**Course Project**
**DeVry University**
**College of Engineering and Information Sciences**
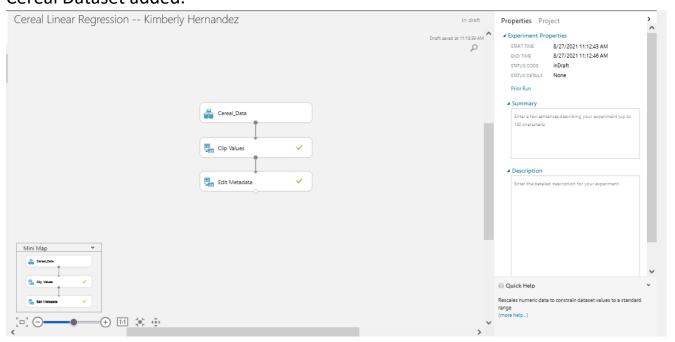
**Course Number: CEIS312**
Module 8 Report
Kimberly Hernandez

# Introduction to the problem

A cereal data set from Kaggle.com is given in order to determine predictions of consumer report ratings on different cereal. This dataset will be added to the Azure Machine learning environment and linear regression will be performed on it. Several steps of the machine learning process will be taken to create a model accurate enough for future predictions of ratings on cereal. Some of the project's focus will be on data preparation and the iteration process. The model must reach at least a 0.7 Coefficient of determination or $R^2$ and to consider error and bias of the model the Relative Squared Error will be taken into consideration as well.

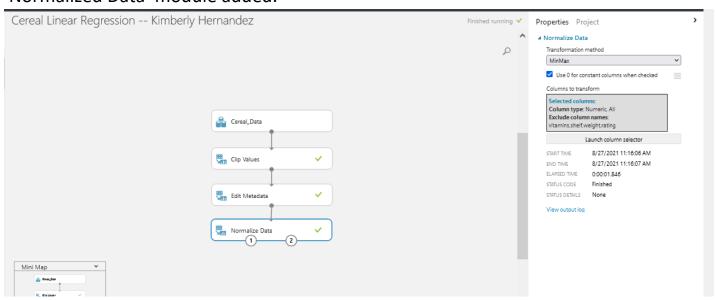# Uploading dataset
Cereal Dataset added.

## Data preparation (normalization)

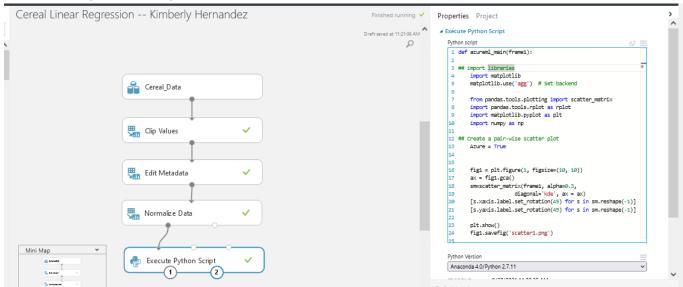'Edit Metadata' module added to make category fields 'Categorical'.
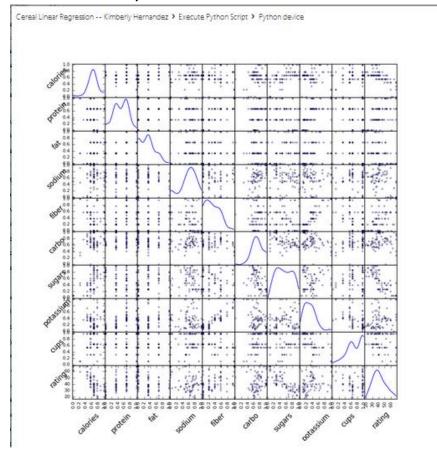


'Normalized Data' module added.

## Data Visualization (python script)

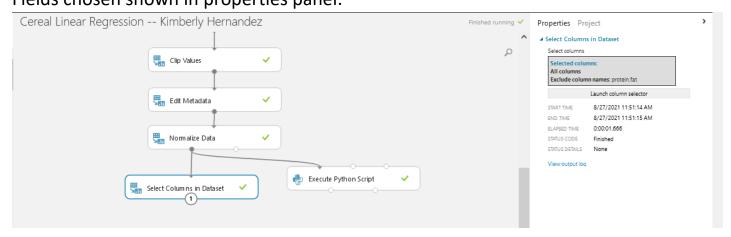'Execute Python Script' module added.
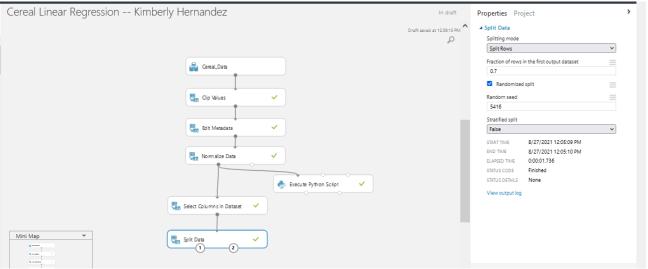


## Code created Python charts.

## Selecting features

Selected columns in dataset module added.

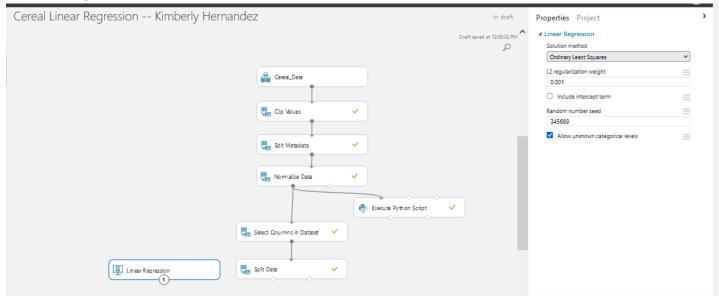Fields chosen shown in properties panel.



## Splitting data
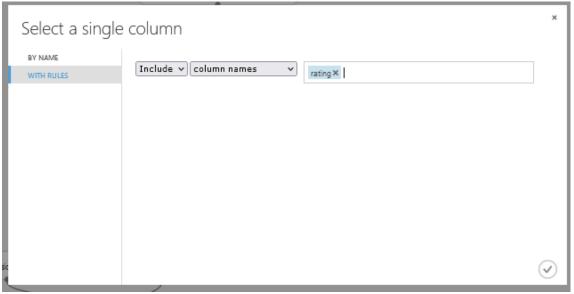
Split module added.

## Linear regression model

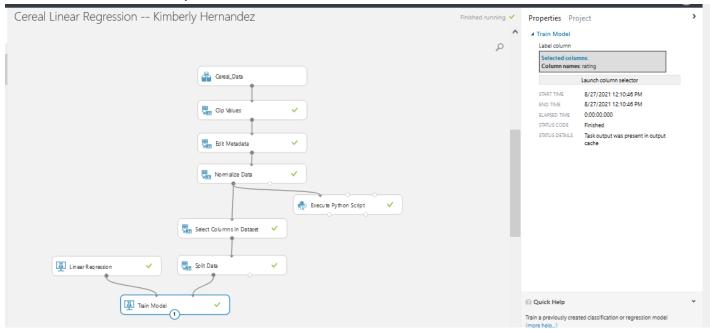Linear Regression Module added.



## Training the model

'Rating' feature being selected under the 'Train Model' module.
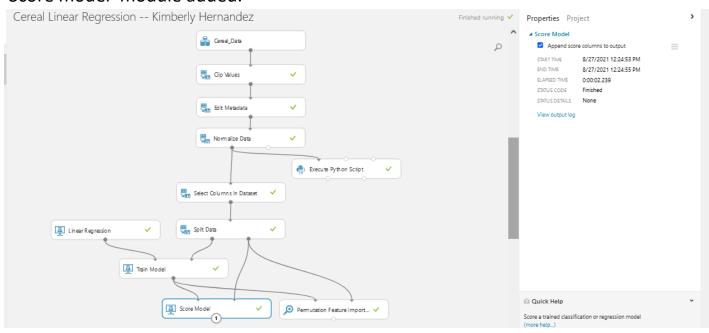
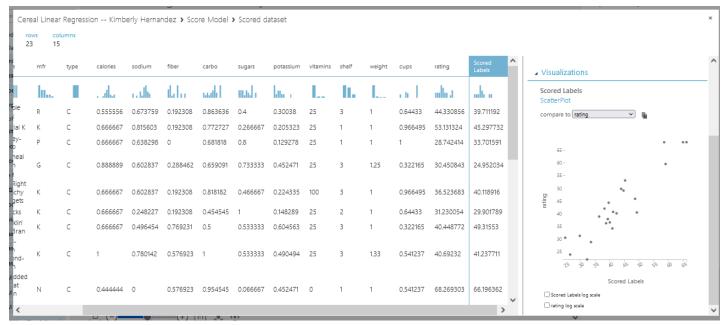'Train model' module present.



## Scoring the model – show scored labels
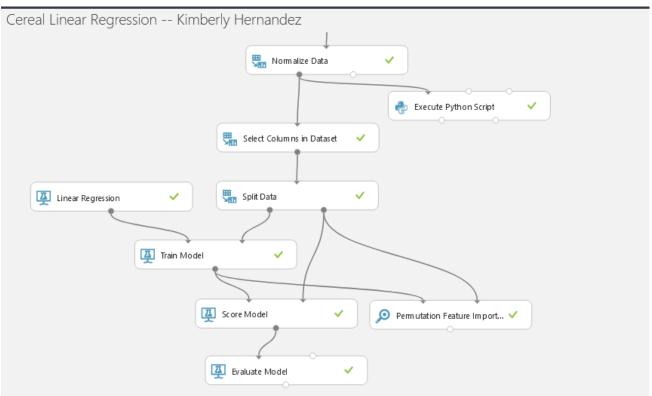'Score model' module added.

Scored labels visualization.



## Evaluating the model
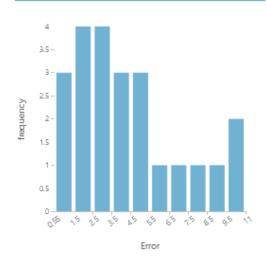'Evaluate Model' module added.

# All metrics and histograms.

### ◢ Metrics

| | |
|---|---|
| Mean Absolute Error | 4.287005 |
| Root Mean Squared Error | 5.106688 |
| Relative Absolute Error | 0.412733 |
| Relative Squared Error | 0.152489 |
| Coefficient of Determination | 0.847511 |

### ◢ Error Histogram



**Explanation:** This model excluded protein and fat because the scatter plot values from the python script showed scattered values in these two features. The current $R^2$ value is 0.847511 which is relatively close to the value of 1.0 which is considered a perfect model. Although there could be some improvement with the $R^2$ and Relative Squared Error, so an iteration process must take place.

## Iteration process

**Explanation:** Since we want to consider error and bias of model some features must be put back in and some will be removed. Doing this will also take into consideration of the overall accuracy of model.

- ### Why were certain features excluded?

**Explanation:** At this point it has been determined through 'Permutation Feature Importance' that the model needs the name, type, shelf, weight, and cups to be excluded from the model. In addition to this process, protein and fat were added back in to evaluate

where these values could bring the model's $R^2$ and Relative Squared Error. (Image below)



- **New evaluation when those features were excluded.**

Cereal Linear Regression -- Kimberly Hernandez > Evaluate Model > Evaluation results

▲ Metrics

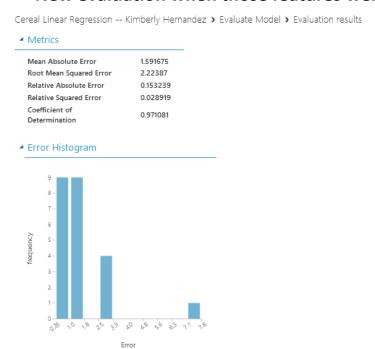| | |
|---|---|
| Mean Absolute Error | 1.591675 |
| Root Mean Squared Error | 2.22387 |
| Relative Absolute Error | 0.153239 |
| Relative Squared Error | 0.028919 |
| Coefficient of Determination | 0.971081 |

▲ Error Histogram



**Explanation:** The model at this point includes protein & fat and excludes the ones listed above. Now the model has a $R^2$ of 0.971081 which is significantly closer to the value of 1.0 which would be considered a perfect model. There is also improvement of the Relative Squared Error which is at 0.028919 which is significantly closer to the value of 0.0 which would be considered a perfect model that has all model errors of 0. It seems that making these alterations to the model increasing its overall accuracy, bias, and error.

- **What features are most influential on the rating?**

  **Explanation:** When the model incorporated features such as Fiber, Sugars, mfr, carbo, vitamins, potassium, calories, and sodium the $R^2$ was at 0.847511. After adding protein and fat the model's $R^2$ significantly increased to 0.971081. It is safe to say after looking at the 'Permutation Feature Importance' module that the most influential features are sugars, calories, protein, carbo, and sodium.

## Conclusion

The model reached an $R^2$ of 0.971081 and a Relative Squared Error(RSE) of 0.028919. It is determined that the model can predict cereal ratings due to the $R^2$ and RSE of the model reaching close to perfect parameters. It was important that the model undergo areas of importance such as the data preparation and iteration process. Without these two areas the model would not have reached an accuracy closer to 1.0.

## Challenges Faced

An error occurred during normalization of model. Error was fixed when 'Vitamins', 'Shelf', and 'Weight' columns was excluded from the 'Normalized data' module.

## Career skills obtained

*Software Skills* – These skills were obtained by using the Azure Machine Learning Studio.
*Machine Learning Algorithms* – Creating the model using linear regression.
*Data Modeling and Evaluation* – Finding patterns of the data to see which features make the model as accurate as possible. This also incorporates going through the evaluation and iteration process.
*Problem Solving skills* – Figuring out how to resolve an error during the process.
*Communication skills* – Creating a report on the dataset for others to understand the process of model development.