# CEIS312 Introduction to Artificial Intelligence & Machine Learning

Final Course Project

Developed by James Garlie

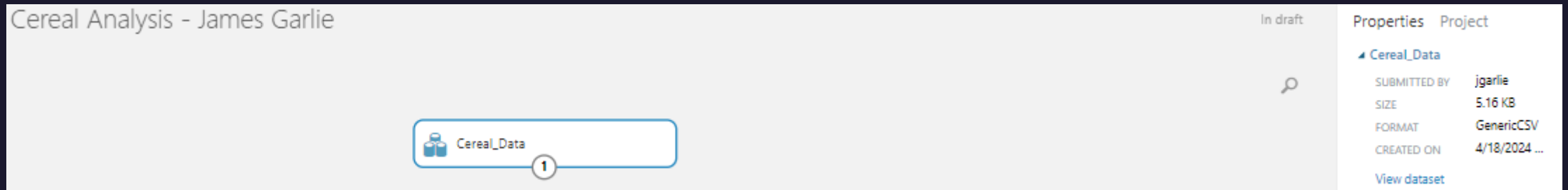DeVry University

April 2024

# Introduction

Artificial intelligence and machine learning are frontiers in the technology field. These areas are often used to address common problems that require difficult tools or skills. AI and ML professionals work with SQL, R, Python, and other tools specific to data science. Different algorithms are used to solve problems and choosing the correct algorithm can be challenging. This project will use Azure Machine Learning, which is a cloud-based service from Microsoft. Azure ML allows you to create and run experiments based on datasets and integrate custom code in SQL, R, or Python.

The presentation concludes with the Conclusion, Challenges, and Career Skills obtained.

# Uploaded Dataset

The first image shows that an experiment has been created titled Cereal Analysis and that a .csv file, titled Cereal Data has been uploaded from my local device and added to the experiment. Note the GenericCSV format and file size shown on the right of the image.
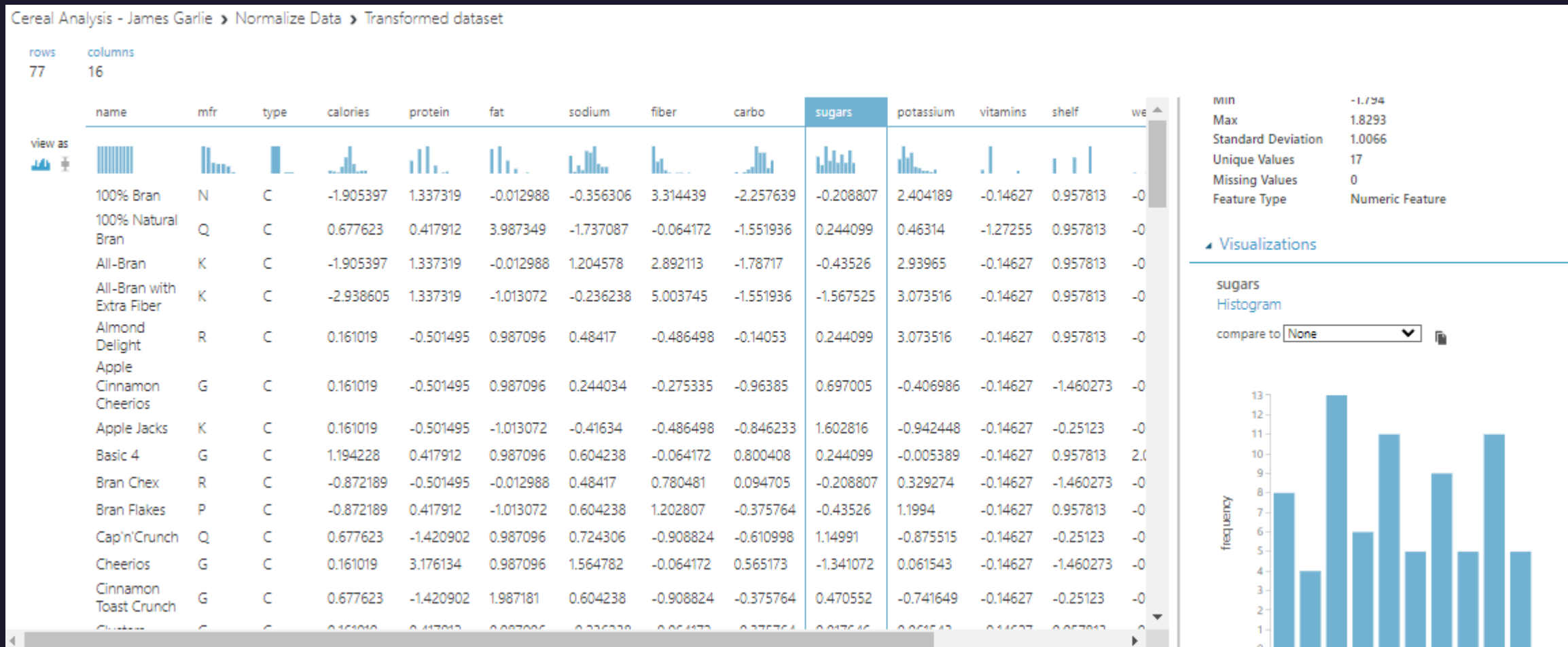


This second image shows the initial visualization of the Cereal Dataset with 77 rows & 16 columns.

# Data Preparation & Normalization

This slide shows filtering the data, formatting, and normalization.

# Data Visualization

This slide shows the Python script used for the visualization of the data.

```python
def azureml_main(frame1):

## import libraries
    import matplotlib
    matplotlib.use('agg')  # Set backend

    from pandas.tools.plotting import scatter_matrix
    import pandas.tools.rplot as rplot
    import matplotlib.pyplot as plt
    import numpy as np

## Create a pair-wise scatter plot
    Azure = True


    fig1 = plt.figure(1, figsize=(10, 10))
    ax = fig1.gca()
    sm=scatter_matrix(frame1, alpha=0.3,
                diagonal='kde', ax = ax)
    [s.xaxis.label.set_rotation(45) for s in sm.reshape(-1)]
    [s.yaxis.label.set_rotation(45) for s in sm.reshape(-1)]

    plt.show()
    fig1.savefig('scatter1.png')
```
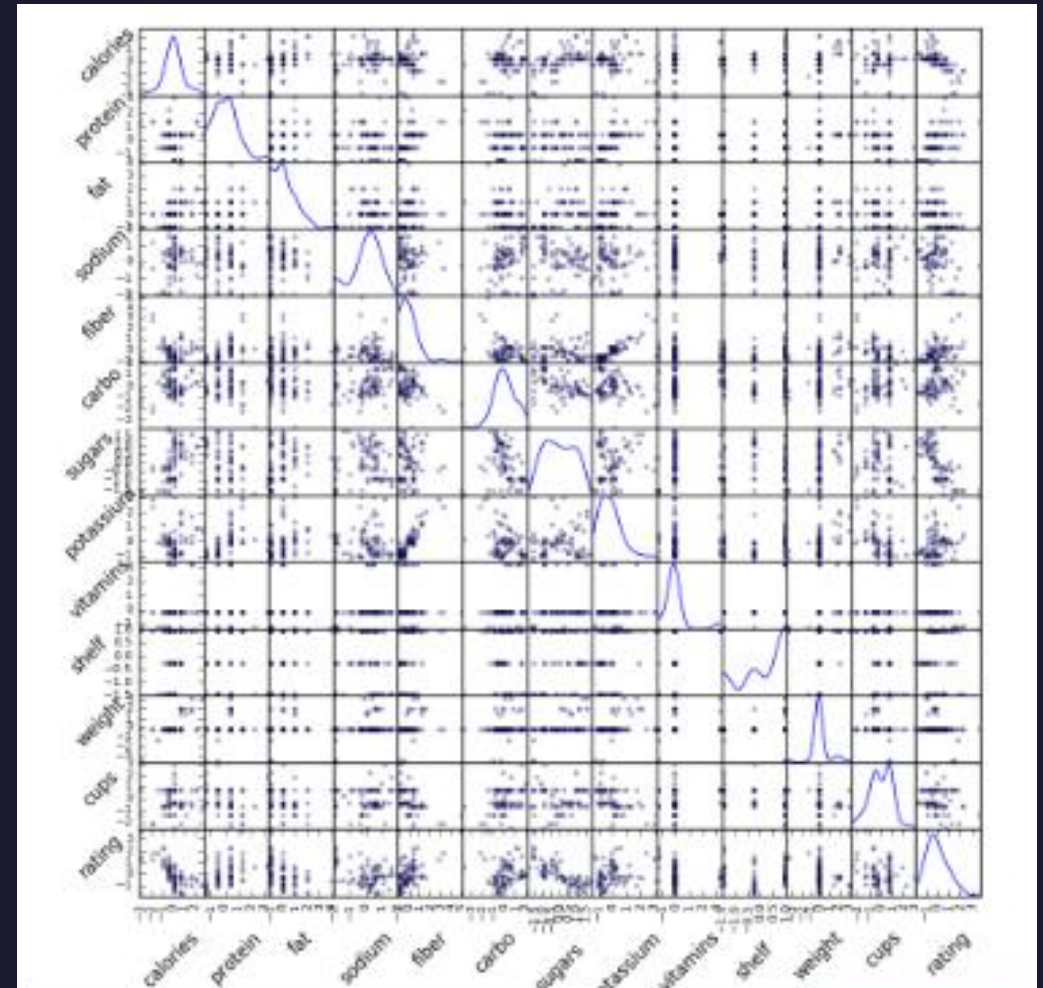
# Selecting Features

This slide shows that I added "Select Columns in Dataset", then selected calories, protein, fiber, and vitamins. I then saved and ran the experiment. The visualization shows the 4 features I selected. Also note we are down to 4 columns verses the original 16.
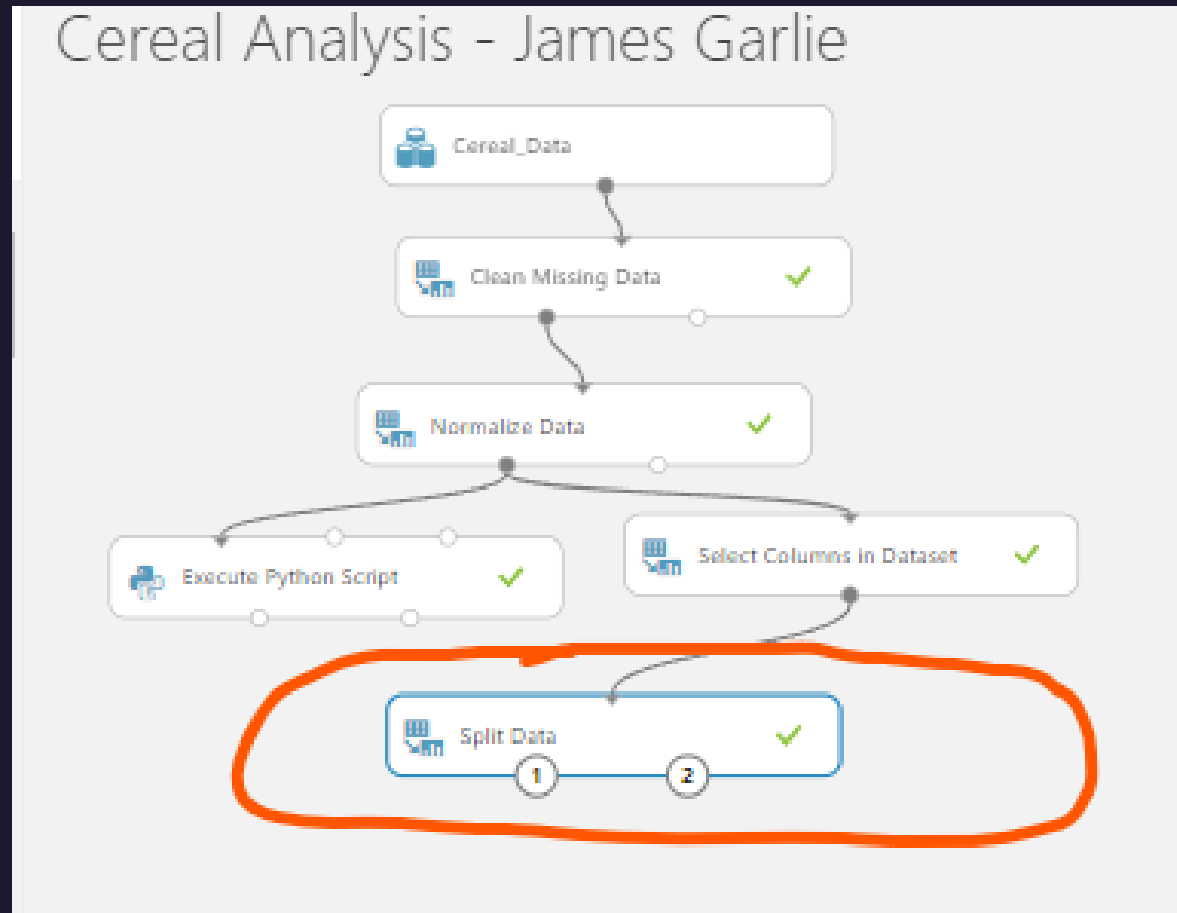
# Splitting Data

This slide shows the addition of the Split Data module and the results of dataset1. Notice we now show 46 rows with 4 columns.
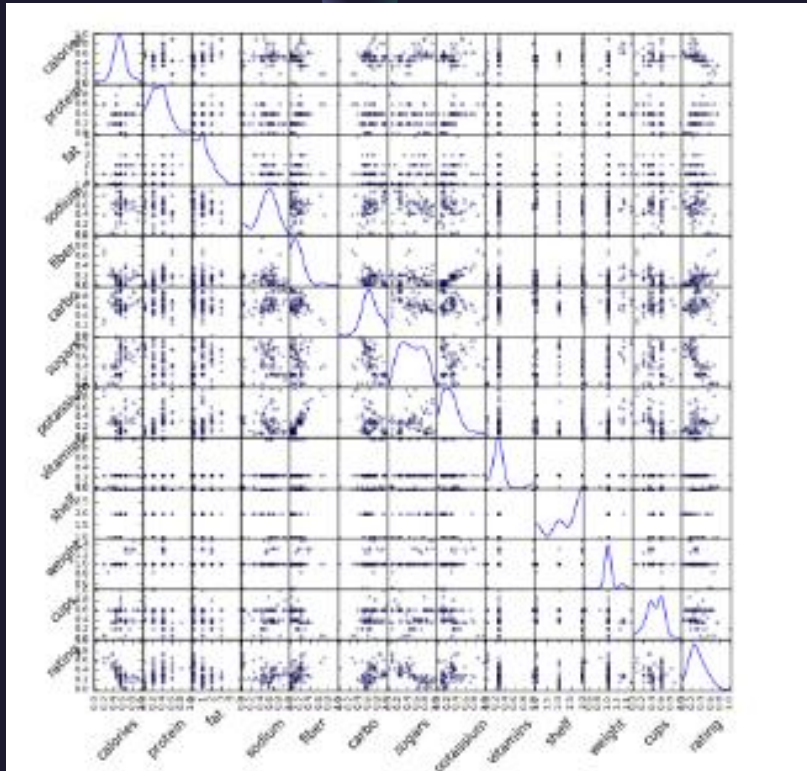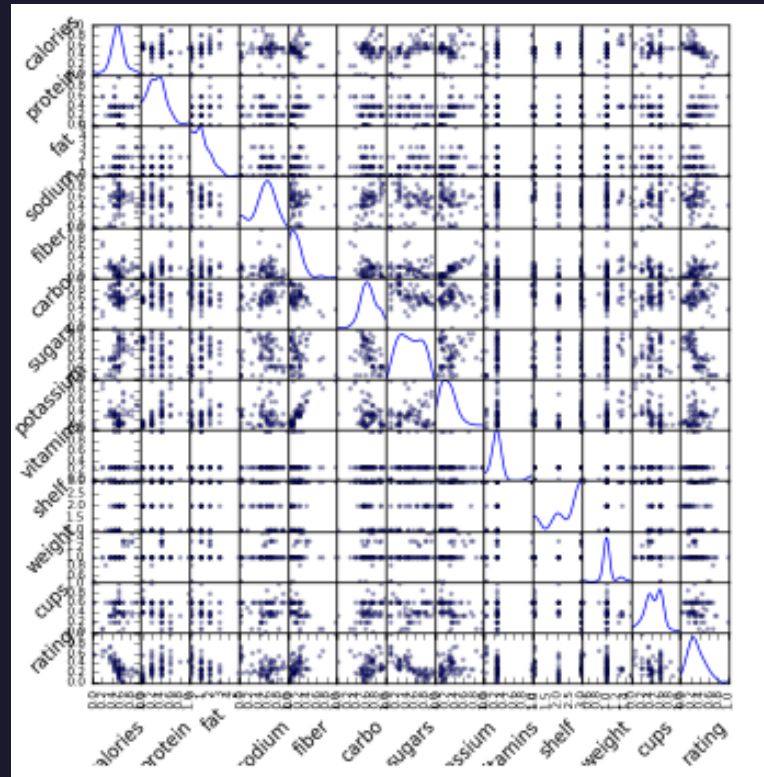
# Iteration Process Threshold

This slide shows three different iterations of the code. The first is the original (10, 10), the second is (7.5, 7.5), and the third is (5, 5). When enlarged, I like the second or middle the best.
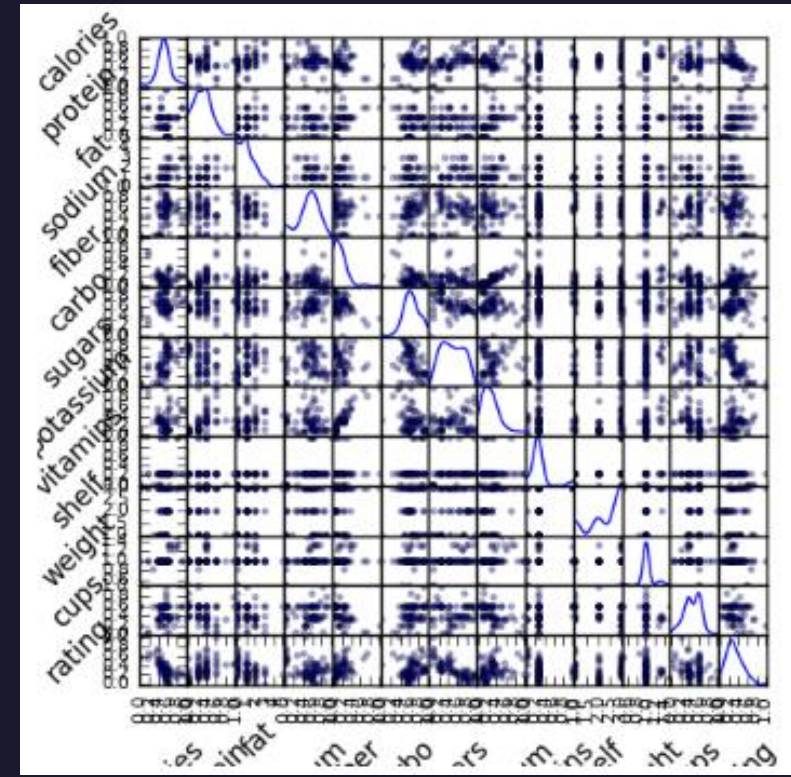
fig1 = plt.figure(1, figsize=(10, 10))

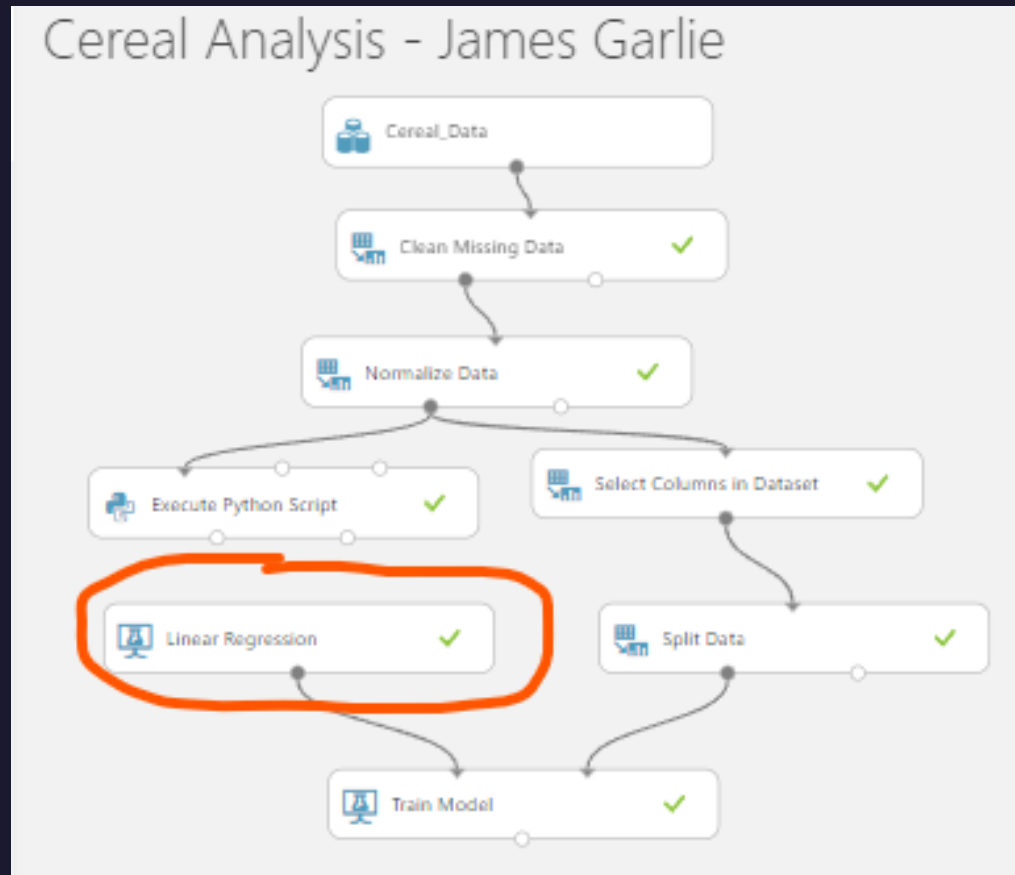fig1 = plt.figure(1, figsize=(7.5, 7.5))

fig1 = plt.figure(1, figsize=(5, 5))

# Linear Regression

This slide shows the addition of the Linear Regression module and the results of the Untrained model.
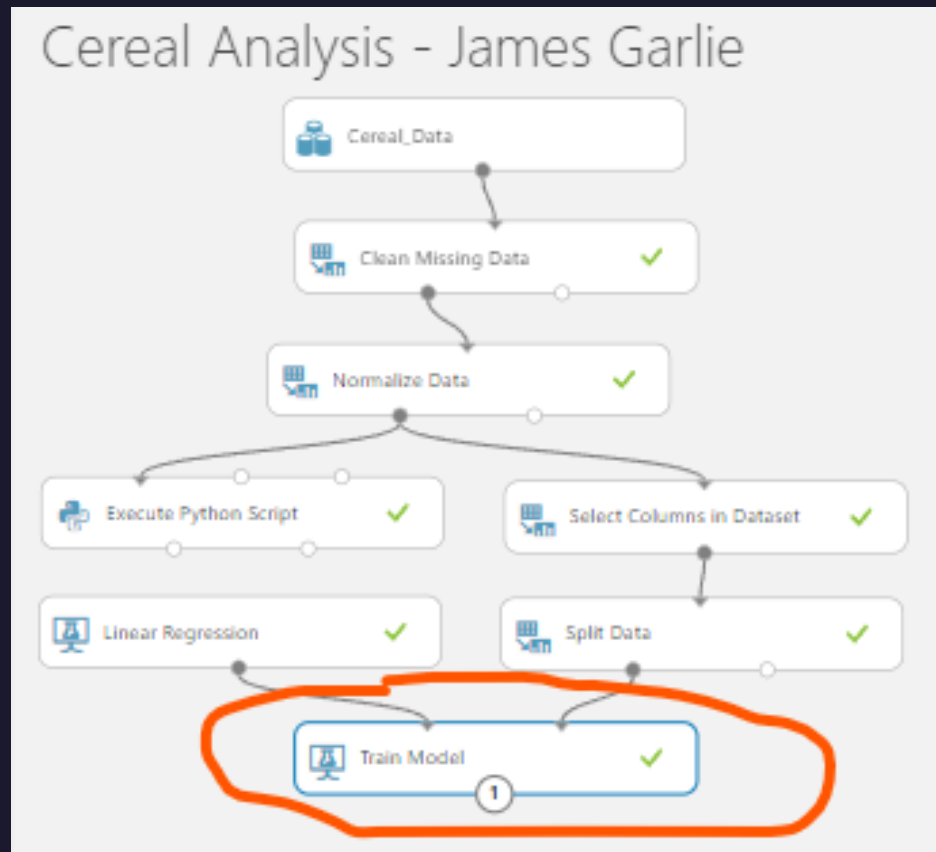


Cereal Analysis - James Garlie

- Cereal_Data
- Clean Missing Data ✓
- Normalize Data ✓
- Execute Python Script ✓
- Select Columns in Dataset ✓
- Linear Regression ✓
- Split Data ✓
- Train Model ✓



Cereal Analysis - James Garlie ❯ Linear Regression ❯ Untrained model

**Batch Linear Regressor**

### Settings

| Setting | Value |
| --- | --- |
| Bias | False |
| Regularization | 0.001 |
| Allow Unknown Levels | True |
| Random Number Seed | |

# Training Model

This slide shows the addition of the Train Model module where I deleted the calories feature, and the results of the Trained Model showing no calories.

# Scoring Model

Here I added the Score Model with results showing 31 rows and 5 columns. Notice calories has been included the new feature showing Scored Lables.
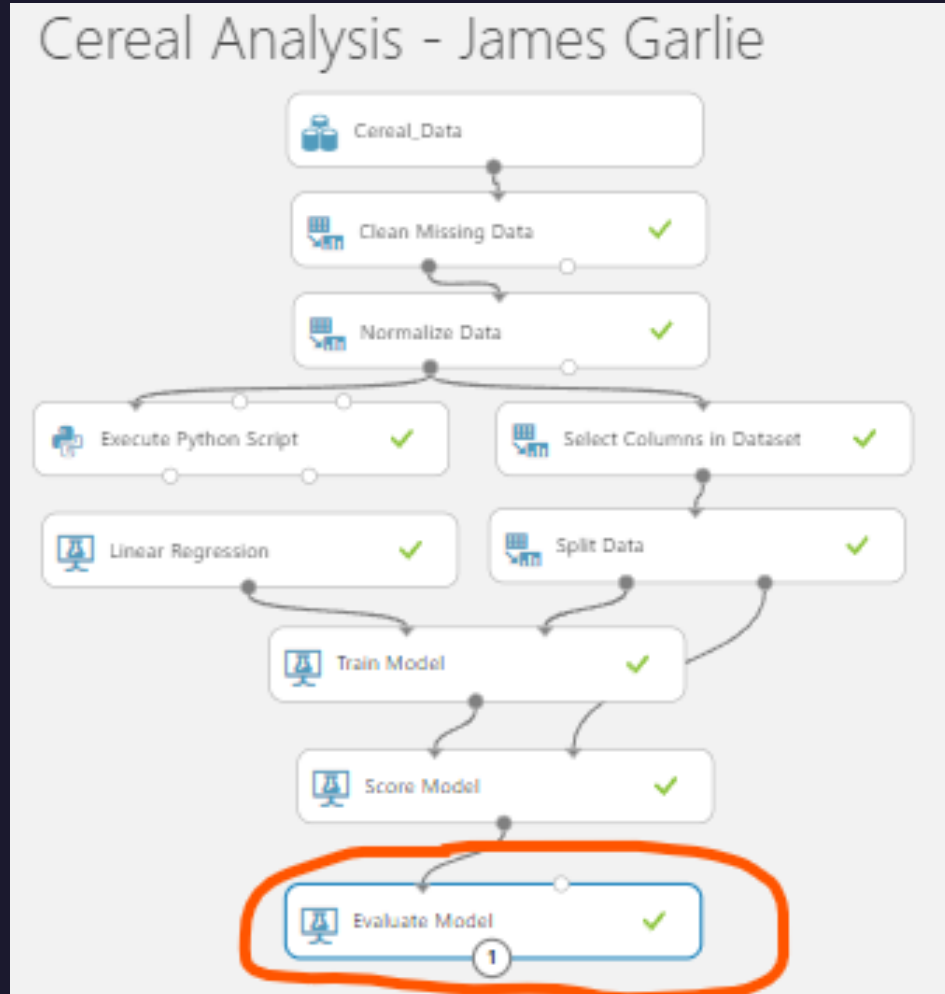
# Evaluating the Model
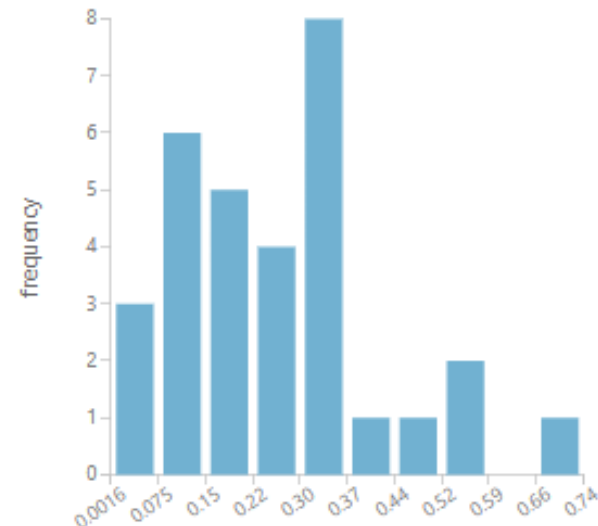
This slide shows the addition of the Evaluate Model module and the results with the Coefficient of Determination of -1.08025.

# Conclusion

Learning the cleansing of the data is vital for Machine-learning.

I found the aspects of connecting modules, changing the data I was looking for, and visualizing the results to be very rewarding.

This project will be of tremendous benefit in the

# Career Skills

➢ Probability and statistics

Programming skills (Python or R)

Data skills (data processing, SQL data analysis, visualization skills)

Machine-learning algorithms

Lifelong learning

TensorFlow (neural networks)

Apache Spark

The analysis we have done using Azure Machine Learning can also be performed in some of these programming languages. For example, using some Python modules, you can receive the same Energy Efficiency Regression in Python.

# Challenges

The biggest challenge I faced
was finding the correct modules
to choose from to assemble the
data required.