As we had completed the process of data cleaning.

The desired model would be REGRESSION according to the data obtained.
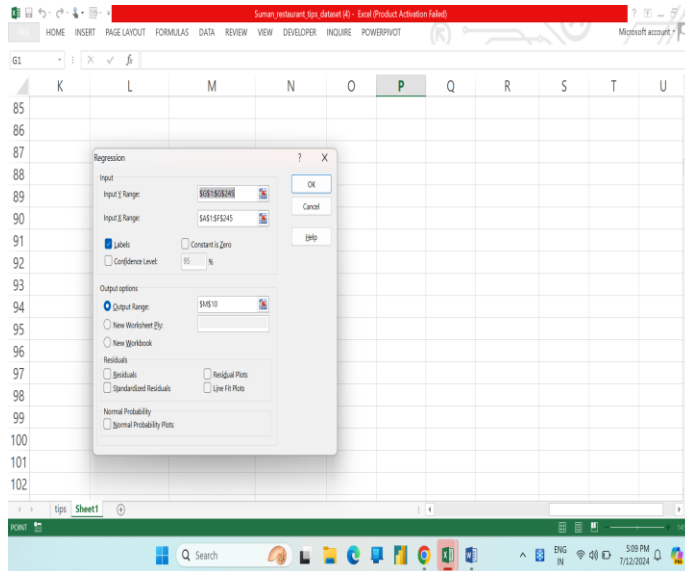
## Problem 6:

Build an appropriate model with the new table.

**Process:**

We have created the Regression Model

**Output:**



After applying the Regression Model.

**Output:**

SUMMARY OUTPUT

| Regression Statistics | |
|---|---|
| Multiple R | 0.684009729 |
| R Square | 0.467869309 |
| Adjusted R Square | 0.463453286 |
| Standard Error | 1.013505967 |
| Observations | 244 |

ANOVA

| | df | SS | MS | F | Significance F |
|---|---|---|---|---|---|
| Regression | 2 | 217.6586 | 108.8293 | 105.9481 | 9.6651E-34 |
| Residual | 241 | 247.5538 | 1.027194 | | |
| Total | 243 | 465.2125 | | | |

| | Coefficients | andard Err | t Stat | P-value | Lower 95% | Upper 95% | ower 95.0% | pper 95.0% |
|---|---|---|---|---|---|---|---|---|
| Intercept | 0.668944741 | 0.193609 | 3.455127 | 0.00065 | 0.2875622 | 1.050327 | 0.287562 | 1.050327 |
| Size | 0.192597794 | 0.085315 | 2.257502 | 0.024872 | 0.02454038 | 0.360655 | 0.02454 | 0.360655 |
| Totalbill | 0.092713337 | 0.009115 | 10.17187 | 1.88E-20 | 0.07475872 | 0.110668 | 0.074759 | 0.110668 |

## Problem 7:

Calculate the Predicted Tips.

Output:



## Problem 8:

Find the Error in the tips according to the Predicted tips.



Find the Sum Square:

# Find the COUNT:



Formula bar: `=COUNT(I2:I245)`

| | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Sum Square | 247.553837 |
| COUNT | 244 |

# Find the RMSE:



Formula bar: `=L3/L4`

| | |
|---|---|
| Sum Square | 247.553837 |
| COUNT | 244 |
| RMSE | 1.014564906 |