# CIS5560 Term Project Tutorial

**Authors: Gregory D. Mamoyac, Tejas Agara Chandrakumar, Nitesh Kamboj, Julia Stachurska**

**Instructor:** [Jongwook Woo](#)

**Date: 05/06/2018**

# Lab Tutorial

jstachu ([jstachu@calstatela.edu](mailto:jstachu@calstatela.edu))

nkamboj([nkamboj@calstatela.edu](mailto:nkamboj@calstatela.edu))

afnu3([afnu3@calstatela.edu](mailto:afnu3@calstatela.edu))

gmamoya[(gmamoja@calstatela.edu)](mailto:gmamoja@calstatela.edu)

05/06/2018

# Azure ML and Spark ML Analysis of Stack 0verflow

For R Language Question and Answers

**Azure K Means**

## Objectives

**List what your objectives are.** In this hands-on lab, you will learn how to:

- Get data manually
- Create Spark cluster
- Train KMeans  system
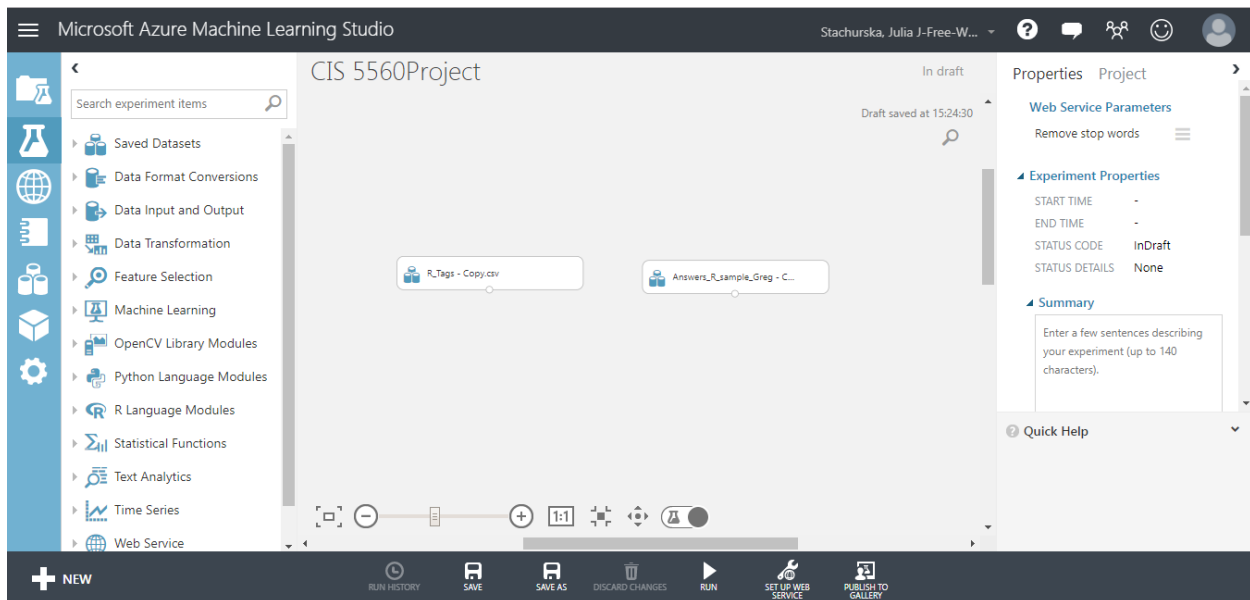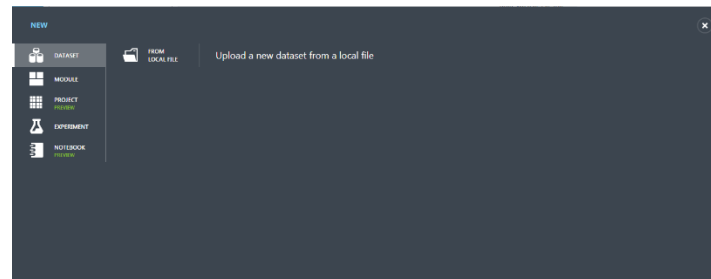- K-Means Clustrering for predicting if an answer is acceptable

- Visualization
- https://gallery.cortanaintelligence.com/Experiment/Clustering-kmeans-2

## Platform Spec

- Microsoft Azure ML
- CPU Speed: ~3.4GHz
- # of nodes: 1
- Total Memory Size: 10GB

# Step 1: Upload the Data Set from local file

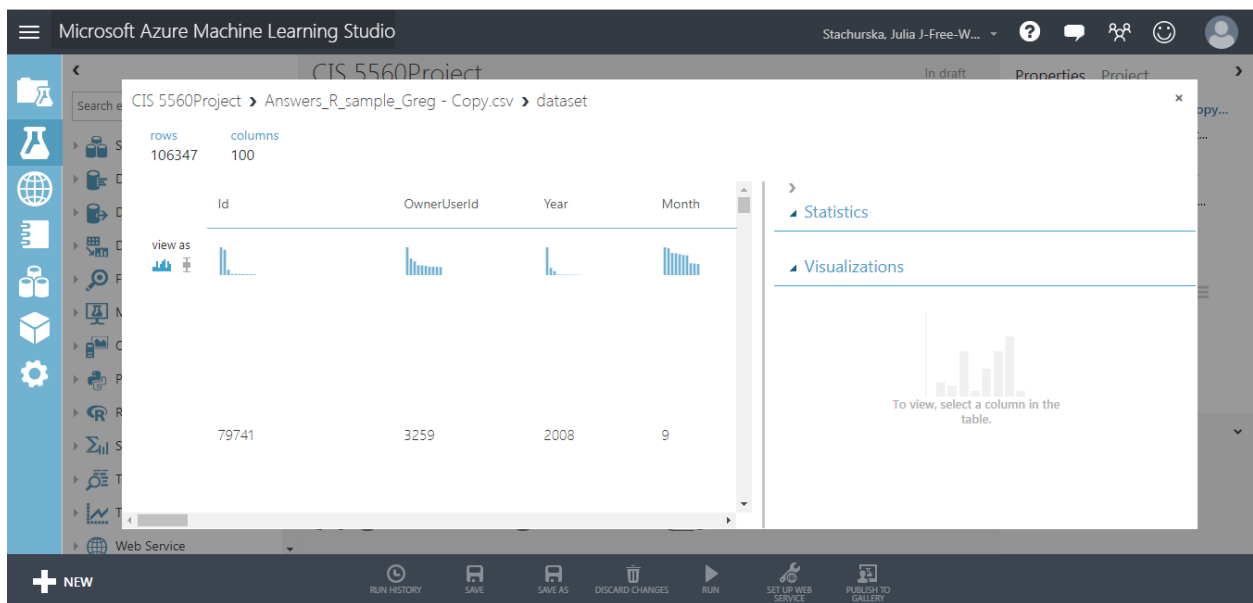This step is to upload two datasets:  R_tags.csv and Anwers_R_Sample.csv from local file.
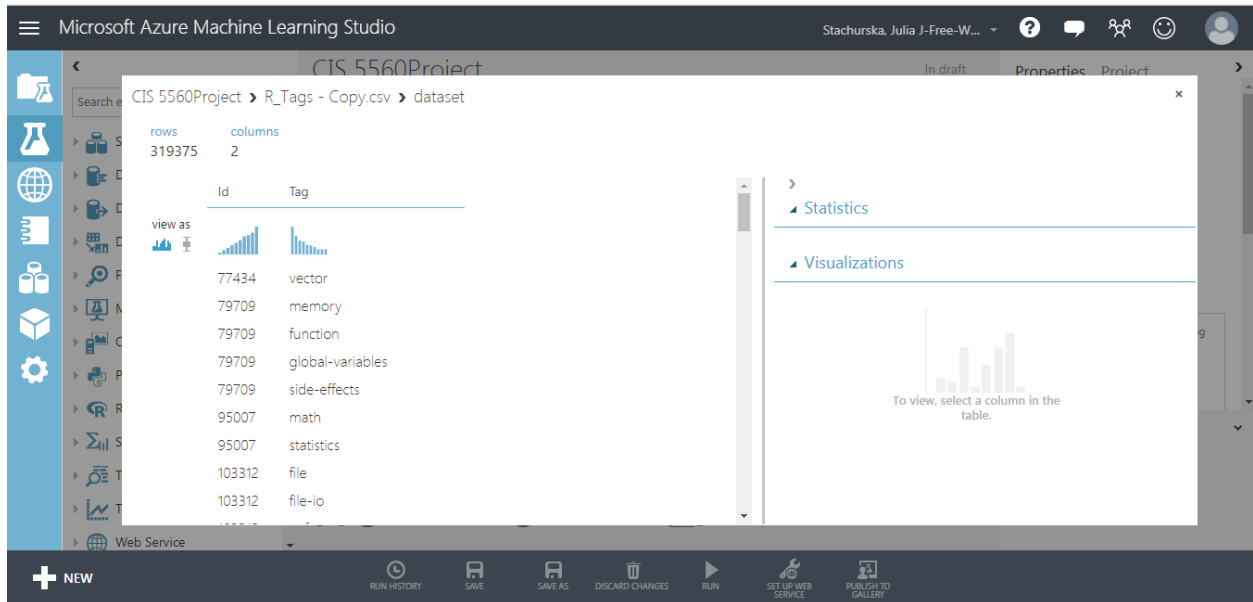




Properties:

- This is a new version of an existing dataset: Unselected
- Enter a name for the new dataset: R_Answers_LDA_Sample
- Select a type for the new dataset: Generic CSV file with a header(.csv)
- Provide an optional description: R Answers RDA.

# Step 2: Visualize the Dataset in Azure ML

This step is to verify if the data set uploaded contains all the data from the source file.

# Step 3: Join Data

This step is to join two datasets into one, using columns "Id" and "ParentId".



Properties:

- Join key columns for L: column name "Id"
- Join key columns for R: column name "ParentId"
- Match case: selected
- Join type: Inner join
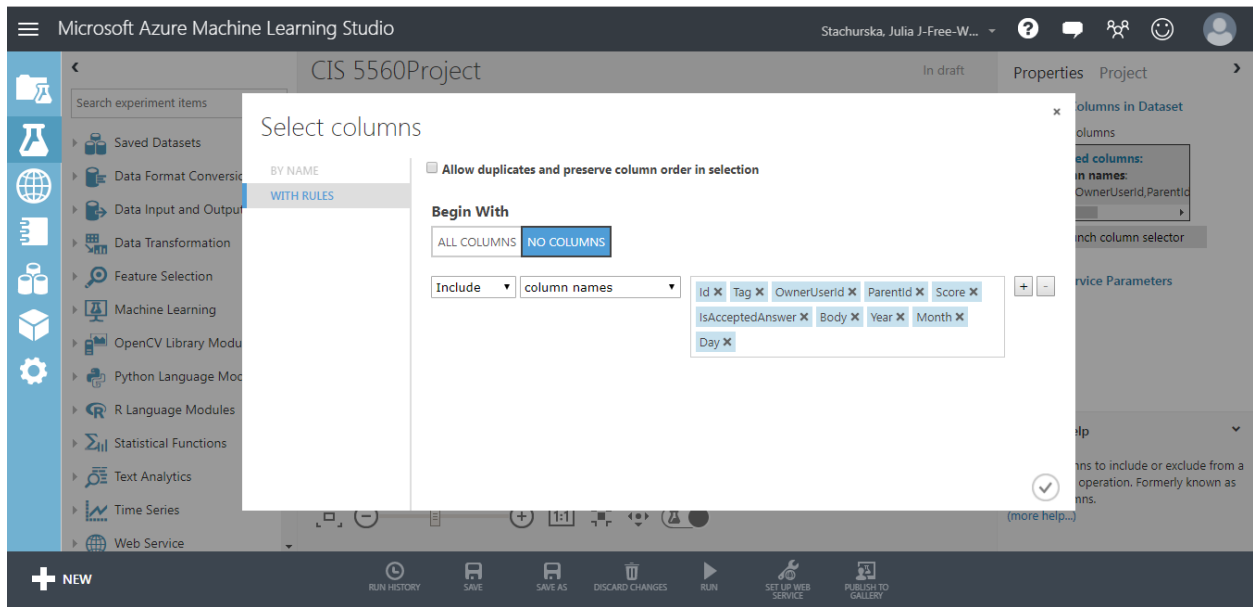- Keep right key columns in joined table: checked

# Step 4: Select columns in Dataset

This is a common user interface element in Azure ML modules to enable selecting the columns you want to use in the module.
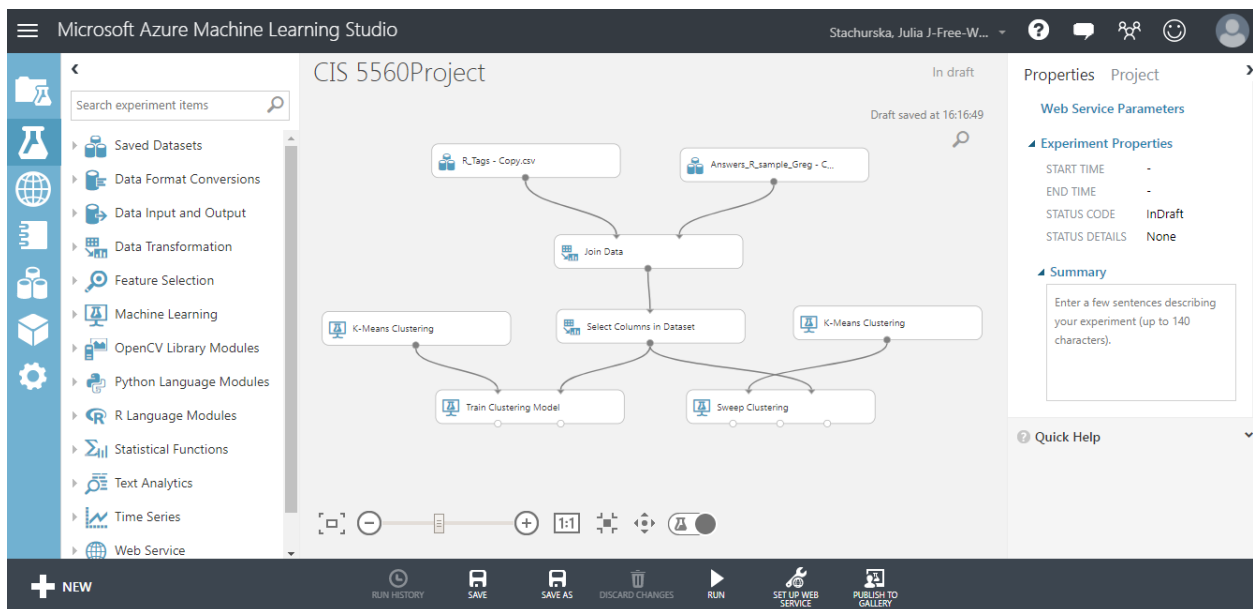


In the Select columns dialog box, select option With Rules to begin with no columns, and include:

- Id,
- Tag,
- OwnerUserId,
- ParentId,
- Score,
- IsAcceptedAnswer,
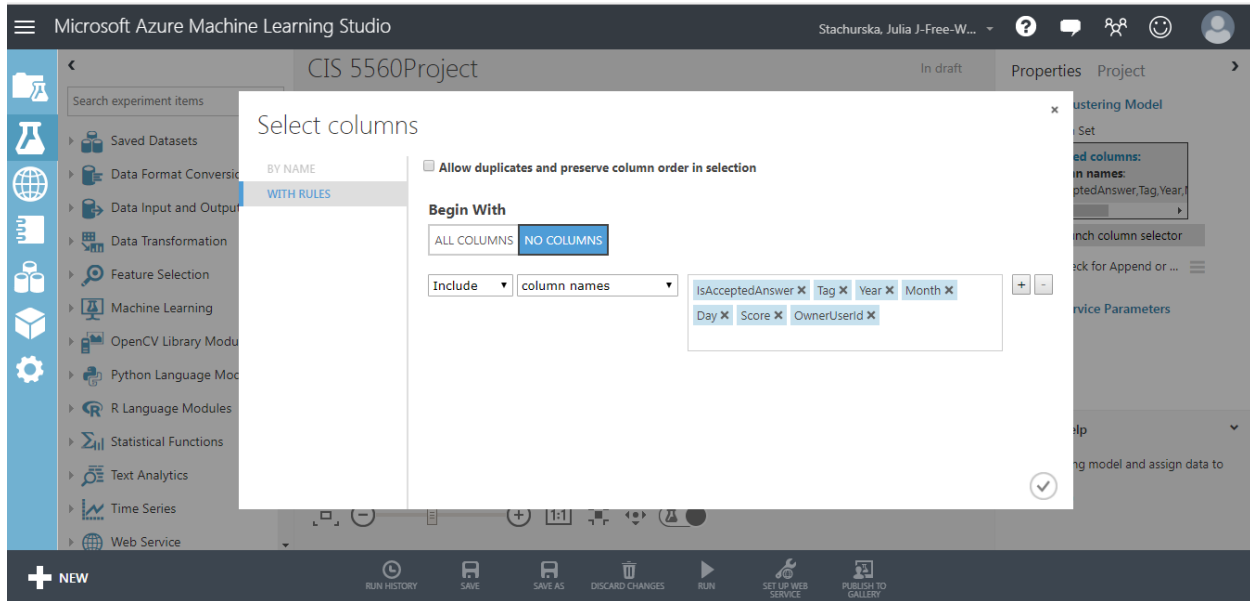- Body,
- Year,
- Month,
- Day

# Step 5: Train Clustering Model and Sweep Clustering (K-Means)

This step is to create clusters (groups in the data) with the number of groups represented by the variable *K*.
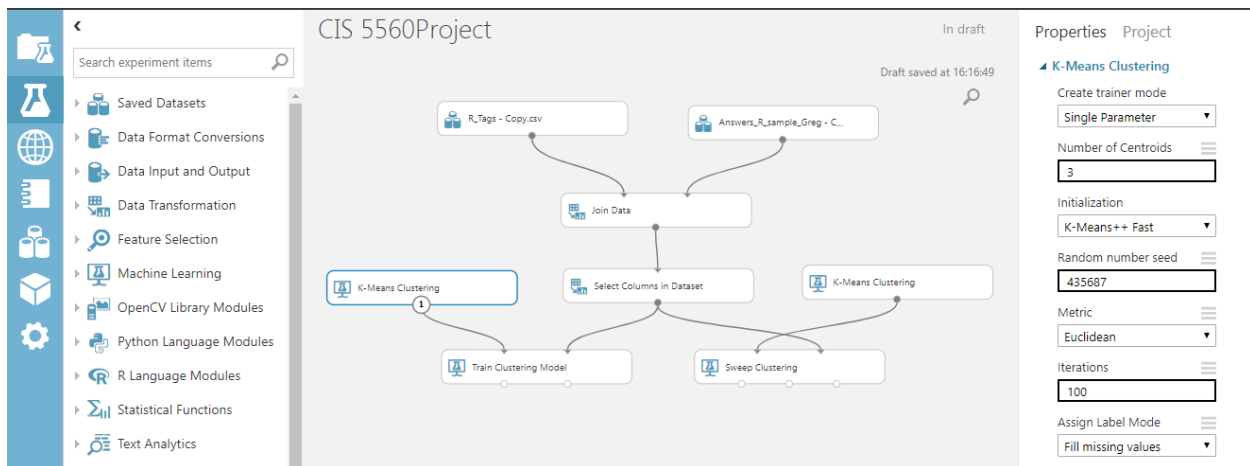
# 5A: Train Clustering Model



Properties:

- Columns used: IsAcceptedAnswer, Tag, Year, Month, Day, Score, OwnerUSerID



K-Means Clustering Properties:

- Create trainer mode: Single Parameter
- Number of centroids: 3
- Initialization: K-Means++ Fast
- Random number speed: 435687
- Metric: Euclidean
- Iterations: 100
- Assign Label Model: Fill Missing Values

## 5B: Sweep Clustering



Properties:

- Metric for measuring clustering result: Simplified Silhouette
- Specify parameter sweeping mode: Random Sweep
- Maximum number of runs on random sweep: 5
- Random seed: 0
- Column Set: Tag Month Year Day OwnerUserId Score IsAcceptedAnswer
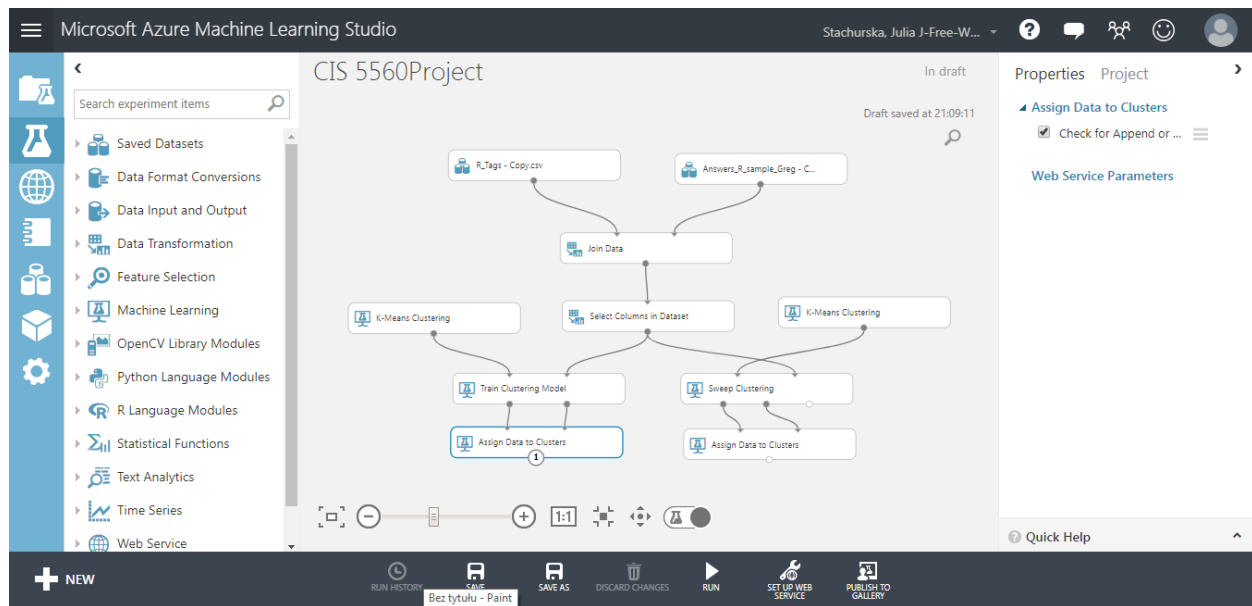- Check for Append or Uncheck for Result Only: Checked



K-Means Clustering Properties:

- Create trainer mode: Parameter Range
- Range for Number of Centroids: 2, 3, 4, 5 (Use Range Builder: Not Checked)
- Initialization for sweep: K-Means++

- Random number seed:435687
- Number of seeds to sweep:1
- Metric: Euclidean
- Iterations: 100
- Assign Label Mode: Fill Missing Values

# Step 6: Assign Data to Clusters

This step is to assign data to clusters using an existing trained clustering models.



## 6A: Assign Data to Clusters

Properties:
- Check for Append or Uncheck for Result Only:  Checked
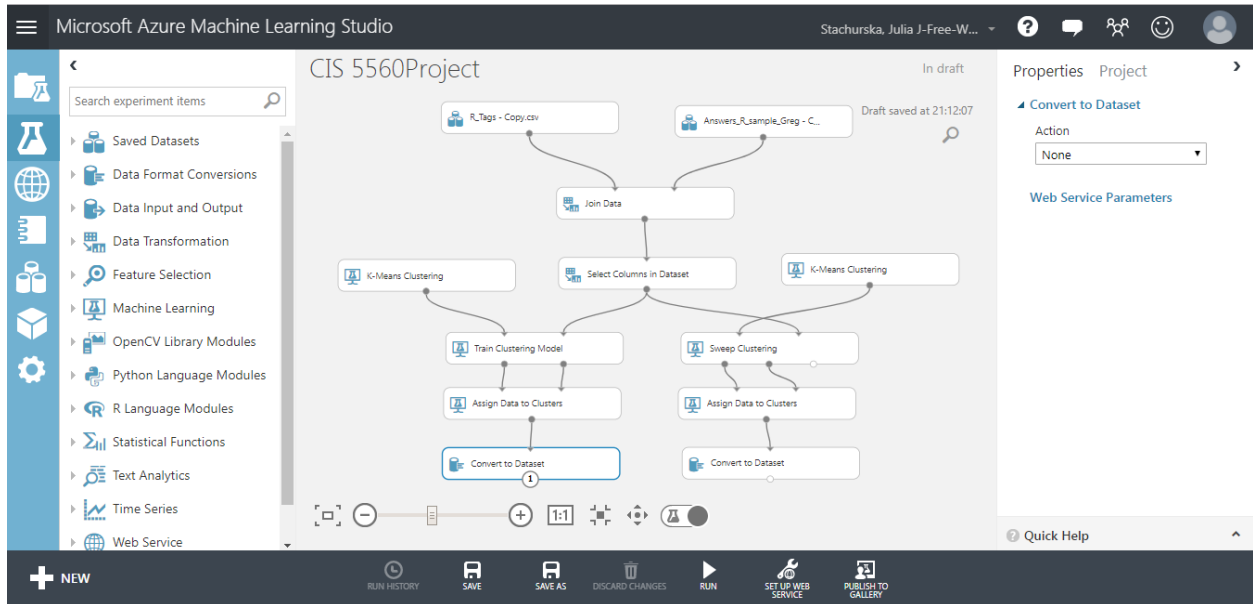
## 6B: Assign Data to Clusters

Properties:
- Check for Append or Uncheck for Result Only:  Checked

# Step 7: Convert to Dataset

This step is to convert data input to the internal Dataset format used by Microsoft Azure Machine Learning



## 7A: Convert to Dataset:
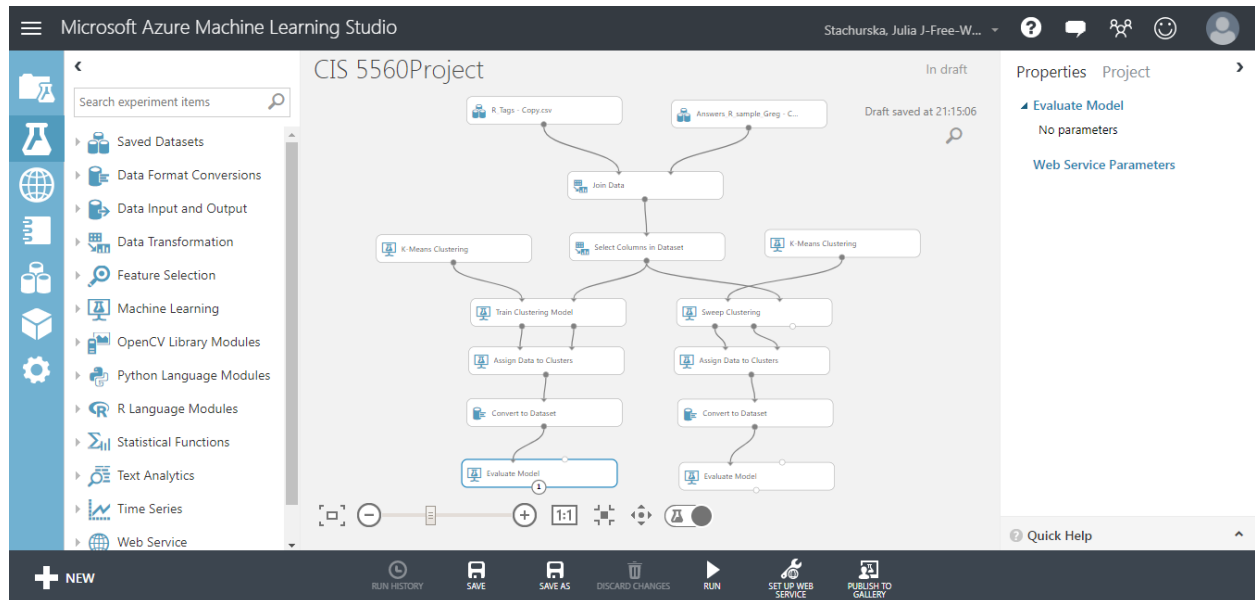
Properties:

- Action: None

## 7B: Convert to Dataset:

Properties:

- Action: None

# Step 8: Evaluate Model

This step is to evaluate a scored model with standard metrics.



# References

- [https://www.kaggle.com/stackoverflow/rquestions](https://www.kaggle.com/stackoverflow/rquestions) 2GB

- [https://gallery.cortanaintelligence.com/Experiment/Clustering-kmeans-2](https://gallery.cortanaintelligence.com/Experiment/Clustering-kmeans-2)