



AWS Builders Day

Getting Started with Amazon Machine Learning

February 2018

Table of Contents

Overview.....	3
Download and Review the Bike Sharing Dataset.....	4
Difference between BI and Machine Learning	5
Download the Dataset and Upload to Amazon S3	7
Create a Datasource and a Machine Learning Model.....	9
Try Real-time Prediction.....	15
Evaluate the Model.....	16
Improve Prediction Accuracy.....	19
Additional Available Models on Amazon Machine Learning Service.....	26
Clean Up	27

Overview

Amazon Machine Learning is an easy to use and robust platform for developers, analysts, and business domain experts of all skill levels. This service allows users to build and train predictive models, using user-friendly visualization tools and wizards, without having to learn complex machine learning algorithms and technology. Users can focus on experimenting and improving models from convert ideas to facts, without worrying about operating and maintaining a highly scalable and reliable infrastructure. Amazon Machine Learning offers a pay-as-you-go model where no upfront hardware or software investment is required, and it has the capability to scale with demand to billions of predictions.

In this lab, we will walk you through creating, training, improving and finally using the machine learning models to perform real-time predictions with Amazon Machine Learning. For the purpose of the tutorial, we will work with the bike sharing dataset from [UCI Machine Learning Repository](#). The sample dataset contains data generated by bike sharing systems, where bikes can be rented via a network of kiosk locations. The goal for this lab is to forecast bike rental demand based on historical usage pattern.

The steps we need to follow to build a machine learning model include:

1. Download the bike sharing dataset and upload it into an Amazon S3 bucket
2. Create a new datasource in the Amazon Machine Learning console using the dataset in the S3 bucket
3. Define a schema on the datasource and build a predictive model
4. Make real-time predictions using the model trained

Download and Review the Bike Sharing Dataset

Before we start training the data, let's review it first by downloading the biking sharing dataset from [here](#). Open the **bike_share_data.csv** file with a text editor such as Notepad. When working with machine learning, majority of the time is spent in gathering, preparing, cleaning and analyzing the data before training it with machine learning. For example, you should look for features or information that have higher predictive power. You should also consider how to deal with missing values and outliers in data. It is also recommended to randomly shuffle the data, especially for [online machine learning algorithm](#) to get a better performing model. It is an iterative process to develop a successful predictive model. The dataset from UCI is already cleaned and formatted, we also have downloaded and shuffled the dataset in the interest of time and keeping the focus on the Amazon Machine Learning service.

The first line in the dataset is column headers and they represent the following –

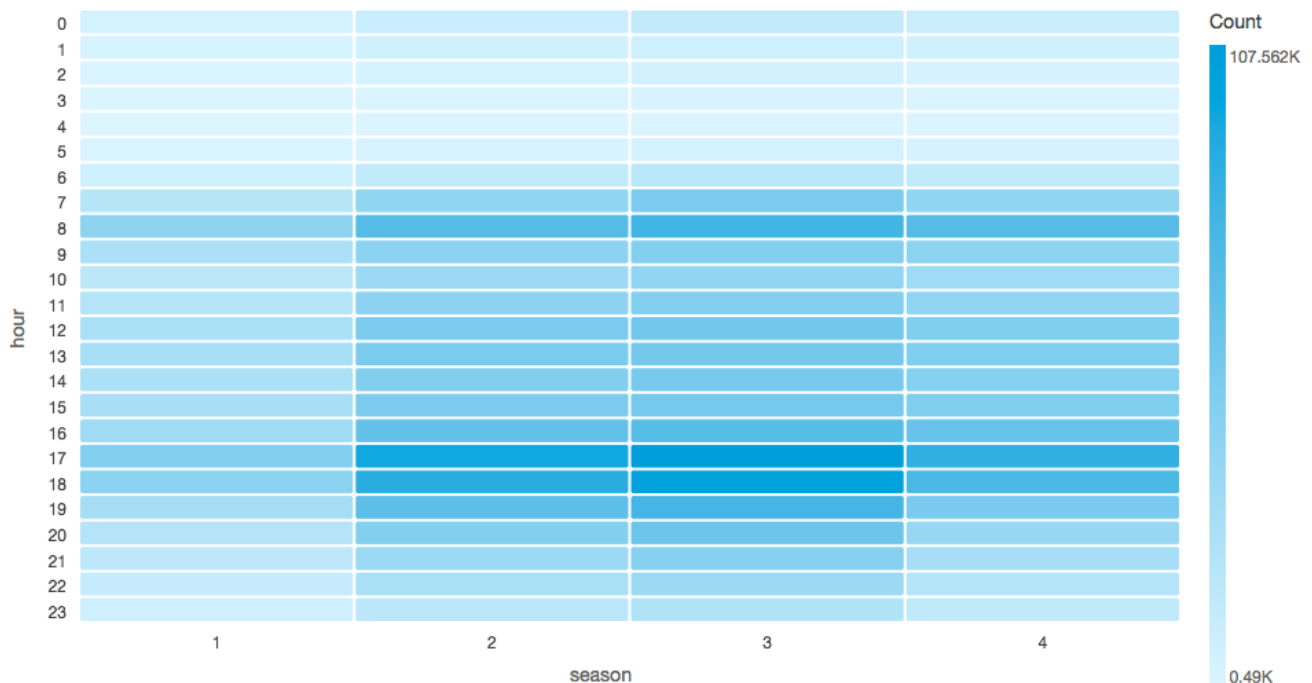
Column Name	Description
season	1-Spring, 2-Summer, 3-Fall, 4-Winter
holiday	1-Yes, 0-No
weekday	0-Sunday, 1-Monday, 2-Tuesday, 3-Wednesday, 4-Thursday, 5-Friday, 6-Saturday
workingday	1-Yes, 0-No
weather	1-Clear, Few clouds, Partly cloudy, Partly cloudy 2-Misty and Cloudy, Misty with Broken clouds, Misty with Few clouds, Misty 3-Light Snow, Light Rain and Thunderstorm, Light Rain and Scatter clouds 4-Heavy Rain and Ice Pallets, Thunderstorm, Snow and Fog
temp	Normalized temperature in Celsius
atemp	Normalized feeling temperature in Celsius
humidity	Normalized humidity
count	Count of total rental bikes aggregated in one hour
datetime	The hour and date

Difference between BI and Machine Learning

Traditionally, business typically analyze and discover trends based on Business Intelligence (BI) Tools with many different reports and charts. This is still useful to analyze and identifying trends. For this lab, we leveraged Amazon QuickSight to visualize and discover usage patterns in the dataset. Amazon QuickSight is a fast business analytics service you can use to build visualizations perform ad-hoc analysis, and quickly gain business insights from your data.

In the screenshot below, we can see that demands are higher during summer (season=2) and fall (season=3) seasons. We can also see that peak time rentals are around 7-8 a.m. and 4-6 p.m.

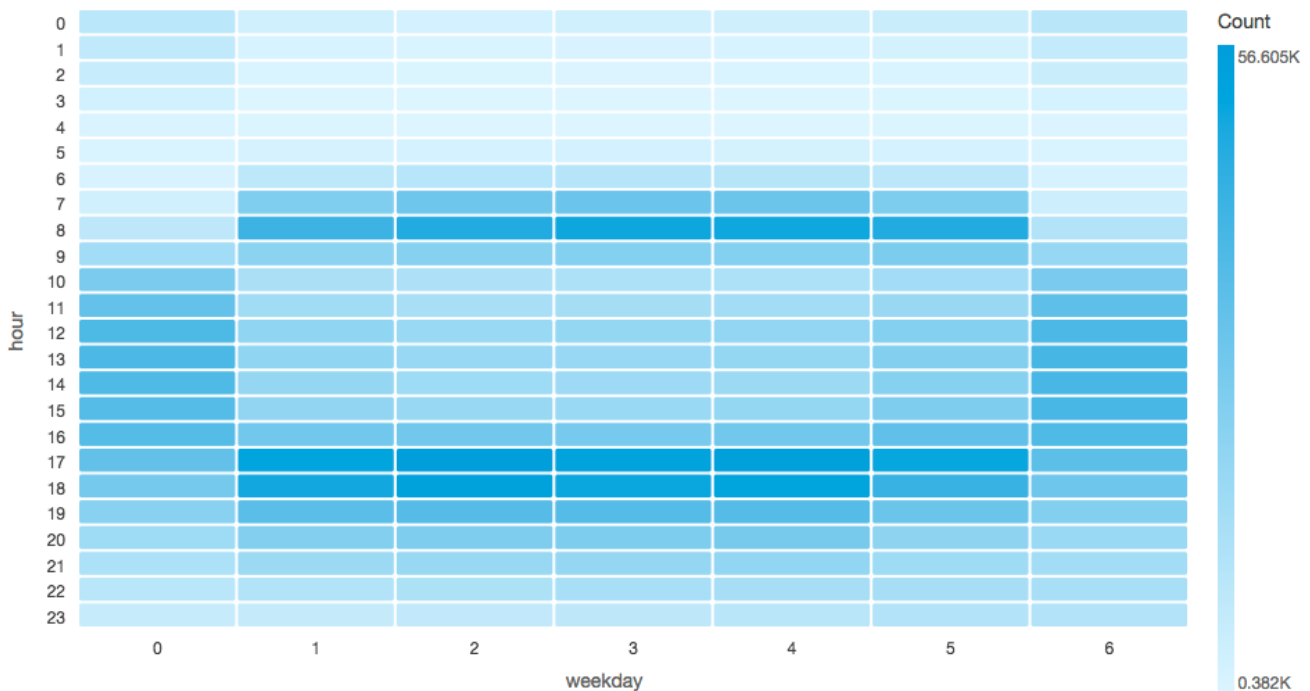
Bike Sharing Density by Season and Hour



Let's take a look at a different chart. In the screenshot below, we can see that during the work days (weekday=1-5), peak times are 7-8 a.m. and 5-6 p.m. And during the weekends

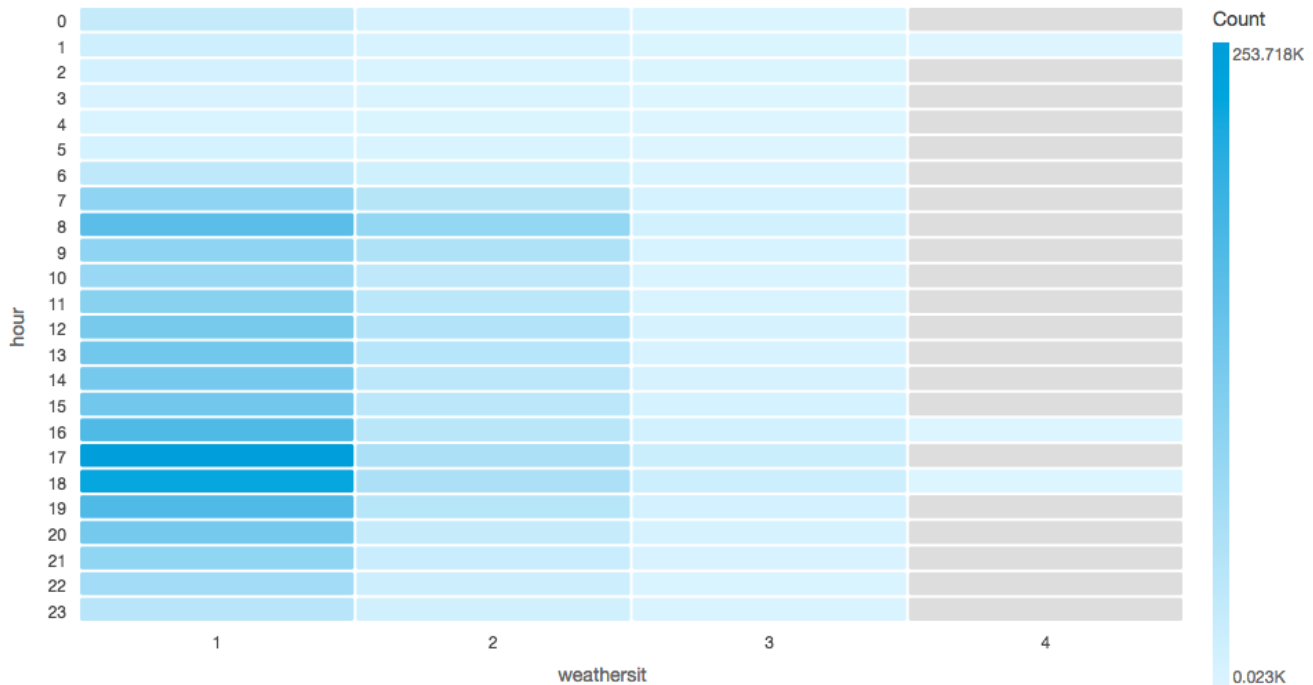
(weekday=0 and 6), peak times are between 10 a.m. – 5 p.m.

Bike Sharing Demand by Weekday and Hour



Let's look at another chart. In the screenshot below, we can see that when weather is sunny and clear (weathersit=1), we have a healthy usage. But when there is heavy rain, thunderstorm or snow (weathersit=4), we see there is hardly any usage.

Bike Sharing Demand by Weather and Hour

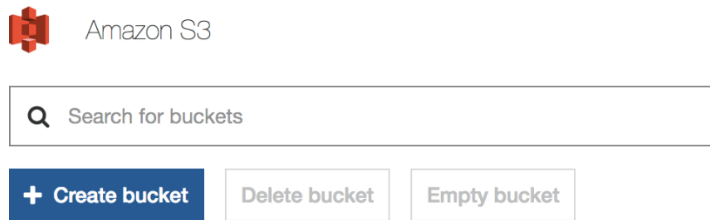


Using a BI tool such as Amazon QuickSight provides the ability to quickly analyze and discover trends in the dataset. This is very helpful; however, it is hard to provide a specific answer given a set of parameters such as specific date and time, weather forecast, day of the week, season and etc in an automated fashion. Machine learning is able to find a pattern in the dataset and predict an answer given the same set of parameters. In the following sections, we will leverage Amazon Machine Learning service to find answers.

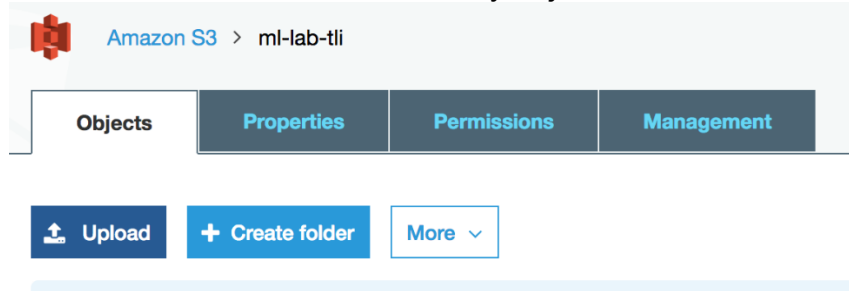
Download the Dataset and Upload to Amazon S3

In this section, we will download the bike sharing datasets and load them into an Amazon S3 bucket. Amazon Machine Learning currently supports S3, Redshift and Relational Database Service (RDS) as the datasource.

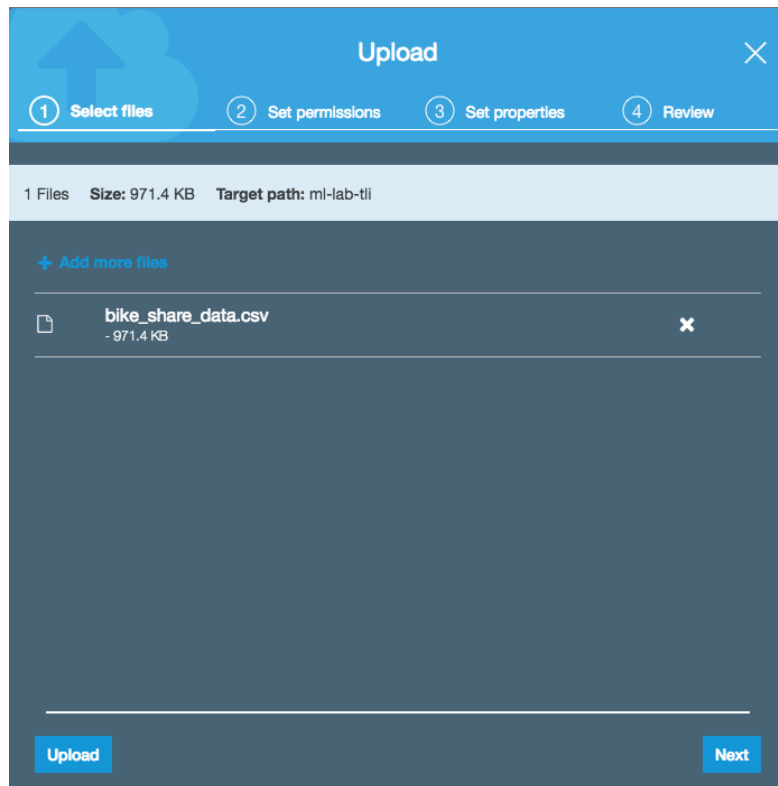
1. If you have not downloaded the bike sharing dataset yet, be sure to download it from [here](#).
2. Sign into the AWS Management Console and open the Amazon S3 console at <https://console.aws.amazon.com/s3>
3. In the upper-right corner of the AWS Management Console, confirm you are in the desired AWS region (e.g., Oregon).
4. Now, we will need to create a bucket. In the S3 console, click the **Create Bucket** button.



5. For the **Bucket Name**, type “ml-lab-**<your-initials>**” in the text box and click **Next** (take note of the bucket name, it will be needed later for creating the datasource). Leave everything default in the next 2 pages and click **Create Bucket** in **Review** page.
6. Click the link for the bucket name you just created, then click **Upload**.



7. Click **Add Files**, find and select the **bike_share_data.csv** file and click **Upload**.



Create a Datasource and a Machine Learning Model

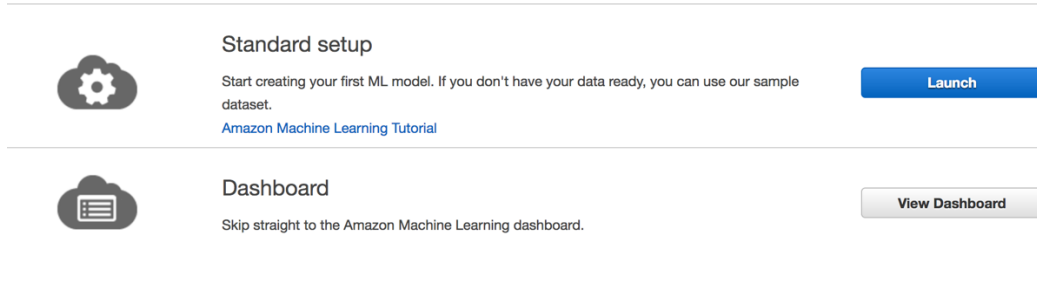
In this section, we will create the datasource needed to train the model with Amazon ML.

1. Open the Amazon Machine Learning console at <https://console.aws.amazon.com/machinelearning/home>
2. If you see the Getting Started page, click on **Get started** button. If you do not see it, go to step 4.



3. Click on **View Dashboard** button

Get started with Amazon Machine Learning



4. Click on **Create New...** button and select “**Datasource and ML Model**”

Getting Started with Amazon Machine Learning

Objects

Create new... Actions Refresh

Datasource and ML model

Datasource

ML model

Evaluation

Batch prediction

You have no objects yet. To get started, choose **Create new**, and then select **Datasource and ML model**.

Machine Learning Concepts

Amazon Machine Learning (Amazon ML) can solve business problems by finding and learning the patterns in your historical data and using the patterns to generate predictions. To get started, you provide Amazon ML with your data. Next, you use Amazon ML to train your ML model, and then you evaluate the model's performance. Finally, you use the model to generate predictions on new data.

[Learn how to use Amazon Machine Learning by walking through the Amazon ML tutorial.](#)

- Under “**Where is your data?**” section, make sure the S3 icon is selected. In the S3 location text box, input the bucket name “ml-lab-**<your-initials>**”. Auto-complete will prompt a list of matches. Select the one that you created. Auto-complete will then prompt you with list of files in the bucket, select **bike_share_data.csv** from the list.

Input data

Import your data to create an Amazon ML datasource. Amazon ML can use your datasource to create and evaluate an ML model, and you can use the data.

Where is your data?



Amazon Redshift

S3 data access

Tell Amazon ML how to access your data and give it permission to access it.

S3 location *

s3:// ml-lab-tli/

Enter the [Learn more](#)

ml-lab-tli/bike_share_data.csv

ml-lab-tli/bike_share_engineered_data.csv

If you already have a schema, Amazon ML will help you create one on the next page.

Datasource name

* Required

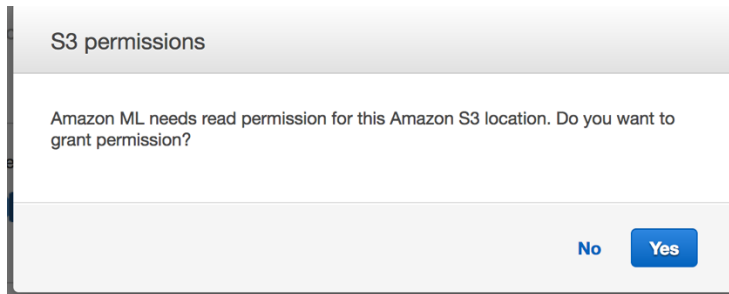
Reset

Cancel

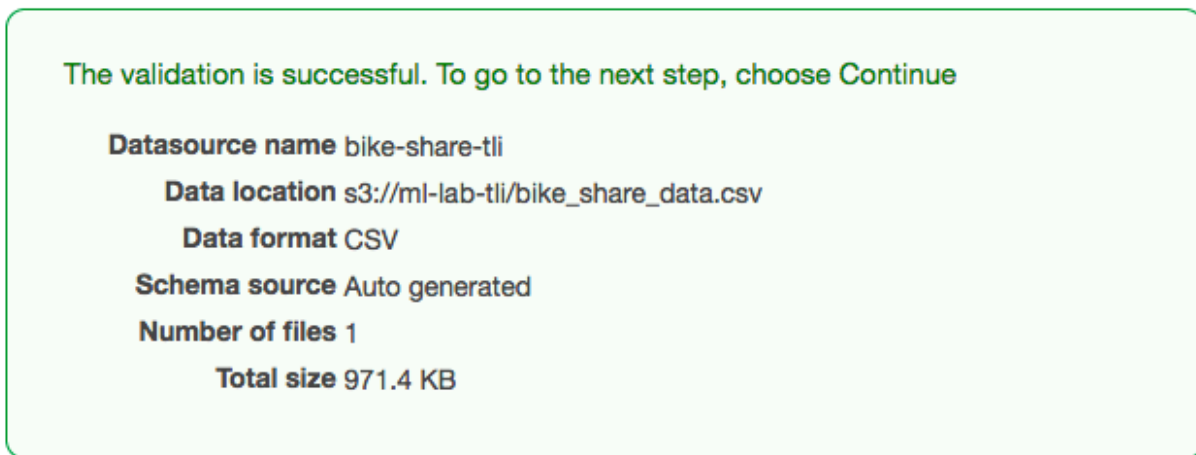
Verify

- In the Datasource name text box, input “bike-share-**<your-initials>**”. Write down the Datasource name as you will need it later to create the model.
- Click the **Verify** button. When prompted, give Amazon ML access to your S3 bucket by clicking **Yes**.

Getting Started with Amazon Machine Learning



8. You should see a message indicating validation is successful. Click the **Continue** button to move to **Schema** page.



9. Select the **Yes** radio button in the “**Does the first line in your CSV contain the column names?**” section. The page will refresh and display column names.

Schema



Amazon ML scanned your input data and inferred the column names and data type for each of the columns in your dataset. Review and edit the data type for each column to ensure that it accurately represents the data. This enables Amazon ML to read the input data correctly and to produce accurate predictions. [Learn more.](#)

Does the first line in your CSV contain the column names? ☒ Yes ☐ No ⓘ

ACTION: Change type ▾

Search by attribute name		Items per page: 10 ▾ << < 1 - 10 of 10 > >>			
<input type="checkbox"/>	Name	Data type	Sample field value 1	Sample field value 2	Sample field value 3
<input type="checkbox"/>	1 datetime	Text ▾	2011-02-08 04:00:00	2012-11-06 00:00:00	2012-05-15 06:00:00

10. Amazon Machine Learning will automatically classify the data it detects into one of four categories - Text, Categorical, Binary, and Numeric. Review the inferred data types and verify they are valid. For example, **season**, **weekday**, and **weather** are encoded as numbers, hence Amazon Machine Learning will infer as such. While Amazon Machine Learning does a good job on handling the data type even when the scheme is not completely accurate; however, if you know the data type, it is best to define it as such. In the **season**, **weekday**, and **weather** column, change **Data type** from **Numeric** to **Categorical**. With numeric data type, order and distance have typically have a meaning, where categorical data type does not.

Getting Started with Amazon Machine Learning

Q Search by attribute name Items per page: 10 « < 1 - 10 of 11 > »

<input type="checkbox"/>	▲	Name	Data type	Sample field value 1	Sample field value 2	Sample field value 3
<input type="checkbox"/>	1	season	Categorical	2	2	2
<input type="checkbox"/>	2	holiday	Binary		0	0
<input type="checkbox"/>	3	weekday	Categorical		5	4
<input type="checkbox"/>	4	workingday	Binary		1	1
<input type="checkbox"/>	5	weather	Categorical	1	3	1
<input type="checkbox"/>	6	temp	Numeric	0.7	0.54	0.7
<input type="checkbox"/>	7	atemp	Numeric	0.6364	0.5152	0.6515
<input type="checkbox"/>	8	humidity	Numeric	0.45	0.52	0.58
<input type="checkbox"/>	9	windspeed	Numeric	0.1642	0.3582	0.1045
<input type="checkbox"/>	10	count	Numeric	830	470	289

« < 1 - 10 of 11 > »

11.

12. Make sure the schemas definition matches the following screenshot then click on the right arrow at the bottom of the table.

Q Search by attribute name Items per page: 10 « < 1 - 10 of 11 > »

<input type="checkbox"/>	▲	Name	Data type	Sample field value 1	Sample field value 2	Sample field value 3
<input type="checkbox"/>	1	season	Categorical	2	2	2
<input type="checkbox"/>	2	holiday	Binary	0	0	0
<input type="checkbox"/>	3	weekday	Categorical	5	5	4
<input type="checkbox"/>	4	workingday	Binary	1	1	1
<input type="checkbox"/>	5	weather	Categorical	1	3	1
<input type="checkbox"/>	6	temp	Numeric	0.7	0.54	0.7
<input type="checkbox"/>	7	atemp	Numeric	0.6364	0.5152	0.6515
<input type="checkbox"/>	8	humidity	Numeric	0.45	0.52	0.58
<input type="checkbox"/>	9	windspeed	Numeric	0.1642	0.3582	0.1045
<input type="checkbox"/>	10	count	Numeric	830	470	289

« < 1 - 10 of 11 > »

13. Make sure the schemas definition matches the following screenshot and note that **datetime** is set to **Text** as the data type. Click **Continue** move to the **Target** page.

Getting Started with Amazon Machine Learning

Search by attribute name Items per page: 10 « < 11 - 11 of 11 > »

<input type="checkbox"/>	▲	Name	▼	Data type	▼	Sample field value 1	Sample field value 2	Sample field value 3
<input type="checkbox"/>		11		datetime	Text ▼	2012-03-23 18:00:00	2011-05-06 18:00:00	2012-05-24 12:00:00

« < 11 - 11 of 11 > »

[Cancel](#)
[Previous](#)
[Continue](#)

14. Select the **count** column as the target. This will forecast the bike rentals based on other variables such as weather, season, temperature and etc. Click **Continue** to move to **Row ID** page.

Target	Name	▲	Data type	▼	Sample field value 1	Sample field value 2	Sample field value 3
<input type="radio"/>	atemp		Numeric		0.6364	0.5152	0.6515
<input checked="" type="radio"/>	count		Numeric		830	470	289
<input type="radio"/>	holiday		Binary		0	0	0
<input type="radio"/>	humidity		Numeric		0.45	0.52	0.58
<input type="radio"/>	season		Categorical		2	2	2
<input type="radio"/>	temp		Numeric		0.7	0.54	0.7
<input type="radio"/>	weather		Categorical		1	3	1
<input type="radio"/>	weekday		Categorical		5	5	4
<input type="radio"/>	windspeed		Numeric		0.1642	0.3582	0.1045
<input type="radio"/>	workingday		Binary		1	1	1

« < 1 - 10 of 10 > »

15. Select **No** in the “**Does your data contain an identifier?**” section. Click **Review** button to move the **Review** page.
Row identifier (optional) ?

An optional row identifier helps you understand how prediction rows correspond to observation rows from the input data. If you choose to make an attribute the row identifier, Amazon ML will add that column to the prediction output. A row identifier is intended for reference purposes only. Amazon ML does not include the row identifier when training ML models.

Does your data contain an identifier? ☐ Yes ☒ No

[Cancel](#)
[Previous](#)
[Review](#)

16. Click **Continue** to move to the **ML Model Settings** page. Leave everything as default on the page and click **Review**. Note the ML model type is **REGRESSION**, which is industry-standard learning algorithm known as Linear Regression.

Getting Started with Amazon Machine Learning

ML model settings

You can use the automatically suggested ML model settings, or you can choose to customize.

ML model type REGRESSION ⓘ

ML model target count

ML model name (Optional)

Select training and evaluation settings Recipes and training parameters control the ML model training process. You can select these settings for your ML model or use the defaults provided by Amazon ML. In either case, you can choose to have Amazon ML reserve a portion of the input data for evaluation. [Learn more.](#)

☒ **Default (Recommended)**

- Generate a default recipe
- Use default training parameters
- Set aside 30% of your training data to evaluate the training
- Split the evaluation data sequentially ⓘ

☐ **Custom**

- Modify the recipe Amazon ML generates
- Modify training parameters
- Randomly or sequentially split your evaluation data ⓘ

Evaluation Name

[Cancel](#) [Previous](#) [Review](#)

17. Review the configuration parameters then scroll down and click the **Create ML Model** button.

Advanced settings

Maximum ML model Size 100MB
 Maximum number of data p... 10
 Shuffle type for training data Auto
 Regularization type L2
 Regularization amount 1e-6 - Mild

Tags ⓘ

Amazon ML copies a maximum of 10 tags from parent objects. Edit the list to keep the tags you need.

No tags

[Cancel](#) [Previous](#) [Create ML model](#)

18. After a short while, the datasource is imported and the wizard will take you to the ML model summary page. Click on the **Amazon Machine Learning** drop down and select **Dashboard**.

Amazon Machine Learning ▾ ML models > ml-U5oLMF9CkQF

Dashboard

Datasources

ML models

Evaluations

Batch Predictions

ML model summary

ID ml-U5oLMF9CkQF

Name ML model: bike-share-tli ✎

Type Numerical regression

19. Please wait for the status changes from **Pending**, **In progress** to **Completed** on all items before continue to the next section. You can click on **Refresh** (on upper right hand corner of the table) to refresh the statuses. The process typically takes about 3 – 5 minutes.

Try Real-time Prediction

- Now that the model is trained, we are ready to try some real-time predictions. Click on the machine learning model to go to the summary page.

Filter: All types ▾		Q Object name or ID		Items per page: 10 ▾		« < 1 - 5 of 5 Objects > »	
	Name	Type	ID	Status	Creation time	Completion time	
<input type="checkbox"/>	▶ Evaluation: ML model: bike-share-tli	Evaluation	ev-x7U49RhgiJg	Completed	May 18, 2017 9:44:37 PM	3 mins.	
<input type="checkbox"/>	▶ ML model: bike-share-tli	ML model	ml-U5oLMF9CkQF	Completed	May 18, 2017 9:44:37 PM	2 mins.	
<input type="checkbox"/>	▶ bike-share-tli_percentBegin=70, percentEnd=...	Datasource	ds-PIHUIHHY8zp	Completed	May 18, 2017 9:44:36 PM	3 mins.	
<input type="checkbox"/>	▶ bike-share-tli_percentBegin=0, percentEnd=7...	Datasource	ds-PSVnKqwENVa	Completed	May 18, 2017 9:44:36 PM	5 mins.	
<input type="checkbox"/>	▶ bike-share-tli	Datasource	ds-zEiYCYvUI4	Completed	May 18, 2017 9:42:33 PM	4 mins.	

- Click on **Try real-time predictions** on the left pane, enter the values shown below, then click the right arrow at the bottom right corner.

ML model report
Summary
Settings
Monitoring
Tools
Try real-time predictions
Evaluations
▶ Evaluation: ML mod...

Try real-time predictions

You submitted 10 out of 10 data values for this prediction. ✕

Try generating real-time predictions for free using the web browser on this page. To request a real-time prediction, complete the following form or provide a single data record in CSV format. To provide a data record, choose the **Paste a record** button. [Paste a record](#)

	Name	Type	Value
1	season	Categorical	3
2	holiday	Binary	0
3	weekday	Categorical	0
4	workingday	Binary	1
5	weather	Categorical	1
6	temp	Numeric	30.34
7	atemp	Numeric	34.09
8	humidity	Numeric	58
9	windspeed	Numeric	7.0015
10	count	Numeric	Target

« < 1 - 10 of 11 > »

[Clear data](#) [Create prediction](#)

Prediction results

Target name count

ML model type NUMERIC

Predicted value 170.36264038085938

```
{
  "Prediction": {
    "details": {
      "Algorithm": "SGD",
      "PredictiveModelType": "REGRESSION"
    },
    "predictedScores": [],
    "predictedValue": 170.36264038085938
  }
}
```

Next steps

3. Enter a value in the datetime field and click **Create prediction**. Given the parameters entered, the predicted number of bikes will be rented is 170.

The screenshot displays the Amazon Machine Learning console interface for creating a prediction. On the left, a table lists attributes, with the 'datetime' attribute selected and its value '2017-09-11 13:00:00' entered in the 'Value' field. Below the table are 'Clear data' and 'Create prediction' buttons. On the right, the 'ML model type' is 'NUMERIC' and the 'Predicted value' is '170.36264038085938'. Below this, a JSON response is shown, with the 'predictedValue' field highlighted in a red box.

Name	Type	Value
11	datetime	Text

Items per page: 10 « < 11 - 11 of 11 > »

Clear data Create prediction

ML model type: NUMERIC

Predicted value: 170.36264038085938

```
{
  "Prediction": {
    "details": {
      "Algorithm": "SGD",
      "PredictiveModelType": "REGRESSION"
    },
    "predictedScores": {},
    "predictedValue": 170.36264038085938
  }
}
```

Evaluate the Model

1. Part of the machine learning process is to evaluate the generated model. The Amazon Machine Learning automatically evaluate its model's accuracy based on industry-stand 70/30 data split. Which means by default, it will split 70% of the dataset to train the model but leaves out the 30% of the data for validating the model. During the validation process, Amazon ML evaluates row by row with its prediction against the actual known value. This is a way to prevent the model from learning too much about the training data and not generalized enough to predict on data that has not seen. This phenomenon is called "over-fitting". There are many techniques such as regularization to prevent this from happening, Amazon Machine Learning offers regularization in advanced mode. However, this is beyond the scope of this tutorial.
2. To visualize the model performance, click on **Amazon Machine Learning** drop down and select **Dashboard**. Click on the machine learning model to go to the summary page.

Getting Started with Amazon Machine Learning

Filter: All types ▾

Items per page: 10 ▾

« < 1 - 5 of 5 Objects > »

Name	Type	ID	Status	Creation time	Completion time
<input type="checkbox"/> ▶ Evaluation: ML model: bike-share-tli	Evaluation	ev-x7U49Rhgiig	Completed	May 18, 2017 9:44:37 PM	3 mins.
<input type="checkbox"/> ▶ ML model: bike-share-tli	ML model	ml-U5oLMF9CkQF	Completed	May 18, 2017 9:44:37 PM	2 mins.
<input type="checkbox"/> ▶ bike-share-tli_[percentBegin=70, percentEnd=...	Datasource	ds-PIHUIHHY8zp	Completed	May 18, 2017 9:44:36 PM	3 mins.
<input type="checkbox"/> ▶ bike-share-tli_[percentBegin=0, percentEnd=7...	Datasource	ds-PSVnKqwENVa	Completed	May 18, 2017 9:44:36 PM	5 mins.
<input type="checkbox"/> ▶ bike-share-tli	Datasource	ds-zEiYYCyvUI4	Completed	May 18, 2017 9:42:33 PM	4 mins.

- Under the **Evaluation** section, click on the down arrow to expand the options and click on **Summary**.

ML model report

- Summary**
- Settings
- Monitoring

Tools

- Try real-time predictions

Evaluations

- ▼ Evaluation: ML mode
 - Summary**
 - Alerts
 - Explore performance

ML model summary

ID	ml-yilspuBe1OH
Name	ML model: ds-casual-tli
Type	Numerical regression
Creation time	May 3, 2017 1:14:17 PM
Completion time	2 mins.
Compute Time (Approximate)	1 min.
Status	Completed
Log	Download log

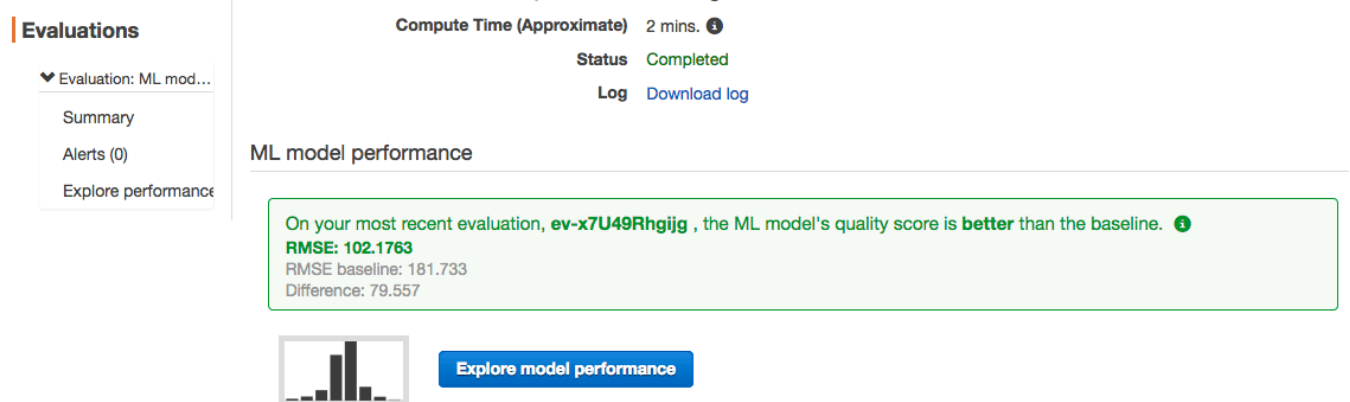
Datasource (training)

Datasource ID	ds-0Uc2Q4qGNU0
Target	casual
Input schema	View input schema

- The accuracy of an regression model is evaluated by root-mean-square-error (**RMSE**). The lower the **RMSE** is, the better quality the model is. Amazon ML provides this as a baseline metric which takes the average of the target values. In the screen shot below, the model quality score (102) is performing better than the baseline (181) with a

Getting Started with Amazon Machine Learning

difference of 80.



5. It is also possible to explore the datasource and evaluate each feature's (column) correlation to the target value (count). Features with higher value means they are more correlated to the target value, hence higher predictive power. This is very helpful when the datasource has a lot of features, you should select only the features with higher predictive power to reduce noise and improve the model's performance. This process is called **feature selection**.
6. To visualize the data correlations to target, click on **Amazon Machine Learning** drop down and select **Dashboard**. Click on one of the datasource to go to the Data insights page.

Filter: All types ▾

Q Object name or ID

Items per

Name	Type	ID	Status
<input type="checkbox"/> ▶ Evaluation: ML model: bike-share-tli	Evaluation	ev-x7U49RhgiJg	Completed
<input type="checkbox"/> ▶ ML model: bike-share-tli	ML model	ml-U5oLMF9CkQF	Completed
<input type="checkbox"/> ▶ bike-share-tli_[percentBegin=70, percentEnd=...	Datasource	ds-PIHUiHHY8zp	Completed
<input type="checkbox"/> ▶ bike-share-tli_[percentBegin=0, percentEnd=7...	Datasource	ds-PSVnKqwENVa	Completed
<input type="checkbox"/> ▶ bike-share-tli	Datasource	ds-zEiYYCvUI4	Completed

7. Under the **Attributes** section, click the **Numeric** attribute. In the screenshot below, we can see the atemp (feels like temperature) feature has a relatively higher correlation to the target. This makes sense as people may not want to rent a bike when it is too cold

or hot.

Data summary
Target distributions
Missing values
Attributes
Binary
Categorical
Numeric
Text

Numeric attributes

Attributes ^	Correlations to target *	Missing values	Invalid values	Range	Mean	Median	Preview
atemp	0.20186	0 (0%)	0 (0%)	0 - 1	0.47459574549635863	0.4848	
count	Not available	0 (0%)	0 (0%)	1 - 968	189.79513223457263	143	
humidity	0.15457	0 (0%)	0 (0%)	0 - 1	0.6287102338060557	0.64	
temp	0.19143	0 (0%)	0 (0%)	0.02 - 0.98	0.4953123802223075	0.5	
windspeed	0.022	0 (0%)	0 (0%)	0 - 0.8507	0.18819597546952888	0.1642	

* Correlations to Target is an approximate statistic for numeric attributes.

Improve Prediction Accuracy

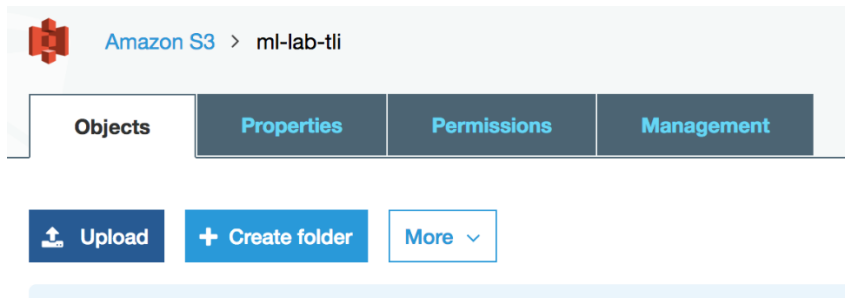
When the model accuracy is not satisfactory, you can improve it through feature engineering and/or apply expert domain knowledge to the data such as feature selection. Dataset quality is an extremely important part of the model training process, remember the rule of “Gold-In-Gold-Out”. For example, recall that the **datetime** column is text. You could break it apart into separate year, month, date and time columns so that Amazon ML could use the new information to improve its accuracy. Recall that in the Amazon QuickSight charts, usage varies throughout the day. So, hour could be a key feature to improve the model’s accuracy. Another suggestion is to add actual text description to the weather feature instead of just a categorical number. Sometimes text can provide more meaning therefore a higher predictive power.

In this section, we will train a new model with a dataset that has datetime broken out to year, month, day, and hour. The step is fairly straight forward, you can use a simple Unix awk command or Excel to extract the data. For example, we used the following awk command to extract the date and time - `awk '{gsub("-", "");gsub(" ", "");gsub(":00:00", "");print }' bike_share_data.csv > bike_share_engineered_data.csv`

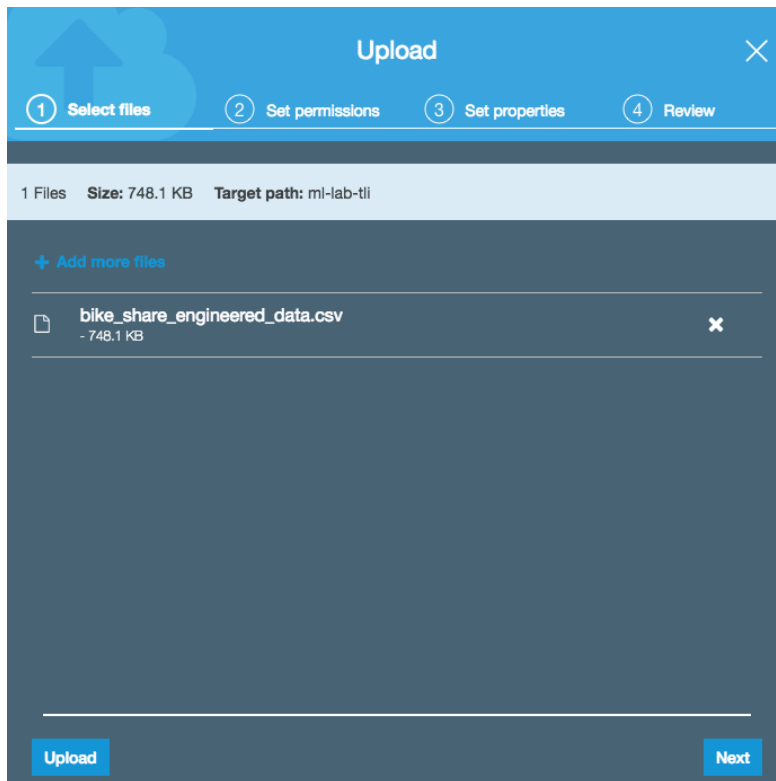
In the interest of time and keeping the focus on the Amazon Machine Learning service, we have broken out the datetime and made it available for download below.

1. Download the engineered biking sharing dataset from [here](#).
2. Sign into the AWS Management Console and open the Amazon S3 console at <https://console.aws.amazon.com/s3>
3. In the upper-right corner of the AWS Management Console, confirm you are in the desired AWS region (e.g., Oregon).
4. Click the bucket name link with “ml-lab-**<your-initials>**”, then click **Upload**.

Getting Started with Amazon Machine Learning



5. Click **Add Files**, find and select the **bike_share_engineered_data.csv** file and click **Upload**.



20. Open the Amazon Machine Learning console at <https://console.aws.amazon.com/machinelearning/home>
21. Click on **Create New...** button and select **"Datasource and ML Model"**

Getting Started with Amazon Machine Learning

Amazon Machine Learning (Amazon ML) can solve business problems by finding and learning the patterns in your historical data and using the patterns to generate predictions. To get started, you provide Amazon ML with your data. Next, you use Amazon ML to train your ML model, and then you evaluate the model's performance. Finally, you use the model to generate predictions on new data.

[Learn how to use Amazon Machine Learning by walking through the Amazon ML tutorial.](#)

22. Under “**Where is your data?**” section, make sure the S3 icon is selected. In the S3 location text box, enter “ml-lab-**<your-initials>**” for the bucket name. Auto-complete should prompt a list of matches. Select the one that you created. Auto-complete will then prompt you with list of files in the bucket, select **bike_share_engineered_data.csv** from the list.

Input data

Import your data to create an Amazon ML datasource. Amazon ML can use your datasource to create and evaluate an ML model, and you can use the data.

Where is your data?



Amazon Redshift

S3 data access

Tell Amazon ML how to access your data and give it permission to access it.

S3 location *

s3:// ml-lab-ttl/

Enter the
[Learn more](#)

ml-lab-ttl/bike_share_data.csv

grant Amazon ML permission to read this data.

If you already

ml-lab-ttl/bike_share_engineered_data.csv

s3://<path-of-input-data>.schema. If you don't have a schema, Amazon ML will help you create one on the next page. [?](#)

Datasource name

* Required

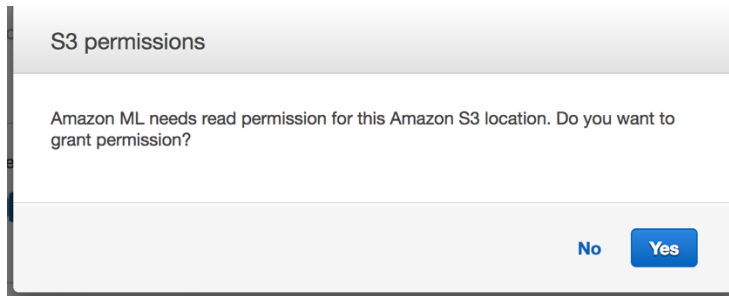
Reset

Cancel

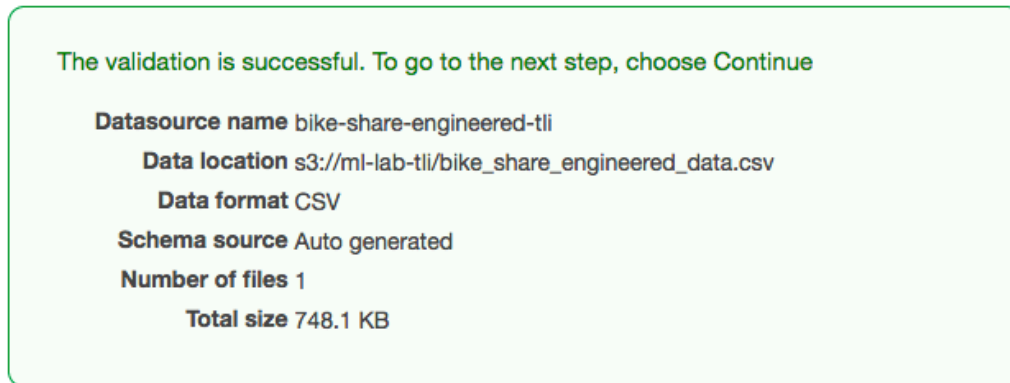
Verify

23. In the Datasource name text box, enter “bike-share-engineered-**<your-initials>**”. Write down the Datasource name as you will need it later to create the model.
24. Click the **Verify** button. If prompted, give Amazon ML access to your S3 bucket by clicking **Yes**.

Getting Started with Amazon Machine Learning



25. You should see a message indicating validation is successful. Click the **Continue** button to move to **Schema** page.



26. Select the **Yes** radio button in the “**Does the first line in your CSV contain the column names?**” section. The page will refresh and display column names.

Schema



Amazon ML scanned your input data and inferred the column names and data type for each of the columns in your dataset. Review and edit the data type for each column to ensure that it accurately represents the data. This enables Amazon ML to read the input data correctly and to produce accurate predictions. [Learn more.](#)

Does the first line in your CSV contain the column names? ☒ Yes ☐ No ⓘ

ACTION: Change type ▾

Q Search by attribute name		Items per page: 10 ▾ << < 1 - 10 of 10 > >>			
<input type="checkbox"/>	Name	Data type	Sample field value 1	Sample field value 2	Sample field value 3
<input type="checkbox"/>	1	datetime	2011-02-08 04:00:00	2012-11-06 00:00:00	2012-05-15 06:00:00

27. In the **season**, **weekday**, **weather**, **month**, and **hour** column, change **Data type** from **Numeric** to **Categorical**. In the **year** column, change **Binary** to **Numeric**. Make sure the schemas definition matches the following screenshot then click on the right arrow at the bottom right of the table.

Getting Started with Amazon Machine Learning

Q Search by attribute name				Items per page: 10 « 1 - 10 of 13 »		
<input type="checkbox"/>	▲	Name	Data type	Sample field value 1	Sample field value 2	Sample field value 3
<input type="checkbox"/>	1	season	Categorical	1	3	2
<input type="checkbox"/>	2	year	Numeric	0	1	1
<input type="checkbox"/>	3	month	Categorical	2	7	4
<input type="checkbox"/>	4	hour	Categorical	2	21	5
<input type="checkbox"/>	5	holiday	Binary	0	0	0
<input type="checkbox"/>	6	weekday	Categorical	1	2	4
<input type="checkbox"/>	7	workingday	Binary	1	1	1
<input type="checkbox"/>	8	weather	Categorical	1	1	1
<input type="checkbox"/>	9	temp	Numeric	0.2	0.74	0.3
<input type="checkbox"/>	10	atemp	Numeric	0.2576	0.697	0.2879

28. Make sure the schemas definition matches the following screenshot and click **Continue** move to the **Target** page.

Q Search by attribute name				Items per page: 10 « 11 - 13 of 13 »		
<input type="checkbox"/>	▲	Name	Data type	Sample field value 1	Sample field value 2	Sample field value 3
<input type="checkbox"/>	11	humidity	Numeric	0.8	0.66	0.61
<input type="checkbox"/>	12	windspeed	Numeric	0	0.2239	0.2239
<input type="checkbox"/>	13	count	Numeric	3	378	21

Cancel Previous **Continue**

29. Select the **count** column as the target. This will forecast the bike rentals based on other variables such as weather, season, temperature and etc. Click **Continue** to move to the **Row ID** page.

Q Search by attribute name				Items per page: 10 « 1 - 10 of 13 »		
Target	▲	Name	Data type	Sample field value 1	Sample field value 2	Sample field value 3
<input type="radio"/>		atemp	Numeric	0.2576	0.697	0.2879
<input checked="" type="radio"/>		count	Numeric	3	378	21
<input type="radio"/>		holiday	Binary	0	0	0
<input type="radio"/>		hour	Categorical	2	21	5

30. Select **No** in the “Does your data contain an identifier?” section. Click **Review** button to move the **Review** page.

Getting Started with Amazon Machine Learning

Row identifier (optional) ?

An optional row identifier helps you understand how prediction rows correspond to observation rows from the input data. If you choose to make an attribute the row identifier, Amazon ML will add that column to the prediction output. A row identifier is intended for reference purposes only. Amazon ML does not include the row identifier when training ML models.

Does your data contain an identifier? ☐ Yes ☒ No

[Cancel](#)[Previous](#)[Review](#)

31. Click **Continue** to move to the **ML Model Settings** page. Leave everything as default on the page and click **Review**.

ML model settings

You can use the automatically suggested ML model settings, or you can choose to customize.

ML model type REGRESSION ⓘ

ML model target count

ML model name (Optional)
ML model: bike-share-ttl

Select training and evaluation settings Recipes and training parameters control the ML model training process. You can select these settings for your ML model or use the defaults provided by Amazon ML. In either case, you can choose to have Amazon ML reserve a portion of the input data for evaluation. [Learn more](#).

☒ **Default (Recommended)**

- Generate a default recipe
- Use default training parameters
- Set aside 30% of your training data to evaluate the training
- Split the evaluation data sequentially ⓘ

☐ **Custom**

- Modify the recipe Amazon ML generates
- Modify training parameters
- Randomly or sequentially split your evaluation data ⓘ

Evaluation Name Evaluation: ML model: bike-share-ttl

[Cancel](#)[Previous](#)[Review](#)

32. Review the configuration parameters then scroll down and click the **Create ML Model** button.

Advanced settings

Maximum ML model Size	100MB
Maximum number of data p...	10
Shuffle type for training data	Auto
Regularization type	L2
Regularization amount	1e-6 - Mild

Tags ⓘ

Amazon ML copies a maximum of 10 tags from parent objects. Edit the list to keep the tags you need.

No tags

[Cancel](#)[Previous](#)[Create ML model](#)

33. After a short while, the datasource is imported and the wizard will take you to the ML model summary page. Click on the **Amazon Machine Learning** drop down and select **Dashboard**.

Getting Started with Amazon Machine Learning

Amazon Machine Learning ▾ ML models > ml-U5oLMF9CkQF

Dashboard

Datasources

ML models

Evaluations

Batch Predictions

ML model summary

ID ml-U5oLMF9CkQF

Name ML model: bike-share-tli ✎

Type Numerical regression

34. Please wait for the status changes from **Pending**, **In progress** to **Completed** on all items before continue to the next step. You can click on **Refresh** (on upper right hand corner of the table) to refresh the statuses. The process typically takes about 3 – 5 minutes.
35. Let's go ahead and visualize the model performance, click on **Amazon Machine Learning** drop down and select **Dashboard**. Click on the "ML model: bike-share-engineered-<your-initials>" model to go to the summary page.

Filter: All types ▾ Items per page: 10 ▾ << < 1 - 10 of 10 Objects > >>

	Name	Type	ID	Status	Creation time	Completion time
<input type="checkbox"/>	▶ Evaluation: ML model: bike-share-engineered-tli	Evaluation	ev-dtxUIN067ZI	Completed	May 19, 2017 8:12:31 AM	4 mins.
<input type="checkbox"/>	▶ ML model: bike-share-engineered-tli	ML model	ml-JH9nx7N2VeW	Completed	May 19, 2017 8:12:31 AM	3 mins.
<input type="checkbox"/>	▶ bike-share-engineered-tli_[percentBegin=70, ...	Datasource	ds-f65iyWYMA56	Completed	May 19, 2017 8:12:31 AM	4 mins.
<input type="checkbox"/>	▶ bike-share-engineered-tli_[percentBegin=0, p...	Datasource	ds-UuQplmwNurq	Completed	May 19, 2017 8:12:30 AM	4 mins.
<input type="checkbox"/>	▶ bike-share-engineered-tli	Datasource	ds-FxZT5UrkEcX	Completed	May 19, 2017 8:10:54 AM	3 mins.

36. Under the **Evaluation** section, click on the down arrow to expand the options and click on **Summary**.

Getting Started with Amazon Machine Learning

ML model report

- Summary**
- Settings
- Monitoring

Tools

- Try real-time predictions

Evaluations

- ▼ Evaluation: ML mode
 - Summary
 - Alerts
 - Explore performance

ML model summary

ID	ml-yilspuBe1OH
Name	ML model: ds-casual-tli
Type	Numerical regression
Creation time	May 3, 2017 1:14:17 PM
Completion time	2 mins.
Compute Time (Approximate)	1 min.
Status	Completed
Log	Download log

Datasource (training)

Datasource ID	ds-0Uc2Q4qGNU0
Target	casual
Input schema	View input schema

Evaluations

37. In the screen shot below, the model quality score is 97 and baseline is 176 with difference of 78. Recall the first model with datetime not broken out had model quality score of 102 and baseline of 181 with difference of 80. The model with datetime broken out performed better overall than the model without. As demonstrated here, feature engineering is a critical component to improve predictive accuracy of a machine learning model.

▼ Evaluation: ML mod...

- Summary
- Alerts (0)
- Explore performance

status Completed
Log [Download log](#)

ML model performance

On your most recent evaluation, **ev-dtxUIN067Zi**, the ML model's quality score is **better** than the baseline.

RMSE: 97.6985
RMSE baseline: 175.985
Difference: 78.287

[Explore model performance](#)

Additional Available Models on Amazon Machine Learning Service

Recall the model created in this lab is called Regression, which is typically used to predict a numeric value such as demand for bike rental in this lab. There are two other models that Amazon ML supports, binary classification and multiclass classification. Binary classification is

an industry-standard learning algorithm known as Logistic Regression. This algorithm is great for predicting yes/no type of answers. For example, fraud detection and spam filtering. Multiclass classification is an industry-standard learning algorithm known as multinomial logistic regression. This algorithm is typically used for solving problems such as product categorization.

Clean Up

1. Delete the bucket with your initials in S3
2. Go to Amazon Machine Learning Dashboard and delete all resources with your initials

"Acknowledgement: Lichman, M. (2013). UCI Machine Learning Repository [<http://archive.ics.uci.edu/ml>]. Irvine, CA: University of California, School of Information and Computer Science."