



## SageMaker Immersion Day



### ► Prerequisites

### ▼ Lab 1. Feature Engineering

Option 1: Amazon SageMaker Data Wrangler and Feature Store

#### Option 2: Numpy and Pandas

Option 3: Amazon SageMaker Processing

Lab 2. Train, Tune and Deploy XGBoost

### ► Lab 3. Bring your own model

### ► Lab 4. Autopilot, Debugger and Model Monitor

### ► Lab 5. Bias and Explainability

### ► Lab 6. SageMaker Pipelines

Lab 7. Real Time ML inference on Streaming Data

Lab 8. Build ML Model with No Code Using Sagemaker Canvas

### ► Lab 9. Amazon SageMaker JumpStart

### ► Lab 10. ML Governance Tools for Amazon SageMaker

### ► Lab 11. SageMaker Notebook Instances

### ▼ Content preferences

Language

English ▼

[SageMaker Immersion Day](#) > [Lab 1. Feature Engineering](#) > **Option 2: Numpy and Pandas**

## Option 2: Numpy and Pandas

- [Data Preparation](#)
- [Conclusion](#)



DO NOT perform this part if you have already executed the feature engineering with Amazon Data Wrangler.

### Data Preparation

In this step you will use your Amazon SageMaker Studio notebook to preprocess the data that you need to train your machine learning model.

1. Click on this “amazon-sagemaker-immersion-day” folder and then double click on the `xgboost_direct_marketing_sagemaker.ipynb` notebook.
2. If you are prompted to choose a Kernel, choose the “Python 3 (Data Science)” kernel and click “Select”.

## SageMaker Immersion Day

### ► Prerequisites

#### ▼ Lab 1. Feature Engineering

Option 1: Amazon SageMaker Data Wrangler and Feature Store

**Option 2: Numpy and Pandas**

Option 3: Amazon SageMaker Processing

Lab 2. Train, Tune and Deploy XGBoost

#### ► Lab 3. Bring your own model

#### ► Lab 4. Autopilot, Debugger and Model Monitor

#### ► Lab 5. Bias and Explainability

#### ► Lab 6. SageMaker Pipelines

Lab 7. Real Time ML inference on Streaming Data

Lab 8. Build ML Model with No Code Using Sagemaker Canvas

#### ► Lab 9. Amazon SageMaker JumpStart

#### ► Lab 10. ML Governance Tools for Amazon SageMaker

#### ► Lab 11. SageMaker Notebook Instances

### ▼ Content preferences

Language

## Set up notebook environment

Set up environment for "xgboost\_direct\_marketing\_sagemaker.ipynb".

Image

Data Science

Kernel

Python 3

Instance type

ml.t3.medium

Start-up script ⓘ

No script

Cancel

Select

3. You will then have the notebook opened. You can verify the Kernel CPU and Memory states on the top right of the notebook.

## SageMaker Immersion Day

### ► Prerequisites

### ▼ Lab 1. Feature Engineering

Option 1: Amazon SageMaker Data Wrangler and Feature Store

#### Option 2: Numpy and Pandas

Option 3: Amazon SageMaker Processing

Lab 2. Train, Tune and Deploy XGBoost

### ► Lab 3. Bring your own model

### ► Lab 4. Autopilot, Debugger and Model Monitor

### ► Lab 5. Bias and Explainability

### ► Lab 6. SageMaker Pipelines

Lab 7. Real Time ML inference on Streaming Data

Lab 8. Build ML Model with No Code Using SageMaker Canvas

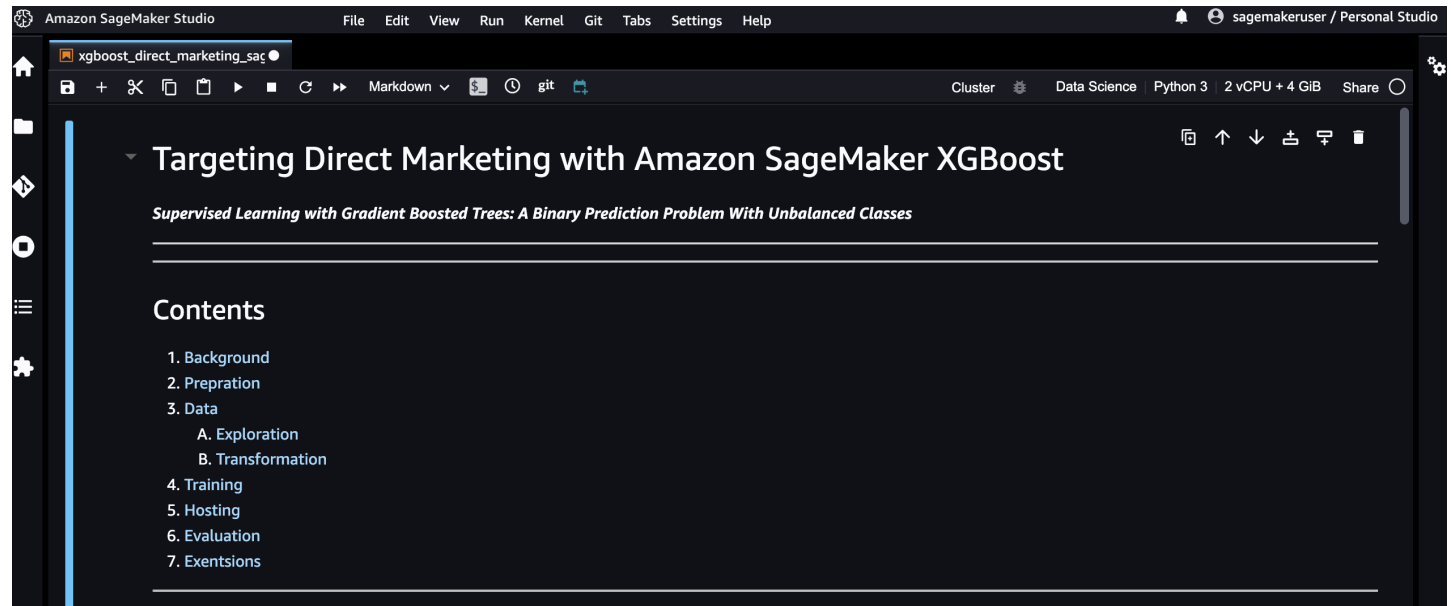
### ► Lab 9. Amazon SageMaker JumpStart

### ► Lab 10. ML Governance Tools for Amazon SageMaker

### ► Lab 11. SageMaker Notebook Instances

### ▼ Content preferences

Language



4. Execute the first four cells by pressing **Shift+Enter** in each of the cells. While the code runs, an \* appears between the square brackets as pictured in the first screenshot to the right. After a few seconds, the code execution will complete, the \* will be replaced with the number 1.

This code will import some libraries and define a few environment variables in your Jupyter notebook environment.

```
# cell 03
import sagemaker_datawrangler          # For interactive data prep widget
import numpy as np                     # For matrix operations and numerical processing
import pandas as pd                    # For munging tabular data
import matplotlib.pyplot as plt        # For charts and visualizations
from IPython.display import Image      # For displaying images in the notebook
from IPython.display import display    # For displaying outputs in the notebook
from time import gmtime, strftime      # For labeling SageMaker models, endpoints, etc.
import sys                             # For writing outputs to notebook
import math                             # For ceiling function
import json                             # For parsing hosting outputs
import os                               # For manipulating filepath names
import sagemaker                       # For writing outputs to notebook
import zipfile                          # For parsing hosting outputs
# Amazon SageMaker's Python SDK provides many helper functions
```

## SageMaker Immersion Day

### ► Prerequisites

### ▼ Lab 1. Feature Engineering

Option 1: Amazon  
SageMaker Data Wrangler  
and Feature Store

**Option 2: Numpy and  
Pandas**

Option 3: Amazon  
SageMaker Processing

Lab 2. Train, Tune and Deploy  
XGBoost

### ► Lab 3. Bring your own model

### ► Lab 4. Autopilot, Debugger and Model Monitor

### ► Lab 5. Bias and Explainability

### ► Lab 6. SageMaker Pipelines

Lab 7. Real Time ML inference  
on Streaming Data

Lab 8. Build ML Model with No  
Code Using Sagemaker Canvas

### ► Lab 9. Amazon SageMaker JumpStart

### ► Lab 10. ML Governance Tools for Amazon SageMaker

### ► Lab 11. SageMaker Notebook Instances

### ▼ Content preferences

Language

5. Now, Let's download the dataset by running the 5th cell.

## Data

Let's start by downloading the [direct marketing dataset](#) from the sample data s3 bucket.

[Moro et al., 2014] S. Moro, P. Cortez and P. Rita. A Data-Driven Approach to Predict the Success of Bank Telemarketing. Decision Support Systems, Elsevier, 62:22-31, June 2014

```
# cell 05
!wget https://sagemaker-sample-data-us-west-2.s3-us-west-2.amazonaws.com/autopilot/direct_marketing/bank-additional.zip

with zipfile.ZipFile('bank-additional.zip', 'r') as zip_ref:
    zip_ref.extractall('.')
```

6. In the next cell you will load the dataset into a pandas dataframe and utilize the `sagemaker_datawrangler` library to explore the data and understand the data distribution for each feature:

## SageMaker Immersion Day

### ► Prerequisites

#### ▼ Lab 1. Feature Engineering

Option 1: Amazon SageMaker Data Wrangler and Feature Store

#### Option 2: Numpy and Pandas

Option 3: Amazon SageMaker Processing

Lab 2. Train, Tune and Deploy XGBoost

#### ► Lab 3. Bring your own model

#### ► Lab 4. Autopilot, Debugger and Model Monitor

#### ► Lab 5. Bias and Explainability

#### ► Lab 6. SageMaker Pipelines

Lab 7. Real Time ML inference on Streaming Data

Lab 8. Build ML Model with No Code Using Sagemaker Canvas

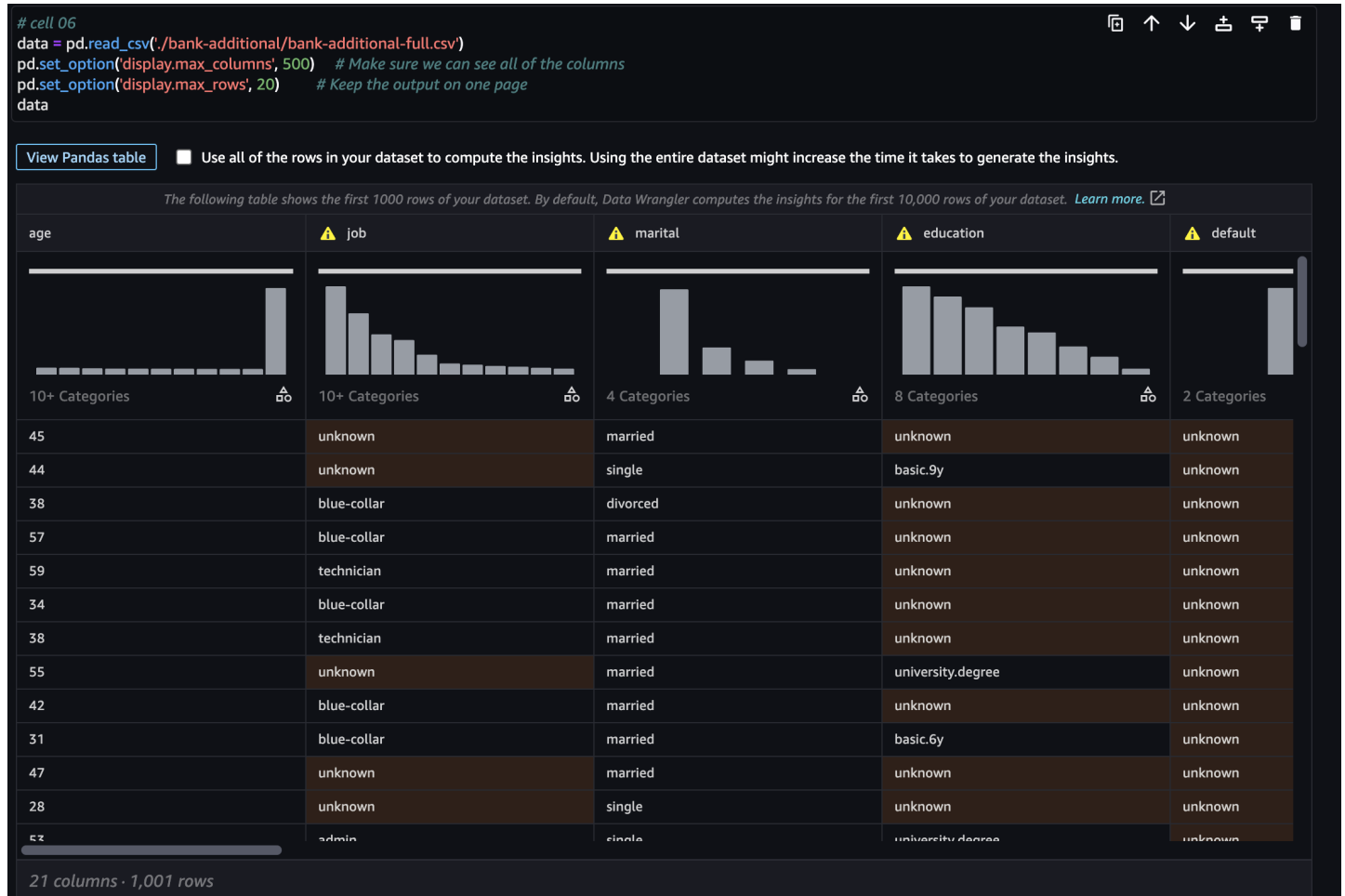
#### ► Lab 9. Amazon SageMaker JumpStart

#### ► Lab 10. ML Governance Tools for Amazon SageMaker

#### ► Lab 11. SageMaker Notebook Instances

### ▼ Content preferences

Language



7. Now let's look at how our features relate to the target that we are attempting to predict.

```
[6]: for column in data.select_dtypes(include=['object']).columns:
      if column != 'y':
          display(pd.crosstab(index=data[column], columns=data['y'], normalize='columns'))

      for column in data.select_dtypes(exclude=['object']).columns:
          print(column)
          hist = data[[column, 'y']].hist(by='y', bins=30)
          plt.show()
```

8. Now let's look at how our features relate to one another

## SageMaker Immersion Day

### ► Prerequisites

### ▼ Lab 1. Feature Engineering

Option 1: Amazon SageMaker Data Wrangler and Feature Store

**Option 2: Numpy and Pandas**

Option 3: Amazon SageMaker Processing

Lab 2. Train, Tune and Deploy XGBoost

### ► Lab 3. Bring your own model

### ► Lab 4. Autopilot, Debugger and Model Monitor

### ► Lab 5. Bias and Explainability

### ► Lab 6. SageMaker Pipelines

Lab 7. Real Time ML inference on Streaming Data

Lab 8. Build ML Model with No Code Using Sagemaker Canvas

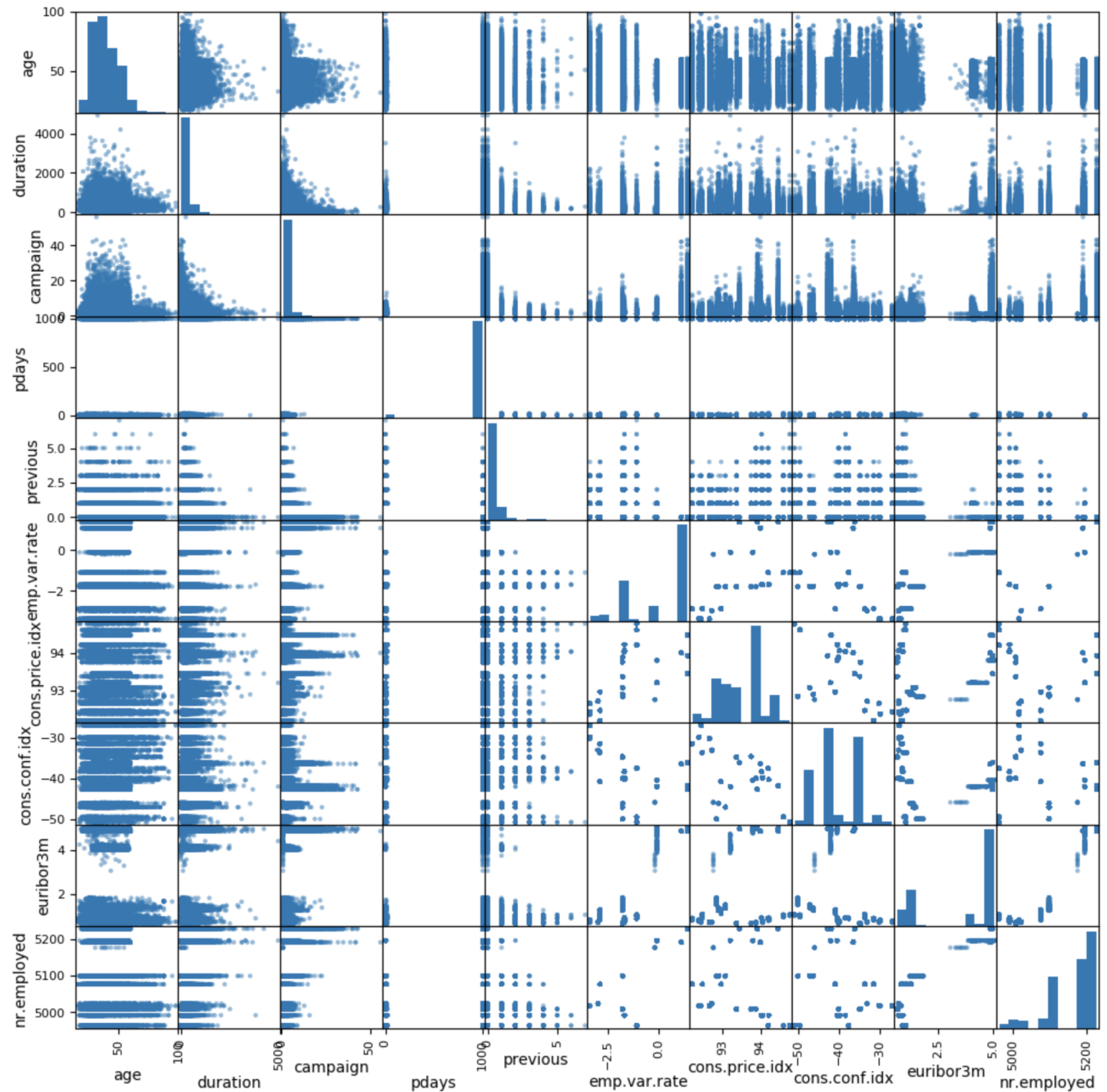
### ► Lab 9. Amazon SageMaker JumpStart

### ► Lab 10. ML Governance Tools for Amazon SageMaker

### ► Lab 11. SageMaker Notebook Instances

### ▼ Content preferences

Language



## SageMaker Immersion Day

### ► Prerequisites

### ▼ Lab 1. Feature Engineering

Option 1: Amazon  
SageMaker Data Wrangler  
and Feature Store

**Option 2: Numpy and  
Pandas**

Option 3: Amazon  
SageMaker Processing

Lab 2. Train, Tune and Deploy  
XGBoost

### ► Lab 3. Bring your own model

### ► Lab 4. Autopilot, Debugger and Model Monitor

### ► Lab 5. Bias and Explainability

### ► Lab 6. SageMaker Pipelines

Lab 7. Real Time ML inference  
on Streaming Data

Lab 8. Build ML Model with No  
Code Using Sagemaker Canvas

### ► Lab 9. Amazon SageMaker JumpStart

### ► Lab 10. ML Governance Tools for Amazon SageMaker

### ► Lab 11. SageMaker Notebook Instances

### ▼ Content preferences

Language

9. Now we will convert all the categorical variables using one hot encoding

```
# cell 09

# Note: These transformations can be done through the graphical widget that we generated above. The data prep widget will automatically generate code for transformations that you do.
data[no_previous_contact] = np.where(data[pdays] == 999, 1, 0) # Indicator variable to capture when pdays takes a value of 999
data[not_working] = np.where(np.in1d(data[job], ['student', 'retired', 'unemployed']), 1, 0) # Indicator for individuals not actively employed
```

```
# cell 10

model_data = pd.get_dummies(data) # Convert categorical variables to sets of indicators
```

10. We will drop few features which are not required

```
# cell 11

model_data = model_data.drop(['duration', 'emp.var.rate', 'cons.price.idx', 'cons.conf.idx', 'euribor3m', 'nr.employed'], axis=1)
```

11. We will split our data set into 3 channels: train, test, validation set:

```
# cell 12

train_data, validation_data, test_data = np.split(model_data.sample(frac=1, random_state=1729), [int(0.7 * len(model_data)), int(0.9 * len(model_data))]) # Randomly sort the data then split
```

12. Amazon SageMaker XGBoost algorithm expects data to be in libSVM or CSV format (without header) and the first column must be the target variable. So in this step we transform the data accordingly

```
# cell 13

pd.concat([train_data['y_yes'], train_data.drop(['y_no', 'y_yes'], axis=1)], axis=1).to_csv('train.csv', index=False, header=False)
pd.concat([validation_data['y_yes'], validation_data.drop(['y_no', 'y_yes'], axis=1)], axis=1).to_csv('validation.csv', index=False, header=False)
```

13. Now we will upload the final data into Amazon S3 bucket.

```
# cell 14

boto3.Session().resource('s3').Bucket(bucket).Object(os.path.join(prefix, 'train/train.csv')).upload_file('train.csv')
boto3.Session().resource('s3').Bucket(bucket).Object(os.path.join(prefix, 'validation/validation.csv')).upload_file('validation.csv')
```

14. Now check the S3 bucket through the console to make sure you have uploaded the train.csv and validation.csv. In the AWS Console Main Page, type "S3".

## SageMaker Immersion Day

### ► Prerequisites

#### ▼ Lab 1. Feature Engineering

Option 1: Amazon  
SageMaker Data Wrangler  
and Feature Store

**Option 2: Numpy and  
Pandas**

Option 3: Amazon  
SageMaker Processing

Lab 2. Train, Tune and Deploy  
XGBoost

#### ► Lab 3. Bring your own model

#### ► Lab 4. Autopilot, Debugger and Model Monitor

#### ► Lab 5. Bias and Explainability

#### ► Lab 6. SageMaker Pipelines

Lab 7. Real Time ML inference  
on Streaming Data

Lab 8. Build ML Model with No  
Code Using SageMaker Canvas

#### ► Lab 9. Amazon SageMaker JumpStart

#### ► Lab 10. ML Governance Tools for Amazon SageMaker

#### ► Lab 11. SageMaker Notebook Instances

### ▼ Content preferences

Language

# AWS Management Console

## AWS services

### Find Services

You can enter names, keywords or acronyms.

Q s3

S3

Scalable Storage in the Cloud

15. Click on your bucket name. If two buckets are present, take the one including your region name ("eu-west-1" as an example in the image below).

Buckets (2)					↺	Copy ARN	Empty	Delete	Create bucket
Buckets are containers for data stored in S3. <a href="#">Learn more</a>					Find buckets by name				
	Name	Region	Access	Creation date					
<input type="radio"/>	sagemaker-eu-west-1-██████████	EU (Ireland) eu-west-1	Objects can be public	January 30, 2021, 13:10:54 (UTC+01:00)					
<input type="radio"/>	sagemaker-studio-██████████	EU (Ireland) eu-west-1	Objects can be public	January 30, 2021, 12:55:04 (UTC+01:00)					

16. Click on "sagemaker/" and then on "DEMO-xgboost-dm/". You will see a "train/" and "validation/" folder.



## SageMaker Immersion Day

### ► Prerequisites

### ▼ Lab 1. Feature Engineering

Option 1: Amazon SageMaker Data Wrangler and Feature Store

#### Option 2: Numpy and Pandas

Option 3: Amazon SageMaker Processing

Lab 2. Train, Tune and Deploy XGBoost

### ► Lab 3. Bring your own model

### ► Lab 4. Autopilot, Debugger and Model Monitor

### ► Lab 5. Bias and Explainability

### ► Lab 6. SageMaker Pipelines

Lab 7. Real Time ML inference on Streaming Data

Lab 8. Build ML Model with No Code Using SageMaker Canvas

### ► Lab 9. Amazon SageMaker JumpStart

### ► Lab 10. ML Governance Tools for Amazon SageMaker

### ► Lab 11. SageMaker Notebook Instances

### ▼ Content preferences

Language

Amazon S3 > sagemaker-studio- > sagemaker/ > DEMO-xgboost-dm/

### DEMO-xgboost-dm/

**Folder overview**

Region EU (Ireland) eu-west-1	S3 URI s3://sagemaker-studio- /sagemaker/DEMO-xgboost-dm/	Amazon resource name (ARN) arn:aws:s3::sagemaker-studio- /sagemaker/DEMO-xgboost-dm/
----------------------------------	--	---

Drag and drop files and folders you want to upload here, or choose **Upload**.

**Objects (2)** Refresh Delete Actions Create folder Upload

Objects are the fundamental entities stored in Amazon S3. For others to access your objects, you'll need to explicitly grant them permissions. [Learn more](#)

< 1 > ⚙

<input type="checkbox"/>	Name ▲	Type ▼	Last modified ▼	Size ▼	Storage class ▼
<input type="checkbox"/>	train/	Folder	-	-	-
<input type="checkbox"/>	validation/	Folder	-	-	-

17. You can have a look inside of each folder to make sure that the “train.csv” file is there as well as the “validation.csv” file.

## SageMaker Immersion Day

### ► Prerequisites

### ▼ Lab 1. Feature Engineering

Option 1: Amazon SageMaker Data Wrangler and Feature Store

### Option 2: Numpy and Pandas

Option 3: Amazon SageMaker Processing

Lab 2. Train, Tune and Deploy XGBoost

### ► Lab 3. Bring your own model

### ► Lab 4. Autopilot, Debugger and Model Monitor

### ► Lab 5. Bias and Explainability

### ► Lab 6. SageMaker Pipelines

Lab 7. Real Time ML inference on Streaming Data

Lab 8. Build ML Model with No Code Using SageMaker Canvas

### ► Lab 9. Amazon SageMaker JumpStart

### ► Lab 10. ML Governance Tools for Amazon SageMaker

### ► Lab 11. SageMaker Notebook Instances

### ▼ Content preferences

Language

Amazon S3 > sagemaker-studio- > sagemaker/ > DEMO-xgboost-dm/ > train/

## train/

### Folder overview

Region  
EU (Ireland) eu-west-1

S3 URI  
s3://sagemaker-studio- /sagemaker/DEMO-xgboost-dm/train/

Amazon resource name (ARN)  
arn:aws:s3::sagemaker-studio- /sagemaker/DEMO-xgboost-dm/train/

Drag and drop files and folders you want to upload here, or choose **Upload**.

### Objects (1)

Objects are the fundamental entities stored in Amazon S3. For others to access your objects, you'll need to explicitly grant them permissions. [Learn more](#)

Find objects by prefix

<input type="checkbox"/>	Name	Type	Last modified	Size	Storage class
<input type="checkbox"/>	train.csv	csv	October 29, 2020, 14:34 (UTC+01:00)	3.4 MB	Standard

Amazon S3 > sagemaker-studio- > sagemaker/ > DEMO-xgboost-dm/ > validation/

## validation/

### Folder overview

Region  
EU (Ireland) eu-west-1

S3 URI  
s3://sagemaker-studio- /sagemaker/DEMO-xgboost-dm/validation/

Amazon resource name (ARN)  
arn:aws:s3::sagemaker-studio- /sagemaker/DEMO-xgboost-dm/validation/

Drag and drop files and folders you want to upload here, or choose **Upload**.

### Objects (1)

Objects are the fundamental entities stored in Amazon S3. For others to access your objects, you'll need to explicitly grant them permissions. [Learn more](#)

Find objects by prefix

<input type="checkbox"/>	Name	Type	Last modified	Size	Storage class
<input type="checkbox"/>	validation.csv	csv	October 29, 2020, 14:34 (UTC+01:00)	989.2 KB	Standard

**Congratulations!!** You have successfully prepared the data to train an XGBoost model.

## Conclusion

In this lab you have walked through the process of required environment setup and data engineering to clean and prepare your data for model building and training. In the next lab you will learn how to Train , Tune and Deploy an XGBoost model using SageMaker's Built-in XGBoost algorithm.

Previous

Next

## SageMaker Immersion Day



### ► Prerequisites

### ▼ Lab 1. Feature Engineering

Option 1: Amazon  
SageMaker Data Wrangler  
and Feature Store

**Option 2: Numpy and  
Pandas**

Option 3: Amazon  
SageMaker Processing

Lab 2. Train, Tune and Deploy  
XGBoost

### ► Lab 3. Bring your own model

### ► Lab 4. Autopilot, Debugger and Model Monitor

### ► Lab 5. Bias and Explainability

### ► Lab 6. SageMaker Pipelines

Lab 7. Real Time ML inference  
on Streaming Data

Lab 8. Build ML Model with No  
Code Using Sagemaker Canvas

### ► Lab 9. Amazon SageMaker JumpStart

### ► Lab 10. ML Governance Tools for Amazon SageMaker

### ► Lab 11. SageMaker Notebook Instances

### ▼ Content preferences

Language

