# Tutorial: Use Azure Machine Learning Workbench for advanced data preparation (Bike share data)

6/13/2018 • 24 minutes to read • Edit Online

Azure Machine Learning (preview) is an integrated, end-to-end data science and advanced analytics solution for professional data scientists to prepare data, develop experiments, and deploy models at cloud scale.

In this tutorial, you use Machine Learning (preview) to learn how to:

- Prepare data interactively with the Machine Learning data preparation tool.
- Import, transform, and create a test dataset.
- Generate a data preparation package.
- Run the data preparation package by using Python.
- Generate a training dataset by reusing the data preparation package for additional input files.
- Execute scripts in a local Azure CLI window.
- Execute scripts in a cloud Azure HDInsight environment.

If you don't have an Azure subscription, create a free account before you begin.

## Prerequisites

- A local installation of Azure Machine Learning Workbench. For more information, follow the installation quickstart.
- If you don't have the Azure CLI installed, follow the instructions to install the latest Azure CLI version.
- An HDInsights Spark cluster created in Azure.
- An Azure storage account.
- Familiarity with how to create a new project in Workbench.
- Although it's not required, it's helpful to have Azure Storage Explorer installed so that you can upload, download, and view the blobs in your storage account.

## Data acquisition

This tutorial uses the Boston hubway dataset and Boston weather data from NOAA.

1. Download the data files from the following links to your local development environment:

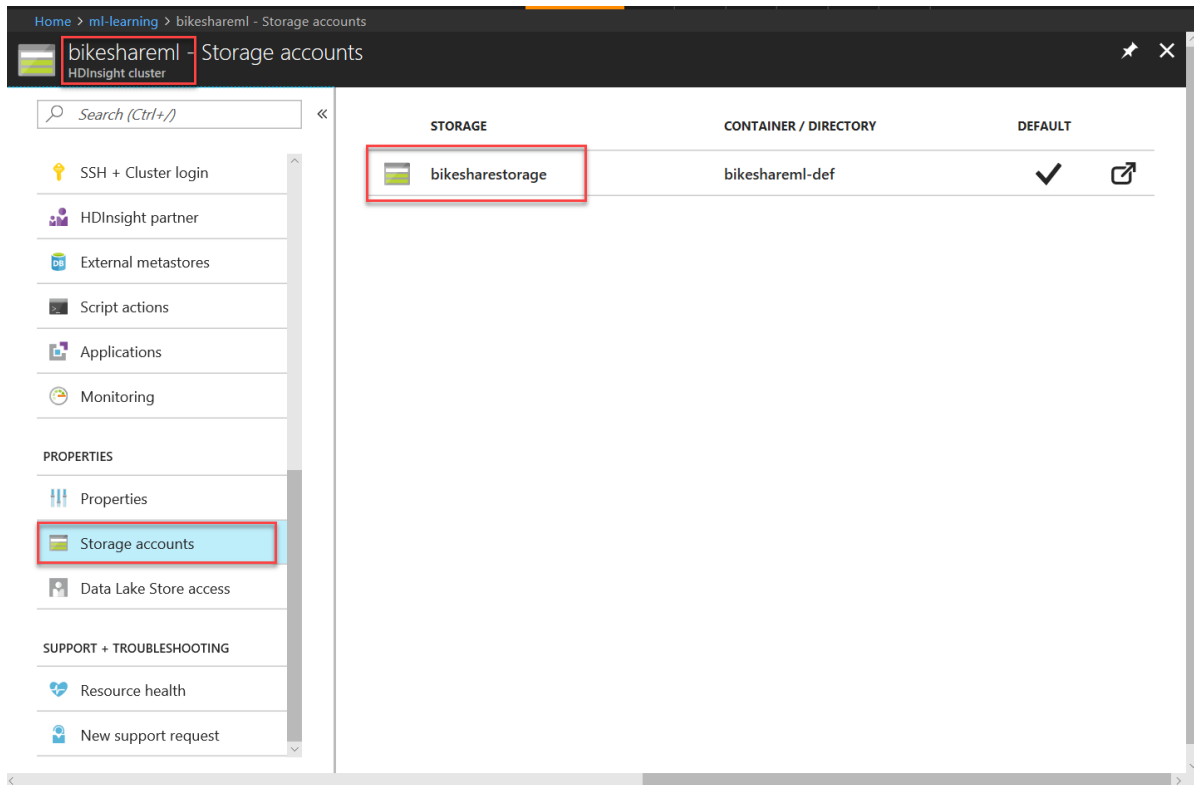   - Boston weather data

   - Hubway trip data from the hubway website:

      - 201501-hubway-tripdata.zip
      - 201504-hubway-tripdata.zip
      - 201510-hubway-tripdata.zip
      - 201601-hubway-tripdata.zip
      - 201604-hubway-tripdata.zip
      - 201610-hubway-tripdata.zip
      - 201701-hubway-tripdata.zip
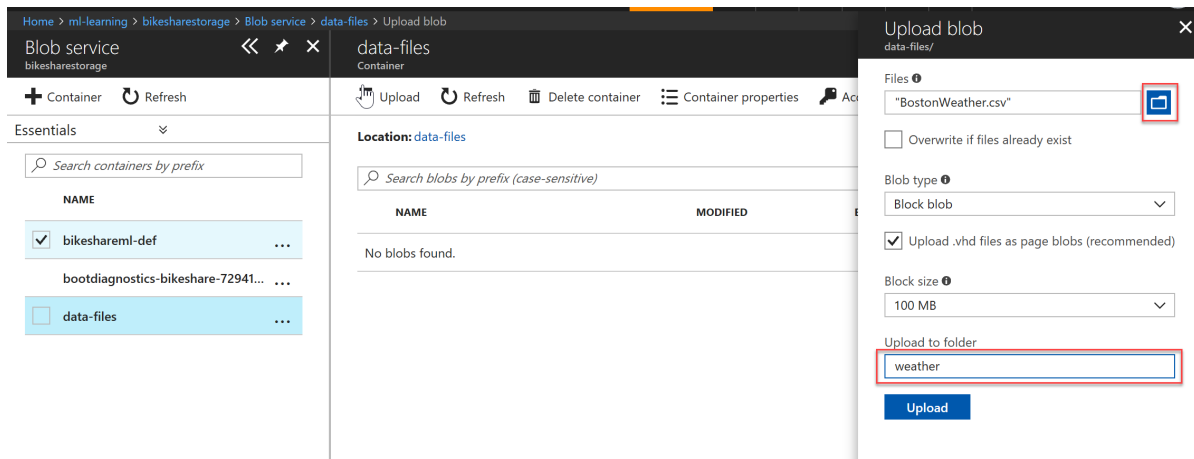
2. Unzip each .zip file after download.

# Upload data files to Azure Blob storage

You can use Azure Blob storage to host your data files.

1. Use the same storage account that is used for the HDInsight cluster you use.



2. Create a new container named **data-files** to store the **BikeShare** data files.

3. Upload the data files. Upload `BostonWeather.csv` to a folder named `weather`. Upload the trip data files to a folder named `tripdata`.



> **TIP**
>
> You also can use Storage Explorer to upload blobs. Use this tool when you want to view the contents of any files generated in the tutorial, too.

# Learn about the datasets

1. The **Boston weather** file contains the following weather-related fields, reported on an hourly basis:

   - **DATE**

- **REPORTTPYE**

- **HOURLYDRYBULBTEMPF**

- **HOURLYRelativeHumidity**

- **HOURLYWindSpeed**

2. The **hubway** data is organized into files by year and month. For example, the file named `201501-hubway-tripdata.zip` contains a .csv file that contains data for January 2015. The data contains the following fields, with each row representing a bike trip:

- **Trip Duration (in seconds)**

- **Start Time and Date**

- **Stop Time and Date**

- **Start Station Name & ID**

- **End Station Name & ID**

- **Bike ID**

- **User Type (Casual = 24-Hour or 72-Hour Pass user; Member = Annual or Monthly Member)**

- **ZIP Code (if user is a member)**

- **Gender (self-reported by member)**

## Create a new project

1. Start **Machine Learning Workbench** from your Start menu or launcher.

2. Create a new Machine Learning project. Select the **+** button on the **Projects** page, or select **File** > **New**.

   - Use the **Bike Share** template.

   - Name your project **BikeShare**.

## Create a new data source

1. Create a new data source. Select the **Data** button (cylinder icon) on the left toolbar to display the **Data** view.

2. Add a data source. Select the **+** icon, and then select **Add Data Source**.



# Add weather data

1. **Data Store**: Select the data store that contains the data. Because you're using files, select **File(s)/Directory**. Select **Next** to continue.

2. **File Selection**: Add the weather data. Browse and select the `BostonWeather.csv` file that you uploaded to Blob Storage earlier. Select **Next**.



3. **File Details**: Verify the file schema that is detected. Machine Learning Workbench analyzes the data in the file and infers the schema to use.

The preview of the data should display the following columns:

- **Path**
- **DATE**
- **REPORTTYPE**
- **HOURLYDRYBULBTEMPF**
- **HOURLYRelativeHumidity**
- **HOURLYWindSpeed**

To continue, select **Next**.

4. **Data Types**: Review the data types that are detected automatically. Machine Learning Workbench analyzes the data in the file and infers the data types to use.

a. For this data, change **DATA TYPE** for all the columns to **String**.

b. To continue, select **Next**.

5. **Sampling**: To create a sampling scheme, select **Edit**. Select the new **Top 10000** row that is added, and then select **Edit**. Set **Sample Strategy** to **Full File**, and then select **Apply**.



To use the **Full File** strategy, select the **Full File** entry, and then select **Set as Active**. A star appears next to **Full File** to indicate that it's the active strategy.

To continue, select **Next**.

6. **Path Column**: Use the **Path Column** section to include the full file path as a column in the imported data. Select **Do Not Include Path Column**.

> **TIP**
>
> Including the path as a column is useful if you're importing a folder of many files with different file names. It's also useful if the file names contain information that you want to extract later.



7. **Finish**: To finish creating the data source, select **Finish**.

A new data source tab named **BostonWeather** opens. A sample of the data is displayed in a grid view. The sample is based on the active sampling scheme specified previously.

Notice that the **Steps** pane on the right side of the screen displays the individual actions taken while creating this data source.

**View data source metrics**

Select **Metrics** at the upper left of the tab's grid view. This view displays the distribution and other aggregated statistics of the sampled data.



> **NOTE**
>
> To configure the visibility of the statistics, use the **Choose Metric** drop-down list. Select and clear metrics there to change the grid view.

To return to the **Data** view, select **Data** in the upper left of the page.

# Add a data source to the data preparation package

1. Select **Prepare** to begin preparing the data.

2. When prompted, enter a name for the data preparation package, such as **BikeShare Data Prep**.

3. Select **OK** to continue.

4. A new package named **BikeShare Data Prep** appears under the **Data Preparation** section of the **Data** tab.

   To display the package, select this entry.

5. Select the **>>** button to expand **Dataflows** and display the dataflows contained in the package. In this example, **BostonWeather** is the only dataflow.

> **IMPORTANT**
> A package can contain multiple dataflows.
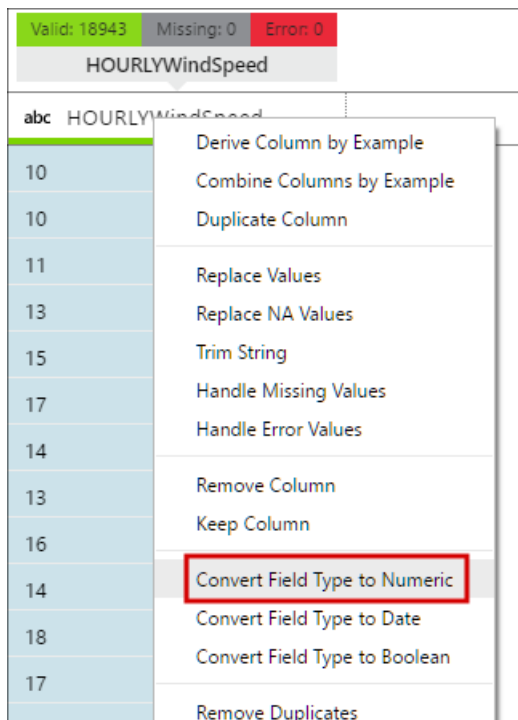


# Filter data by value

1. To filter data, right-click a cell with a certain value, and select **Filter**. Then select the type of filter.

2. For this tutorial, select a cell that contains the value `FM-15`. Then set the filter to **equals**. Now the data is filtered to only return rows where the **REPORTTYPE** is `FM-15`.

> **NOTE**
>
> FM-15 is a type of Meteorological Terminal Aviation Routine (METAR) weather report. The FM-15 reports are empirically observed to be the most complete, with little missing data.

# Remove a column

You no longer need the **REPORTTYPE** column. Right-click the column header, and select **Remove Column**.



# Change datatypes and remove errors

1. Select Ctrl (Command ⌘ on Mac) while you select column headers to select multiple columns at the same time. Use this technique to select the following column headers:

   - **HOURLYDRYBULBTEMPF**

- **HOURLYRelativeHumidity**

- **HOURLYWindSpeed**

2. Right-click one of the selected column headers, and select **Convert Field Type to Numeric**. This option converts the data type for the columns to numeric.



3. Filter out the error values. Some columns have data type conversion problems. This problem is indicated by the red color in the **Data Quality Bar** for the column.

   To remove the rows that have errors, right-click the **HOURLYDRYBULBTEMPF** column heading. Select **Filter Column**. Use the default **I Want To** as **Keep Rows**. Change the **Conditions** drop-down list to select **is not error**. Select **OK** to apply the filter.

4. To eliminate the remaining error rows in the other columns, repeat this filter process for the **HOURLYRelativeHumidity** and **HOURLYWindSpeed** columns.

# Use by example transformations

To use the data in a prediction for two-hour time blocks, you must compute the average weather conditions for two-hour periods. Use the following actions:

- Split the **DATE** column into separate **Date** and **Time** columns. See the following section for the detailed steps.

- Derive an **Hour_Range** column from the **Time** column. See the following section for the detailed steps.

- Derive a **Date_Hour_Range** column from the **DATE** and **Hour_Range** columns. See the following section for the detailed steps.

**Split column by example**

1. Split the **DATE** column into separate **Date** and **Time** columns. Right-click the **DATE** column header, and select **Split Column by Example**.



2. Machine Learning Workbench automatically identifies a meaningful delimiter and creates two columns by splitting the data into date and time values.

3. Select **OK** to accept the split operation results.



**Derive column by example**

1. To derive a two-hour range, right-click the **DATE_2** column header, and select **Derive Column by Example**.



A new empty column is added with null values.

2. Select in the first empty cell in the new column. To provide an example of the time range desired, type

**12AM-2AM** in the new column, and then select Enter.



> **NOTE**
>
> Machine Learning Workbench synthesizes a program based on the examples provided by you and applies the same program on remaining rows. All other rows are automatically populated based on the example you provided. Workbench also analyzes your data and tries to identify edge cases.

> **IMPORTANT**
>
> Identification of edge cases might not work on Mac in the current version of Workbench. Skip the following step 3 and step 4 on Mac. Instead, select **OK** after all the rows are populated with the derived values.

3. The text **Analyzing Data** above the grid indicates that Workbench is trying to detect edge cases. When finished, the status changes to **Review next suggested row** or **No suggestions**. In this example, **Review next suggested row** is returned.

4. To review the suggested changes, select **Review next suggested row**. The cell that you should review and correct, if needed, is highlighted on the display.



Select **OK** to accept the transformation.

5. You are returned to the grid view of data for **BostonWeather**. The grid now contains the three columns added previously.

6. To rename the column, double-click the column header, and type **Hour Range**. Select Enter to save the change.



7. To derive the date and hour range, multi-select the **Date_1** and **Hour Range** columns, right-click, and then select **Derive Column by Example**.



Type **Jan 01, 2015 12AM-2AM** as the example against the first row, and select Enter.

Workbench determines the transformation based on the example you provide. In this example, the result is that the date format is changed and concatenated with the two-hour window.

8. Wait for the status to change from **Analyzing Data** to **Review next suggested row**. This change might take several seconds. Select the status link to go to the suggested row.

| | |
|---|---|
| 8PM-10PM | Dec 01, 2015 8PM-10PM |
| 10PM-12AM | Dec 01, 2015 10PM-12AM |
| 10PM-12AM | Dec 01, 2015 10PM-12AM |
| 12AM-2AM | *null* |
| 12AM-2AM | *null* |
| 2AM-4AM | `null` |

The row has a null value because the source date value can be for either dd/mm/yyyy or mm/dd/yyyy. Type the correct value of **Jan 13, 2015 2AM-4AM**, and select Enter. Workbench uses the two examples to improve the derivation for the remaining rows.

| | |
|---|---|
| 12AM-2AM | Jan 13, 2015 12AM-2AM |
| 2AM-4AM | Jan 13, 2015 2AM-4AM |
| 2AM-4AM | Jan 13, 2015 2AM-4AM |
| 4AM-6AM | Jan 13, 2015 4AM-6AM |

9. Select **OK** to accept the transform.

## BostonWeather  Metrics

| | abc DATE | abc DATE_1 | abc DATE_2 | abc Hour Range | abc Column |
|---|---|---|---|---|---|
| 1 | 1/1/2015 0:54 | 1/1/2015 | 0:54 | 12AM-2AM | Jan 01, 2015 12AM-2AM |
| 2 | 1/1/2015 1:54 | 1/1/2015 | 1:54 | 12AM-2AM | Jan 01, 2015 12AM-2AM |
| 3 | 1/1/2015 2:54 | 1/1/2015 | 2:54 | 2AM-4AM | Jan 01, 2015 2AM-4AM |
| 4 | 1/1/2015 3:54 | 1/1/2015 | 3:54 | 2AM-4AM | Jan 01, 2015 2AM-4AM |
| 5 | 1/1/2015 4:54 | 1/1/2015 | 4:54 | 4AM-6AM | Jan 01, 2015 4AM-6AM |

10. To rename the column, double-click the header. Change the name to **Date Hour Range**, and then select Enter.

11. Multi-select the **DATE**, **DATE_1**, **DATE_2**, and **Hour Range** columns. Right-click, and then select **Remove column**.

## Summarize data (mean)

The next step is to summarize the weather conditions by taking the mean of the values, grouped by hour range.

1. Select the **Date Hour Range** column, and then on the **Transforms** menu, select **Summarize**.

2. To summarize the data, drag columns from the grid at the bottom of the page to the left and right panes at the top. The left pane contains the text **Drag columns here to group data**. The right pane contains the text **Drag columns here to summarize data**.

a. Drag the **Date Hour Range** column from the grid at the bottom to the left pane. Drag **HOURLYDRYBULBTEMPF**, **HOURLYRelativeHumidity**, and **HOURLYWindSpeed** to the right pane.

b. In the right pane, select **Mean** as the **Aggregate** measure for each column. Select **OK** to finish the summarization.

# Transform dataflow by using script

Changing the data in the numeric columns to a range of 0 to 1 allows some models to converge quickly. Currently, there is no built-in transformation to generically do this transformation. Use a Python script to perform this operation.

1. On the **Transform** menu, select **Transform Dataflow (Script)**.

2. Enter the following code in the text box that appears. If you used the column names, the code should work without modification. You are writing a simple min-max normalization logic in Python.

> **WARNING**
>
> The script expects the column names used previously in this tutorial. If you have different column names, you must change the names in the script.

```
maxVal = max(df["HOURLYDRYBULBTEMPF_Mean"])
minVal = min(df["HOURLYDRYBULBTEMPF_Mean"])
df["HOURLYDRYBULBTEMPF_Mean"] = (df["HOURLYDRYBULBTEMPF_Mean"]-minVal)/(maxVal-minVal)
df.rename(columns={"HOURLYDRYBULBTEMPF_Mean":"N_DryBulbTemp"},inplace=True)

maxVal = max(df["HOURLYRelativeHumidity_Mean"])
minVal = min(df["HOURLYRelativeHumidity_Mean"])
df["HOURLYRelativeHumidity_Mean"] = (df["HOURLYRelativeHumidity_Mean"]-minVal)/(maxVal-minVal)
df.rename(columns={"HOURLYRelativeHumidity_Mean":"N_RelativeHumidity"},inplace=True)

maxVal = max(df["HOURLYWindSpeed_Mean"])
minVal = min(df["HOURLYWindSpeed_Mean"])
df["HOURLYWindSpeed_Mean"] = (df["HOURLYWindSpeed_Mean"]-minVal)/(maxVal-minVal)
df.rename(columns={"HOURLYWindSpeed_Mean":"N_WindSpeed"},inplace=True)

df
```

**Transform Dataflow (Script)**                                                          ✕

**Code to Transform Dataflow**

```
1  maxVal = max(df["HOURLYDRYBULBTEMPF_Mean"])
2  minVal = min(df["HOURLYDRYBULBTEMPF_Mean"])
3  df["HOURLYDRYBULBTEMPF_Mean"] = (df["HOURLYDRYBULBTEMPF_Mean"]-minVal)/(maxVal-minVal)
4  df.rename(columns={"HOURLYDRYBULBTEMPF_Mean":"N_DryBulbTemp"},inplace=True)
5
6  maxVal = max(df["HOURLYRelativeHumidity_Mean"])
7  minVal = min(df["HOURLYRelativeHumidity_Mean"])
8  df["HOURLYRelativeHumidity_Mean"] = (df["HOURLYRelativeHumidity_Mean"]-minVal)/(maxVal-minVal)
9  df.rename(columns={"HOURLYRelativeHumidity_Mean":"N_RelativeHumidity"},inplace=True)
10
11 maxVal = max(df["HOURLYWindSpeed_Mean"])
12 minVal = min(df["HOURLYWindSpeed_Mean"])
13 df["HOURLYWindSpeed_Mean"] = (df["HOURLYWindSpeed_Mean"]-minVal)/(maxVal-minVal)
14 df.rename(columns={"HOURLYWindSpeed_Mean":"N_WindSpeed"},inplace=True)
```

**Code Block Type**

| Expression | ▾ |
|---|---|

**Hint**

Please provide Python (3.5) code that transforms your data.
A Pandas DataFrame called 'df' has been made available to your code. Your code should either modify 'df' or reassign a new DataFrame to 'df' before it completes.
The following Python imports are provided: math, numbers, datetime, re, pandas (aliased as pd), numpy (aliased as np), scipy (aliased as sp).
Module signature: def transform(df): return transformed Pandas DataFrame.

[ OK ]     [ Cancel ]

3. Select **OK** to use the script. The numeric columns in the grid now contain values in the range of 0 to 1.

**BostonWeather**

| | abc Date Hour... | # N_DryBulb... | # N_Relative... | # N_WindSp... |
|---|---|---|---|---|
| 1 | Jan 01, 2015 12.. | 0.29383886255... | 0.42528735632... | 0.28169014084... |
| 2 | Jan 01, 2015 2... | 0.30331753554... | 0.40229885057... | 0.33802816901... |
| 3 | Jan 01, 2015 4... | 0.29857819905... | 0.42528735632... | 0.45070422535... |
| 4 | Jan 01, 2015 6... | 0.29857819905... | 0.43678160919... | 0.38028169014... |
| 5 | Jan 01, 2015 8... | 0.33649289099... | 0.32758620689... | 0.42253521126... |

You have finished preparing the weather data. Next, prepare the trip data.
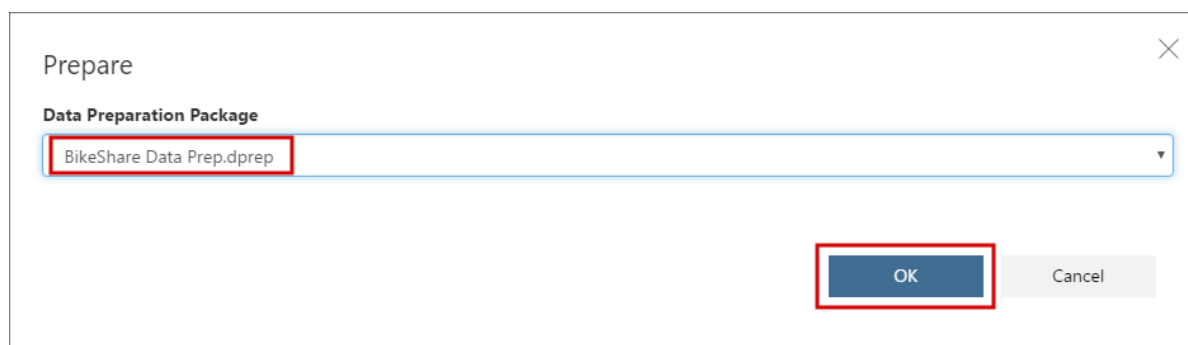
# Load trip data

1. To import the `201701-hubway-tripdata.csv` file, use the steps in the Create a new data source section. Use the following options during the import process:
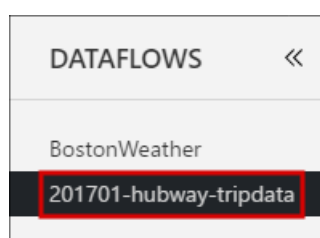
   - **File Selection**: Select **Azure Blob** when you browse to select the file.

   - **Sampling Scheme**: Select **Full File** sampling scheme, and make the sample active.

   - **Data Type**: Accept the defaults.

2. After you import the data, select **Prepare** to begin preparing the data. Select the existing **BikeShare Data Prep.dprep** package, and then select **OK**.

   This process adds a **Dataflow** to the existing **Data Preparation** file rather than creating a new one.



3. After the grid has loaded, expand **DATAFLOWS**. There are now two dataflows: **BostonWeather** and **201701-hubway-tripdata**. Select the **201701-hubway-tripdata** entry.



## Use the map inspector

For data preparation, useful visualizations called inspectors are available for string, numeric, and geographical data. They help you to understand the data better and identify outliers. Follow these steps to use the map inspector.

1. Multi-select the **start station latitude** and **start station longitude** columns. Right-click one of the columns, and then select **Map**.

   > **TIP**
   > To enable multi-select, hold down the Ctrl key (Command ⌘ on Mac), and select the header for each column.

2. To maximize the map visualization, select the **Maximize** icon. To fit the map to the window, select the **E** icon on the upper-left side of the visualization.



3. Select the **Minimize** button to return to the grid view.

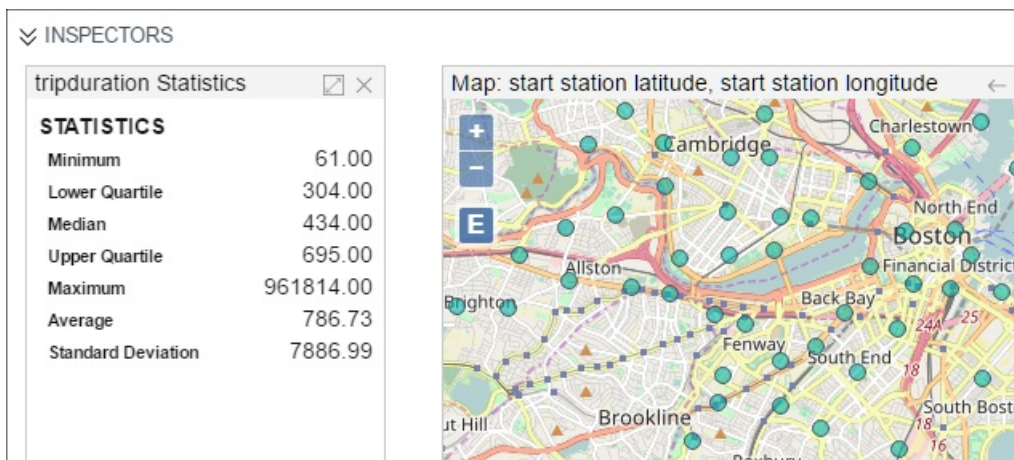## Use the column statistics inspector

To use the column statistics inspector, right-click the **tripduration** column, and select **Column Statistics**.

This process adds a new visualization titled **tripduration Statistics** in the **INSPECTORS** pane.
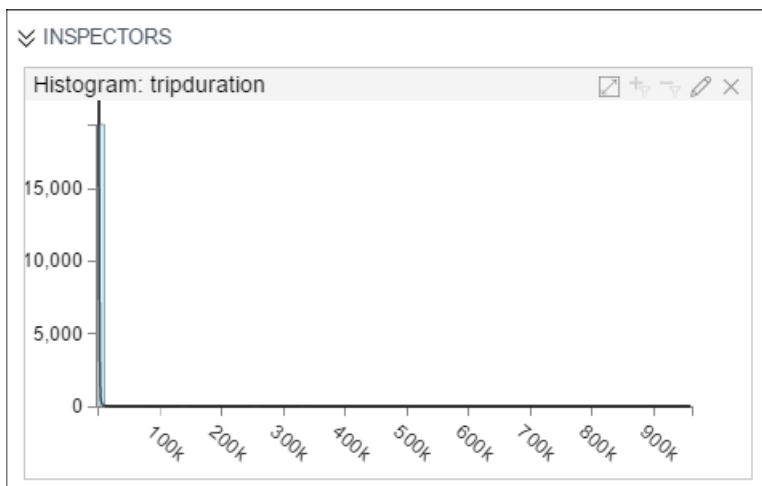


> **IMPORTANT**
>
> The maximum value of the trip duration is 961,814 minutes, which is about two years. It seems there are some outliers in the dataset.

# Use the histogram inspector

To attempt to identify outliers, right-click the **tripduration** column, and select **Histogram**.

The histogram isn't helpful because the outliers skew the graph.

# Add a column by using script

1. Right-click the **tripduration** column, and select **Add Column (Script)**.



2. In the **Add Column (Script)** dialog box, use the following values:

   - **New Column Name**: logtripduration

   - **Insert this New Column After**: tripduration

   - **New Column Code**: `math.log(row.tripduration)`
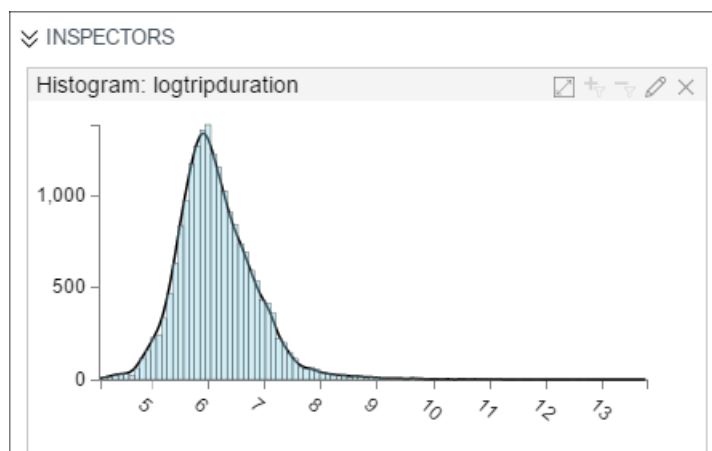
   - **Code Block Type**: Expression

3. Select **OK** to add the **logtripduration** column.

4. Right-click the column, and select **Histogram**.



Visually, this histogram seems like a normal distribution with an abnormal tail.

## Use an advanced filter

Using a filter on the data updates the inspectors with the new distribution.

1. Right-click the **logtripduration** column, and select **Filter Column**.

2. In the **Edit** dialog box, use the following values:

   - **Filter this Number Column**: logtripduration

   - **I Want To**: Keep Rows

   - **When**: Any of the Conditions below are True (logical OR)

   - **If this Column**: less than

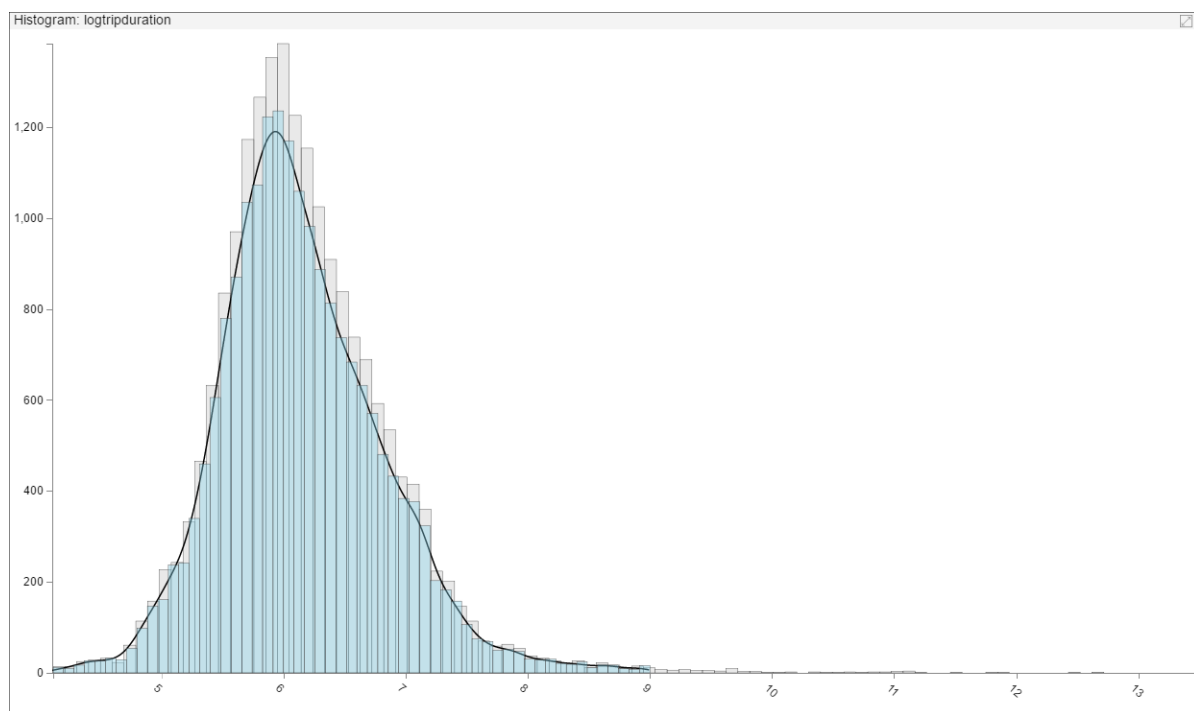- **The Value**: 9



3. Select **OK** to apply the filter.



**The halo effect**

1. Maximize the **logtripduration** histogram. A blue histogram is overlaid on a gray histogram. This display is called the **Halo Effect**:

   - The gray histogram represents the distribution before the operation (in this case, the filtering operation).

   - The blue histogram represents the histogram after the operation.

   The halo effect helps with visualizing the effect of an operation on the data.

Histogram: logtripduration

> **NOTE**
>
> The blue histogram appears shorter compared to the previous one. This difference is due to automatic re-bucketing of data in the new range.

2. To remove the halo, select **Edit** and clear **Show halo**.



Histogram

Column

logtripduration

Minimum Number of Buckets (applies even when default bucketing is checked)

6

☑ Default Number of Buckets (Scott's Rule)
☐ Show halo
☑ Kernel Density Plot Overlay (Gaussian Kernel)

3. Select **OK** to disable the halo effect. Then minimize the histogram.

**Remove columns**

In the trip data, each row represents a bike pickup event. For this tutorial, you need only the **starttime** and **start station id** columns. To remove the other columns, multi-select these two columns, right-click the column header, and then select **Keep Column**. Other columns are removed.

| | # tripduration | 🗓 starttime | 🗓 stoptime | # start statio... | abc start statio... | # start statio... |
|---|---|---|---|---|---|---|
| 1 | 350 | 2017-01-01 00... | 2017-01-01 00... | | | 1 |
| 2 | 891 | 2017-01-01 00... | 2017-01-01 00... | | | 3 |
| 3 | 1672 | 2017-01-01 00... | 2017-01-01 00... | | | 3 |
| 4 | 747 | 2017-01-01 00... | 2017-01-01 00... | | | 4 |
| 5 | 621 | 2017-01-01 00... | 2017-01-01 00... | | | 6 |
| 6 | 664 | 2017-01-01 00... | 2017-01-01 00... | | | 1 |
| 7 | 260 | 2017-01-01 00... | 2017-01-01 00... | | | 1 |
| 8 | 403 | 2017-01-01 00... | 2017-01-01 00... | | | 5 |
| 9 | 642 | 2017-01-01 00... | 2017-01-01 00... | | | 6 |

Context menu:
- Derive Column by Example
- Combine Columns by Example
- Duplicate Column
- Replace NA Values
- Handle Missing Values
- Handle Error Values
- Remove Column
- Keep Column
- Convert Field Type to String
- Convert Unix Timestamp to DateTime

## Summarize data (count)

To summarize bike demand for a two-hour period, use derived columns.

1. Right-click the **starttime** column, and select **Derive Column by Example**.



| | starttime | |
|---|---|---|
| 1 | 2017-01-01 00:06:5 | Derive Column by Example |
| 2 | 2017-01-01 00:13:1 | Split Column by Example |
| 3 | 2017-01-01 00:16:1 | Duplicate Column |
| 4 | 2017-01-01 00:21:2 | Text Clustering |
| 5 | 2017-01-01 00:30:0 | Replace NA Values / Handle Missing Values |

2. For the example, enter a value of **Jan 01, 2017 12AM-2AM** for the first row.

> **IMPORTANT**
>
> In the previous example of deriving columns, you used multiple steps to derive a column that contained the date and time period. In this example, you can see that this operation can be performed as a single step by providing an example of the final output.

> **NOTE**
>
> You can give an example against any of the rows. For this example, the value of **Jan 01, 2017 12AM-2AM** is valid for the first row of data.

DERIVE COLUMN BY EXAMPLE: You have selected 1 source column and provided 0 examples. No suggestions   Advanced mode

| | 🗓 starttime | abc Column | # start stati... |
|---|---|---|---|
| 1 | 2017-01-01 00:06:58 | Jan 01, 2017 12AM-2AM | 67 |
| 2 | 2017-01-01 00:13:16 | null | 36 |
| 3 | 2017-01-01 00:16:17 | null | 36 |
| 4 | 2017-01-01 00:21:22 | null | 46 |

3. Wait until the application computes the values against all the rows. The process might take several seconds. After the analysis is finished, use the **Review next suggested row** link to review data.



Ensure that the computed values are correct. If not, update the value with the expected value, and select Enter. Then wait for the analysis to finish. Complete the **Review next suggested row** process until you see **No suggestions**. **No suggestions** means the application looked at the edge cases and is satisfied with the synthesized program. It's a best practice to perform a visual inspection of the transformed data before you accept the transformation.

4. Select **OK** to accept the transform. Rename the newly created column to **Date Hour Range**.



5. Right-click the **starttime** column header, and select **Remove column**.

6. To summarize the data, on the **Transform** menu, select **Summarize**. To create the transformation, use the following steps:

- Drag **Date Hour Range** and **start station id** to the **Group By** pane on the left.

- Drag **start station id** to the **summarize data** pane on the right.

7. Select **OK** to accept the summary result.

## Join dataflows

To join the weather data with the trip data, use the following steps:

1. On the **Transforms** menu, select **Join**.

2. **Tables**: Select **BostonWeather** as the **Left** dataflow and **201701-hubway-tripdata** as the **Right** dataflow. To continue, select **Next**.

3. **Key Columns**: Select the **Date Hour Range** column in both the tables, and then select **Next**.

4. **Join Type**: Select **Matching rows** as the join type, and then select **Finish**.



This process creates a new dataflow named **Join Result**.

# Create additional features

1. To create a column that contains the day of the week, right-click the **Date Hour Range** column and select **Derive Column by Example**. Use a value of **Sun** for a date that occurred on a Sunday. An example is **Jan 01, 2017 12AM-2AM**. Select Enter, and then select **OK**. Rename this column to **Weekday**.

2.  To create a column that contains the time period for a row, right-click the **Date Hour Range** column, and select **Derive Column by example**. Use a value of **12AM-2AM** for the row that contains **Jan 01, 2017 12AM-2AM**. Select Enter, and then select **OK**. Rename this column to **Period**.



3.  To remove the **Date Hour Range** and **r_Date Hour Range** columns, select Ctrl (Command ⌘ on Mac), and then select each column header. Right-click, and select **Remove Column**.

# Read data from Python

You can run a data preparation package from Python or PySpark and retrieve the result as a **Data Frame**.

To generate an example Python script, right-click **BikeShare Data Prep**, and select **Generate Data Access Code File**. The example Python file is created in your **Project Folder** and is also loaded in a tab within Workbench. The following Python script is an example of the code that is generated:

```python
# Use the Azure Machine Learning data preparation package
from azureml.dataprep import package

# Use the Azure Machine Learning data collector to log various metrics
from azureml.logging import get_azureml_logger
logger = get_azureml_logger()

# This call will load the referenced package and return a DataFrame.
# If run in a PySpark environment, this call returns a
# Spark DataFrame. If not, it will return a Pandas DataFrame.
df = package.run('BikeShare Data Prep.dprep', dataflow_idx=0)

# Remove this line and add code that uses the DataFrame
df.head(10)
```

For this tutorial, the name of the file is `BikeShare Data Prep.py`. This file is used later in the tutorial.

# Save test data as a CSV file

To save the **Join Result** dataflow to a .csv file, you must change the `BikeShare Data Prep.py` script.

1.  Open the project for editing in Visual Studio Code.

2. Update the Python script in the `BikeShare Data Prep.py` file by using the following code:

```python
import pyspark

from azureml.dataprep.package import run
from pyspark.sql.functions import *

# start Spark session
spark = pyspark.sql.SparkSession.builder.appName('BikeShare').getOrCreate()

# dataflow_idx=2 sets the dataflow to the 3rd dataflow (the index starts at 0), the Join Result.
df = run('BikeShare Data Prep.dprep', dataflow_idx=2)
df.show(n=10)
row_count_first = df.count()

# Example file name: 'wasb://data-files@bikesharestorage.blob.core.windows.net/testata'
# 'wasb://<your container name>@<your azure storage name>.blob.core.windows.net/<csv folder name>
blobfolder = 'Your Azure Storage blob path'

df.write.csv(blobfolder, mode='overwrite')

# retrieve csv file parts into one data frame
csvfiles = "<Your Azure Storage blob path>/*.csv"
df = spark.read.option("header", "false").csv(csvfiles)
row_count_result = df.count()
print(row_count_result)
if (row_count_first == row_count_result):
    print('counts match')
else:
    print('counts do not match')
print('done')
```

3. Replace `Your Azure Storage blob path` with the path to the output file to be created. Replace for both the `blobfolder` and `csvfiles` variables.

# Create an HDInsight run configuration

1. In Machine Learning Workbench, open the command-line window, select the **File** menu, and then select **Open Command Prompt**. Your command prompt starts in the project folder with the prompt `C:\Projects\BikeShare>`.



> **IMPORTANT**
> You must use the command-line window (opened from Workbench) to accomplish the steps that follow.

2. Use the command prompt to sign in to Azure.

   The Workbench app and CLI use independent credential caches when you authenticate against Azure resources. You need to do this only once until the cached token expires. The `az account list` command returns the list of subscriptions available to your login. If there is more than one, use the ID value from the desired subscription. Set that subscription as the default account to use with the `az account set -s` command, and then provide the subscription ID value. Then confirm the setting by using the account `show` command.

```
REM login by using the aka.ms/devicelogin site
az login

REM lists all Azure subscriptions you have access to
az account list -o table

REM sets the current Azure subscription to the one you want to use
az account set -s <subscriptionId>

REM verifies that your current subscription is set correctly
az account show
```

3.  Create the HDInsight run config. You need the name of your cluster and the `sshuser` password.

```
az ml computetarget attach cluster --name hdinsight --address <yourclustername>.azurehdinsight.net --
username sshuser --password <your password>
az ml experiment prepare -c hdinsight
```

> **NOTE**
>
> When a blank project is created, the default run configurations are **local** and **docker**. This step creates a new run configuration that is available in Workbench when you run your scripts.

## Run in an HDInsight cluster

Return to the Machine Learning Workbench application to run your script in the HDInsight cluster.

1.  Return to the home screen of your project by selecting the **Home** icon on the left.

2.  Select **hdinsight** from the drop-down list to run your script in the HDInsight cluster.

3.  Select **Run**. The script is submitted as a job. The job status changes to **Completed** after the file is written to the specified location in your storage container.

# Substitute data sources

In the previous steps, you used the `201701-hubway-tripdata.csv` and `BostonWeather.csv` data sources to prepare the test data. To use the package with the other trip data files, use the following steps:

1. Create a new data source by using the steps given previously, with the following changes to the process:

   - **File Selection**: When you select a file, multi-select the six remaining trip tripdata .csv files.

     

     > **NOTE**
     >
     > The **+5** entry indicates that there are five additional files beyond the one that is listed.

   - **File Details**: Set **Promote Headers Mode** to **All Files Have The Same Headers**. This value indicates that each of the files contains the same header.

     

     Save the name of this data source because it's used in later steps.

2. Select the folder icon to view the files in your project. Expand the **aml_config** directory, and then select the `hdinsight.runconfig` file.

3. Select the **Edit** button to open the file in Visual Studio Code.

4. Add the following lines at the end of the `hdinsight.runconfig` file, and then select the disk icon to save the file.

```
DataSourceSubstitutions:
    201701-hubway-tripdata.dsource: 201501-hubway-tripdata.dsource
```

This change replaces the original data source with the one that contains the six trip data files.

## Save training data as a CSV file

1. Browse to the Python file `BikeShare Data Prep.py` that you edited previously. Provide a different file path to save the training data.