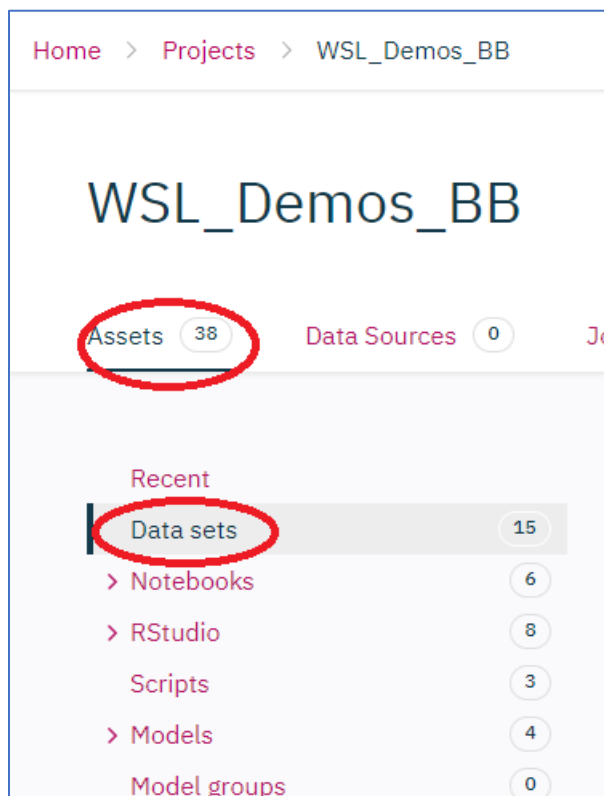# Data Refinery Lab

This lab will use the Titanic data set to demonstrate data profiling, data visualization, and data preparation capabilities of the Data Refinery tool.  The lab consists of the following steps:

1. Use the Data Refinery Tool to:
   a. Profile the data to help determine missing values
   b. Visualize the data to gain a better understanding
   c. Prepare the data for modeling
   d. Run the sequence of data preparation operations on the entire data set.

## Step 1: Add a Data Flow

1. In the Watson Studio Local project, click on **Assets** and click on **Data Sets**.



2. Click on titanic.csv.

3. The Data Refinery panel will display the data set.



## Step 2: Profile the data to help determine missing values.

1. Click on the **Profile** tab.



2. The Profile panel displays the counts of the top 10 count values for each column. Note that you can change 10 to another number if desired. You can also switch to the bottom 10 counts for a column. Scroll to the right to view the cabin column.

3. Note that the cabin column has many missing values and should be removed as part of the data preparation step.



4. In a similar fashion, scroll to the right to examine the boat, body, and home.dest columns. These also have many missing values and should be removed as part of the data preparation step.

5. Age and Embarked also have missing values. Embarked has only 2 missing values. Age has over 100 missing values, but we will keep that column in the analysis. As part of data preparation, we will remove the rows that contain the missing age and embarked values.

6. Click on the **Data** tab. We will add columns that contain more readable values for the survived and pclass columns. The column survived_value will contain a "Y" or "N". The pclass_value column will contain "first", "second", or "third". We will use the mutate (R dpylr function) and ifelse functions to do the conversion. Click on the **Code an operation to cleanse and shape your data.**



7. Type the following:
    mutate(survived_value=ifelse(survived==1, "Y", "N"))

    and then click Apply. If you scroll to the right you should see the new column "survived_value".

8. Type the following to create pclass_value,
   mutate(pclass_value=ifelse(pclass==1,"first",ifelse(pclass==2,"second","third")))



9. The result is shown below. Notice that the right panel will contain a running list of the transformations.



# Step 3: Visualize the data to get a better understanding

1. Click on the **Visualizations** tab.



2. Let's take a look at the breakdown of passengers by passenger class. We will use our new pclass_value field. Enter or select pclass_value and then click **Visualize Data**

3. The result is shown below.



4. We can switch this to a bar chart, by clicking on **Bar**.

5.  A pop-up warning is displayed. Click the **Don't show this again checkbox** and click **Continue**.



6.  The result is shown below.

7. Let's examine the relationship between survival and the passenger class. We will add the survived_value as the **Split by** column and change to **Stacked** view. The result is shown below. We can see that survival probability for first class customers is significantly better.



8. Plot the fare values. The result is shown below. Note that it is highly skewed which affects the performance of some machine learning algorithms. One way to deal with this is to apply a logarithmic transformation. We will do that as part of data preparation.

# Step 4: Prepare the data for modeling

Based on the data analysis, we need to do the following to prepare the data for modeling.

1. Remove columns cabin, boat, body, home.dest
2. Remove rows with missing values of age, and embarked.
3. Create a new column(log_fare) that is the logarithm of the fare column

We will also bin the age, and log_fare fields.

1. Return to the Data panel by clicking on the **Data** tab



2. Remove the cabin column by selecting on the vertical ellipse and then clicking on **Remove**.

| cabin String | embarked String | boat String |
|---|---|---|
| | Remove | |
| B5 | | 2 |
| C22 C26 | Remove duplicates | 11 |
| C22 C26 | Remove empty rows | |
| C22 C26 | | |
| C22 C26 | Sort ascending | |
| E12 | Sort descending | 3 |
| D7 | Substitute | 10 |
| A36 | | |
| C101 | CONVERT COLU... > | D |
| | TEXT > | |
| C62 C64 | View All | |
| C62 C64 | C | 4 |
| B35 | C | 9 |
| | S | 6 |

3. Remove the boat, body, and home.dest columns in a similar manner by selecting on the vertical ellipse adjacent to the column and clicking on **Remove**. Notice the STEPS panel on the right-hand side that provides a running list of the data operations.

**6 STEPS**

Data Source : titanic.csv

Custom code

mutate(survived_value =
ifelse(survived==1,"Y","N"))

Custom code

mutate(pclass_value =
ifelse(pclass==1,"first",ifelse(pclass==
2,"second","third")))

Remove

Removed cabin

Remove

Removed boat

Remove

Removed body

Remove      JUST ADDED

Removed home.dest

4. For the age and embarked columns, click on the vertical ellipse adjacent to the columns, and click on **Remove empty rows**.



5. Convert the fare column from a String to a Decimal by clicking on the vertical ellipse adjacent to the column, click on **Convert Column**, and then click on **Decimal**.



6. Create a new column that is the log to the base 10 of the fare by clicking into the **Code an operation to cleanse and shape your data,** and entering

mutate(log_fare=log10(fare))

then click **Apply**.



7. Convert the age from String to Integer by clicking on the vertical ellipse adjacent to the age column, clicking on **Convert Column**, and clicking on **Integer**.



8. Bin the age column into the following bins by clicking into the **Code an operation to cleanse and shape your data,** and entering

mutate(age_bin=ifelse(age<6,0,ifelse(age<12,1,ifelse(age<18,2,ifelse(age<40,3,ifelse(age<65,4,ifelse(age<80,5,6)))))))

and then click **Apply**.

| Bin | Age Range |
|-----|-----------|
| 0 | 0-5 |
| 1 | 6-11 |
| 2 | 12-17 |
| 3 | 18-39 |
| 4 | 40-64 |
| 5 | 65-79 |
| 6 | Over 79 |
| | |

9. Bin the log_fare column, by clicking into the **Code an operation to cleanse and shape your data,** and entering

mutate(log_fare_bin=ifelse(log_fare<0,0,ifelse(log_fare>8,9,as.integer(log_fare)+1)))

and then clicking **Apply**



10. Now we will drop the age, fare, and log_fare columns as they are no longer needed for modeling purposes. Select the vertical ellipse adjacent to the column and click on **Remove** as shown below.

11. Save the Data Flow by clicking on the Save Data Flow icon .

## Step 5:  Run the sequence of Data Flow operations on the entire data set.

When users are interacting with the Data Refinery tool, the operations are applied to a subset of the data set to facilitate faster response times. To run the data operations on the entire data set (titanic.csv), the user will generate an **R-Script** and run a **Job**.

1. Click on the  icon to **Save R-Script**.



2. Navigate to the project **Scripts** by clicking on the WSL_Demos_XX link and then clicking on **Scripts**.



3. Click on the vertical ellipse for **titanic.R** script, and click on **Create Job**.

4. Enter a **Name**, select **Script Run** for **Type,** select **RStudio with R 3.4.3** for **Worker**, select the **Titanic.R** script for **Source asset**, enter for **arg1** the following **input=datasets/titanic.csv**, click + and enter for **arg2** the following **output=datasets/titanic_shaped_out.csv,** and then click **Create**.

Target host *
Local instance

Source asset *
/scripts/titanic.R                                                    ⌄

Environment variables  ⊕

VARIABLE_1=value 1

Command line arguments  ⊕

input=datasets/titanic.csv

output=datasets/titanic_shaped_out.csv

Scheduled to run *
● On demand
○ Every  day     ⌄   at  12:00 AM ⇕

Cancel    Create

5. Scroll down to the **Runs** section and click **run now.**

**Runs**                                                                      ⊚ run now

| ID | NAME | TARGET HOST | TRIGGERED BY | STARTED AT | DURATION (S) | RESULT |
|----|------|-------------|--------------|------------|--------------|--------|

no runs found

6. Enter a Name for the Run, and then click **Run**.

## Run TitanicDR

**Name** *

Run1

96

**Target host** *

Local instance

**Environment variables** ⊕

*VARIABLE_1=value 1*

**Command line arguments** ⊕

datasets/titanic.csv

datasets/titanic_cleansed_out.csv

Cancel    Run

7.  Wait for the script to complete.



**Run1** ○

No description available.

BB  Bernard Beekman
21 Aug 2018, 11:59 PM

| RUN ID | TARGET HOST | DURATION | RESULT |
|---|---|---|---|
| 1534910392-1001 | Local instance | 3 s | ✓ Success |

8.  When the script completes navigate to Data Sets by clicking on the **WSL_Demos_XX** project link and then click on **Data Sets**. Notice that the **titanic_shaped_out.csv** data set has been created. Click on the vertical ellipse on the right side and click on **Preview**.



| Name | Type | Size | Data Source | Last Modified | |
|---|---|---|---|---|---|
| titanic_shaped_out.csv | CSV | 86.64 KB | Local file | 10 Dec 2018, 7:36 PM | ⋮ |
| titanic.csv | CSV | 105.73 KB | Local file | 10 Dec 2018, 6:47 PM | Preview |
| | | | | | Export |

9. We can see that the new fields defined in the Data Refinery flow have been created. Click on the "x" to close the window.



Preview - titanic_shaped_out.csv

| pclass | survived | name | sex | sibsp | parch | ticket | embarked | survived_value | pclass_value | age_bin | log_fare_bin |
|--------|----------|------|-----|-------|-------|--------|----------|----------------|--------------|---------|--------------|
| 1 | 1 | Allen, Miss. Elisabeth Walton | female | 0 | 0 | 24160 | S | Y | first | 3 | 3 |
| 1 | 1 | Allison, Master. Hudson Trevor | male | 1 | 2 | 113781 | S | Y | first | 0 | 3 |
| 1 | 0 | Allison, Miss. Helen Loraine | female | 1 | 2 | 113781 | S | N | first | 0 | 3 |
| 1 | 0 | Allison, Mr. Hudson Joshua Creighton | male | 1 | 2 | 113781 | S | N | first | 3 | 3 |
| 1 | 0 | Allison, Mrs. Hudson J C (Bessie Waldo Daniels) | female | 1 | 2 | 113781 | S | N | first | 3 | 3 |
| 1 | 1 | Anderson, Mr. Harry | male | 0 | 0 | 19952 | S | Y | first | 4 | 2 |
| 1 | 1 | Andrews, Miss. Kornelia Theodosia | female | 1 | 0 | 13502 | S | Y | first | 4 | 2 |