

Lab 3.1 - Amazon SageMaker - Creating and importing data

The screenshot shows the AWS SageMaker console interface. At the top, there's a navigation bar with 'aws' and a search bar. Below it, the path 'Amazon SageMaker AI > Notebook instances > Create notebook instance' is visible. The main area is titled 'Create notebook instance'. It contains two main sections: 'Notebook instance settings' and 'Permissions and encryption'. In the 'Notebook instance settings' section, the 'Notebook instance name' is set to 'MyNotebook', 'Notebook instance type' is 'ml.m4.xlarge', 'Platform identifier' is 'Amazon Linux 2, Jupiter Lab 4', and 'Lifecycle configuration' is 'ml-pipeline-c195280a5013002113473376t1w19621122561'. Under 'Additional configuration', the 'Volume size in GB' is set to '5', and the 'Minimum IMDS Version' is '2'. In the 'Permissions and encryption' section, the 'IAM role' dropdown is open, showing 'AmazonSageMakerFullAccess'. Below this, a Jupyter notebook session titled 'PythonCheatSheet.ipynb' is running. The session shows a code cell generating a DataFrame and its head, followed by a plot command. A line chart is displayed with four series labeled A, B, C, and D, showing data from January 2000 to July 2002.

```
[#8]: # On a DataFrame, the plot() method is convenient to plot all of the columns with labels
df4 = pd.DataFrame(np.random.randn(1000, 4), index=ts.index,columns=['A','B','C','D'])
df4.cumsum()
df4.head()

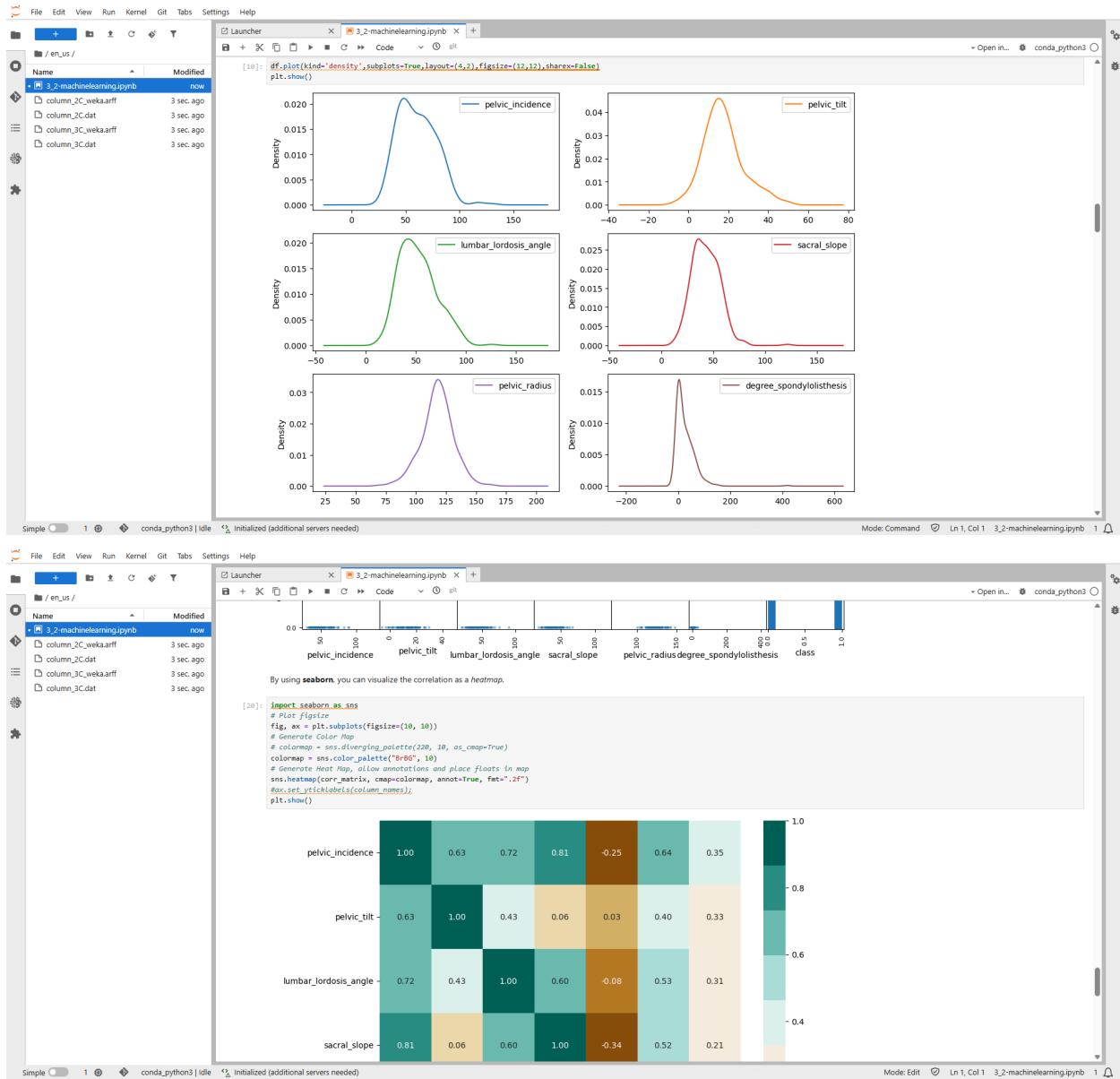
[#9]: df4.plot()
plt.show()
```

2000-01-01 -0.608307 -0.619764 -0.260578 0.765093
2000-01-02 0.352232 -0.694344 -0.784845 1.220520
2000-01-03 0.105910 0.127102 -2.813004 3.441426
2000-01-04 1.461577 -1.350820 -2.698823 -4.474760
2000-01-05 2.760711 -2.984324 -0.278835 3.113061

Jan 2000 Jul Jan Jul Jan Jul

The screenshot shows two Jupyter Notebook interfaces side-by-side. The left notebook, titled 'AMAZON SAGEMAKER SAMPLE NOTEBOOKS', displays a list of files and code snippets related to Amazon SageMaker, including 'Amazon_Tabular_Classification.ipynb', 'Amazon_Tabular_Classification_TabTransformer.ipynb', 'Amazon_Tabular_Regression.ipynb', 'Amazon_Tabular_Regression_TabTransformer.ipynb', 'Amazon_Tensorflow_Object_Detection.ipynb', 'automatic-speech-recognition.ipynb', 'blazingtext_word2vec_text.ipynb', 'code-quality-evaluate-human-eval.ipynb', 'custom_dog_image_generator.ipynb', 'deepar_exmaple.ipynb', 'DeepAR_Electricity.ipynb', 'domain_adaptation-finetuning.ipynb', 'factcheck_mnist_mnist.ipynb', 'falcon_causal_ipynb', 'gemma-fine-tuning.ipynb', 'image-generation-stable-diffusion.ipynb', 'language-models-tutorial.ipynb', 'lamb_jit_start_end_inference.ipynb', 'ljumpstart_text-generation-inference.ipynb', 'k_nearest_neighbors_cvtype.ipynb', 'LDA-Introduction.ipynb', 'linear_learner_mnist.ipynb', 'Linear_learner_Regression.ipynb', 'format.ipynb', 'llama-2-chat-composition.ipynb', 'llama-2-finetuning.ipynb', 'llama-2-text-completion.ipynb', 'llama-3-training-inference-finetuning-deployment.ipynb', 'llama-3-finetuning.ipynb', 'llama-3-text-completion.ipynb', 'llama-guard-text-moderation.ipynb', 'mistral-7b-instruction-domain-adaptation-finetuning.ipynb', 'mm_synthethic.ipynb', 'object-detection_birds.ipynb', 'project2vec_sentence_similarity.ipynb', 'pca_mnist.ipynb', 'question_answering_CoherenLangchainJumpstart.ipynb', 'question_answering_LangchainJumpstart.ipynb', 'question_answering_mm_synthethic_embedding.ipynb', 'llama-2-jump-random-cut_forest.ipynb', 'SageMaker-Seq2Seq-Language-English.ipynb', 'semantic_segmentation_pascalvoc.ipynb', 'textEmbedding-seed-selection.ipynb', 'text-generation-att2pt.ipynb', 'text-generation-falcon.ipynb', 'text-generation-few-shot-learning.ipynb'. The right notebook, titled 'An Introduction to Linear Learner with MNIST', displays code for training a linear model on the MNIST dataset, including imports like 'import warnings', 'import requests', 'import zipfile', 'import pandas as pd', 'from sklearn import linear_model', and code for reading a zip file containing 'column_2C_weka.arff' and loading it into a DataFrame. Both notebooks are running in a 'conda_python3' environment.

Lab 3.2 - Amazon SageMaker - Exploring Data



Lab 3.3 - Amazon SageMaker - Encoding Categorical Data

Step 1: Importing and exploring the data

You will start by examining the data in the dataset.

To get the most out of this lab, read the instructions and code before you run the cells. Take time to experiment!

Start by importing the pandas package and setting some default display options.

```
[1]: import pandas as pd
pd.set_option('display.max_rows', 500)
pd.set_option('display.max_columns', 500)
pd.set_option('display.width', 1000)
```

Next, load the dataset into a pandas DataFrame.

The data doesn't contain a header, so you will define those column names in a variable that's named `col_names` to the attributes listed in the dataset description.

```
[2]: url = "imports-85.csv"
col_names = ['symboling', 'normalized-losses', 'fuel-type', 'aspiration', 'num-of-doors', 'body-style', 'drive-wheels', 'engine-location', 'wheel-base', 'length', 'width', 'height', 'curb-weight', 'engine-type', 'num-of-cylinders', 'engine-size', 'fuel-system', 'bore', 'stroke', 'compression-ratio', 'horsepower', 'peak-rpm', 'city-mpg', 'highway-mpg', 'price']
```

```
[2]: df_car = pd.read_csv(url, sep=',', names=col_names, na_values='?', header=None)
```

First, see to see the number of rows (instances) and columns (features), you will use `.shape`.

```
[3]: df_car.shape
```

```
[3]: (205, 25)
```

Next, examine the data by using the `.head` method.

```
[4]: df_car.head(5)
```

```
[4]:
```

	symboling	fuel-type	aspiration	num-of-doors	body-style	drive-wheels	engine-location	engine-type	num-of-cylinders	fuel-system	price	horsepower	city-mpg
0	3	gas	std	two	convertible	rwd	front	dohc	four	mpfi	13495.0	111.0	21
1	3	gas	std	two	convertible	rwd	front	dohc	four	mpfi	16500.0	111.0	21
2	1	gas	std	two	hatchback	rwd	front	ohcv	six	mpfi	16500.0	154.0	19
3	2	gas	std	four	sedan	fwd	front	ohc	four	mpfi	13950.0	102.0	24
4	2	gas	std	four	sedan	4wd	front	ohc	five	mpfi	17450.0	115.0	18

```
[32]: df_car = pd.get_dummies(df_car, columns=['body-style'], drop_first=False)
df_car = pd.get_dummies(df_car, columns=['engine-location'], drop_first=True)
df_car = pd.get_dummies(df_car, columns=['engine-type'], drop_first=False)
df_car = pd.get_dummies(df_car, columns=['fuel-system'], drop_first=False)
df_car = pd.get_dummies(df_car, columns=['fuel-type'], drop_first=True)

print("Final shape: (df_car.shape)")
df_car.head()
```

```
[33]:
```

	symboling	fuel-type	aspiration	num-of-doors	body-style	drive-wheels	engine-location	engine-type	num-of-cylinders	fuel-system	price	horsepower	city-mpg
0	3	gas	std	two	convertible	rwd	front	dohc	four	mpfi	13495.0	111.0	21
1	3	gas	std	two	convertible	rwd	front	dohc	four	mpfi	16500.0	111.0	21
2	1	gas	std	two	hatchback	rwd	front	ohcv	six	mpfi	16500.0	154.0	19
3	2	gas	std	four	sedan	fwd	front	ohc	four	mpfi	13950.0	102.0	24
4	2	gas	std	four	sedan	4wd	front	ohc	five	mpfi	17450.0	115.0	18

```
[33]: Final shape: (205, 30)
```

File Edit View Run Kernel Git Tabs Settings Help

3.3-machinelearning.ipynb

```
[33]: df_car = pd.get_dummies(df_car, columns=['body-style'], drop_first=False)
df_car = pd.get_dummies(df_car, columns=['engine-location'], drop_first=True)
df_car = pd.get_dummies(df_car, columns=['engine-type'], drop_first=False)
df_car = pd.get_dummies(df_car, columns=['fuel-system'], drop_first=False)
df_car = pd.get_dummies(df_car, columns=['fuel-type'], drop_first=True)

print("Final shape: (df_car.shape)")

Final shape: (205, 30)
```

body-style	engine-location_rear	engine-type_dohc	engine-type_dohcf	engine-type_i	engine-type_ohc	engine-type_ohcf	engine-type_ohcv	engine-type_rotor	fuel-system_1bbl	fuel-system_2bbl	fuel-system_4bbl	fuel-system_idi	fuel-system_mf1	fuel-system_mpfi	fuel-system_spdi	fuel-system_spfi	fuel-type_gas
False	False	True	False	False	False	False	False	False	False	False	False	False	True	False	False	True	True
False	False	True	False	False	False	False	True	False	False	False	False	False	True	False	False	True	True
False	False	False	False	False	True	False	False	False	False	False	False	False	True	False	False	True	True
False	False	False	False	False	True	False	False	False	False	False	False	False	True	False	False	True	True

You now have four columns. These columns all contain text values.

```
[34]: # df_car.head()
```

Most machine learning algorithms require inputs that are numerical values.

- The `num-of-cylinders` and `num-of-doors` features have an ordinal value. You could convert the values of these features into their numerical counterparts.
- However, `aspiration` and `drive-wheels` don't have an ordinal value. These features must be converted differently.

You will explore the ordinal features first.

Step 2: Encoding ordinal features

```
File Edit View Run Kernel Git Tabs Settings Help
```

3.3-machinelearning.ipynb

drive-wheels_fwd
drive-wheels_rwd

The encoding was straightforward. If the value in the `drive-wheels` column is `4wd`, then a `1` is the value in the `drive-wheels_4wd` column. A `0` is the value for the other columns that were generated. If the value in the `drive-wheels` column is `fwd`, then a `1` is the value in the `drive-wheels_fwd` column, and so on.

These binary features enable you to express the information in a numerical way, without implying any order.

Examine the final column that you will encode.

The data in the `aspiration` column only has two values: `std` and `turbo`. You could encode this column into two binary features. However, you could also ignore the `std` value and record whether it's `turbo` or not. To do this, you would still use the `get_dummies` method but specify `drop_first=True`.

```
[46]: df_car['aspiration'].value_counts()
```

```
[46]: aspiration
std    168
turbo   37
Name: count, dtype: int64
```

```
[47]: df_car = pd.get_dummies(df_car,columns=['aspiration'], drop_first=True)
```

```
[48]: df_car.head()
```

engine-type_i	engine-type_ohc	engine-type_ohcf	engine-type_ohcv	engine-type_rotor	fuel-system_1bbl	fuel-system_2bbl	fuel-system_4bbl	fuel-system_idi	fuel-system_mf1	fuel-system_mpfi	fuel-system_spdi	fuel-system_spfi	fuel-type_gas	doors	cylinders	drive-wheels_4wd	drive-wheels_fwd
False	False	False	False	False	False	False	False	False	True	False	False	True	2.0	4	False	False	
False	False	False	False	False	False	False	False	False	True	False	False	True	2.0	4	False	False	
False	False	False	True	False	False	False	False	False	True	False	True	True	2.0	6	False	False	
False	True	False	False	False	False	False	False	False	True	False	False	True	4.0	4	False	True	
False	True	False	False	False	False	False	False	True	False	False	True	True	4.0	5	True	False	

Challenge task: Go back to the beginning of this lab, and add other columns to the dataset. How would you encode the values of each column? Update the code to include some of the other features.

Congratulations!

Lab 3.4 - Amazon SageMaker - Training a model

Step 1: Exploring the data

```
[1]: import warnings, requests, zipfile, io
warnings.filterwarnings('ignore')
import pandas as pd
from scipy.io import arff
import boto3

[2]: f_zip = "http://archive.ics.uci.edu/ml/machine-learning-databases/00212/vertebral_column.data.zip"
r = requests.get(f_zip, stream=True)
Vertebral_zip = zipfile.ZipFile(io.BytesIO(r.content))
Vertebral_zip.extractall()

[3]: data = arff.loadarff('column_2c_weka.arff')
df = pd.DataFrame(data[0])

[4]: class_mapper = {b'Abnormal':1,b'Normal':0}
df['class']=df['class'].replace(class_mapper)
```

You will start with a quick reminder of the data in the dataset.

To get the most out of this lab, carefully read the instructions and code before you run the cells. Take time to experiment!

First, use `shape` to examine the number of rows and columns.

```
[5]: df.shape
```

```
[5]: (316, 7)
```

Next, get a list of the columns.

```
[6]: df.columns
```

```
[6]: Index(['pelvic_incidence', 'pelvic_tilt', 'lumbar_lordosis_angle',
       'sacral_slope', 'pelvic_radius', 'degree_spondylolisthesis', 'class'],
       dtype='object')
```

You can see the six biomechanical features, and that the target column is named `class`.

```
[7]: df.info()
```

```
[7]: <pandas.core.frame.DataFrame: columns=[u'pelvic_incidence', u'pelvic_tilt', u'lumbar_lordosis_angle', u'sacral_slope', u'pelvic_radius', u'degree_spondylolisthesis', u'class'], dtypes=[u'int64', u'int64', u'int64', u'int64', u'int64', u'int64', u'category'], index=RangeIndex(start=0, stop=316, step=1), memory_usage=1000>
```

The estimator needs `channels` to feed data into the model. For training, the `train_channel` and `validate_channel` will be used.

```
[18]: train_channel = sagemaker.inputs.TrainingInput(
    "s3://[1][1]/train".format(bucket,prefix),
    content_type='text/csv')

validate_channel = sagemaker.inputs.TrainingInput(
    "s3://[1][1]/validate".format(bucket,prefix,validate_file),
    content_type='text/csv')

data_channels = {'train': train_channel, 'validation': validate_channel}

Running fit will train the model.
```

Note: This process can take up to 5 minutes.

```
[19]: xgb_model.fit(inputs=data_channels, log=False)
```

```
INFO:sagemaker:Creating training-job with name: sagemaker-xgboost-2026-01-19-22-40-39-136
2026-01-19 22:40:40 Starting - Starting the training job...
2026-01-19 22:40:54 Starting - Preparing the instances for training...
2026-01-19 22:41:18 Downloading - Downloading input data...
2026-01-19 22:42:59 Training - Training image download completed. Training in progress...
2026-01-19 22:43:09 Uploading - Uploading generated training model.
2026-01-19 22:43:22 Completed - Training job completed
```

After the training is complete, you are ready to test and evaluate the model. However, you will do testing and validation in later labs.

Lab 3.5 - Amazon SageMaker - Deploying a model

The image shows two side-by-side Jupyter Notebook interfaces. Both notebooks are titled '3_5-machinelearning.ipynb' and are running in a 'conda_python3' kernel.

Top Notebook:

- Cell 3:** Contains Python code for downloading a dataset from an S3 bucket and preparing it for training. It uses the 'arff' library to load the dataset and 'train_test_split' from 'sklearn.model_selection' to split the data.

```

[3]: bucket='c195288a5813010113522242tlw91022884490-labbucket-zsrwlj27sua2'
[3]: import warnings, requests, zipfile, io
[3]: warnings.simplefilter('ignore')
[3]: import pandas as pd
[3]: from scipy.io import arff

[3]: import os
[3]: import boto3
[3]: import sagemaker
[3]: from sagemaker.image_uris import retrieve
[3]: from sklearn.model_selection import train_test_split
[3]: sagemaker.config INFO - Not applying SDK defaults from location: /etc/dg/sagemaker/config.yaml
[3]: sagemaker.config INFO - Not applying SDK defaults from location: /home/ec2-user/.config/sagemaker/config.yaml

[3]: f_zip = 'http://archive.ics.uci.edu/ml/machine-learning-databases/00212/vertebral_column_data.zip'
[3]: r = requests.get(f_zip, stream=True)
[3]: Vertebral_zip = zipfile.ZipFile(io.BytesIO(r.content))
[3]: Vertebral_zip.extractall()

[3]: data = arff.loadarff('column_2C_weka.arff')
[3]: df = pd.DataFrame(data[0])

[3]: class_mapper = {b'Abnormal':1,b'Normal':0}
[3]: df['class']=df['class'].replace(class_mapper)

[3]: cols = df.columns.tolist()
[3]: cols = cols[1:] + cols[:1]
[3]: df = df[cols]

[3]: train, test_and_validate = train_test_split(df, test_size=0.2, random_state=42, stratify=df['class'])
[3]: test, validate = train_test_split(test_and_validate, test_size=0.5, random_state=42, stratify=test_and_validate['class'])

[3]: prefix='lab3'

[3]: train_file='vertebral_train.csv'
[3]: test_file='vertebral_test.csv'

```

Bottom Notebook:

- Cell 18:** Contains Python code for defining a binary conversion function and printing the first 10 rows of the 'target_predicted' and 'test' datasets.

```

[18]: def binary_convert(x):
[18]:     threshold = 0.65
[18]:     if x > threshold:
[18]:         return 1
[18]:     else:
[18]:         return 0

[18]: def binary_convert(x):
[18]:     threshold = 0.5
[18]:     if x > threshold:
[18]:         return 1
[18]:     else:
[18]:         return 0

[18]: target_predicted['binary'] = target_predicted['class'].apply(binary_convert)

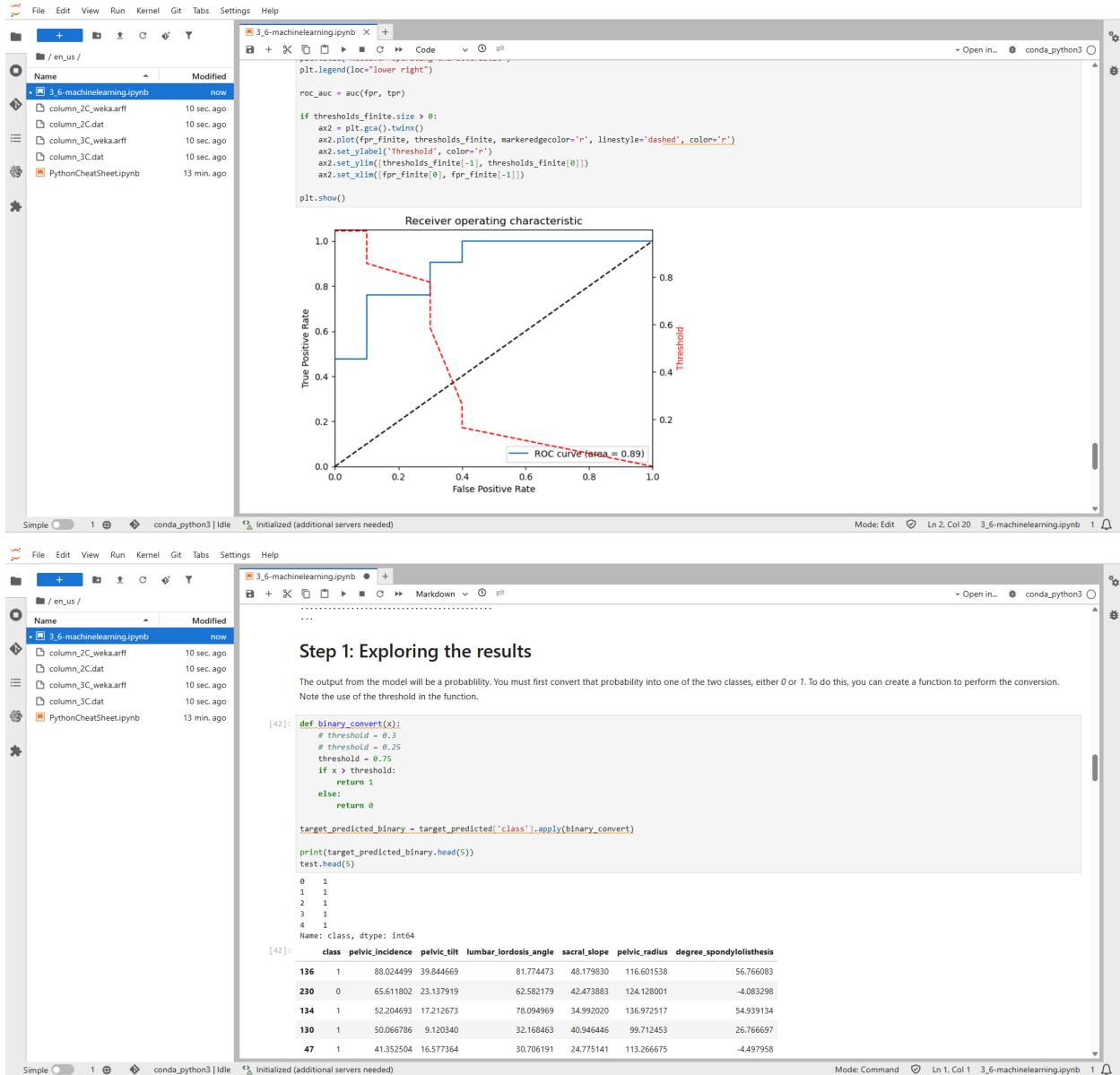
[18]: print(target_predicted.head(10))
[18]: test.head(10)

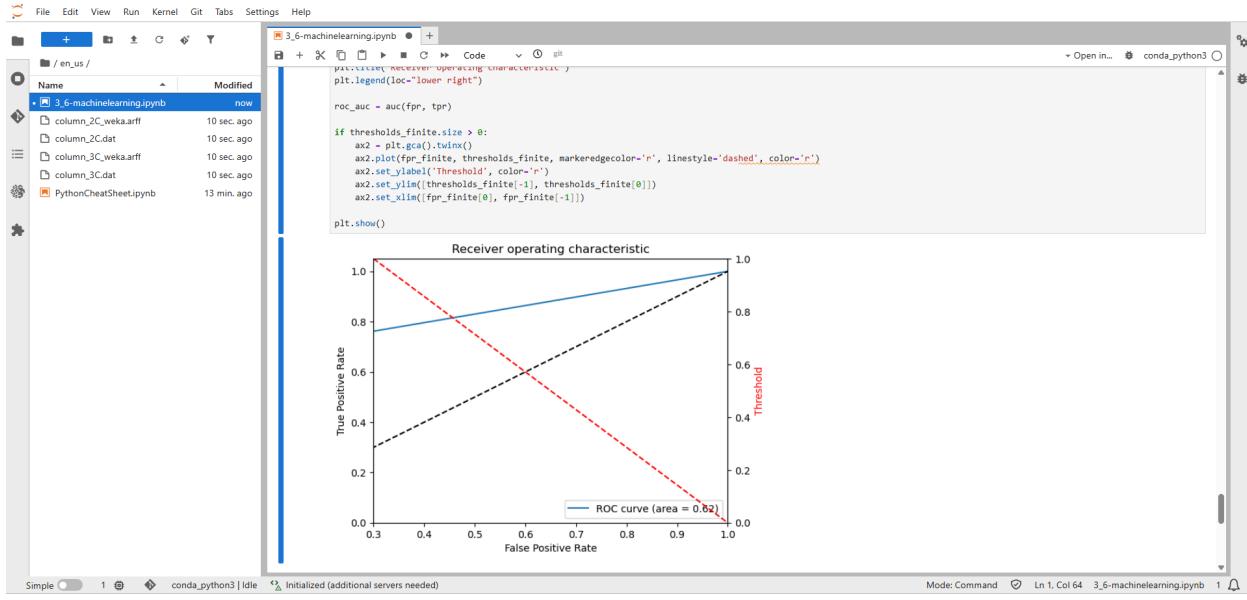
```

	class	binary
0	0.996607	1
1	0.777283	1
2	0.994641	1
3	0.993690	1
4	0.939139	1
5	0.997390	1
6	0.997717	1
7	0.987518	1
8	0.993334	1
9	0.682776	1

	class	pelvic_incidence	pelvic_tilt	lumbar_lordosis_angle	sacral_slope	pelvic_radius	degree_spondylolisthesis
136	1	88.024499	39.844669	81.774473	48.179830	116.601538	56.766083
230	0	65.611802	23.137919	62.582179	42.473883	124.128001	-4.083298
134	1	52.204693	17.212673	78.094969	34.992020	136.972517	54.939134

Lab 3.6 - Amazon SageMaker - Generating model performance metrics





Lab 3.7 - Amazon SageMaker - Hyperparameter Tuning

The screenshot shows a Jupyter Notebook interface with a code cell containing Python code for hyperparameter tuning using Amazon SageMaker. The output of the cell shows logs from the training job, indicating the creation of training and transform jobs, and the upload of generated training models.

```


```

batch_output = "s3://{}//batch-out/".format(bucket,prefix)
batch_input = "s3://{}//batch-in/{}".format(bucket,prefix,batch_X_file)

xgb_transformer = xgb_model.transformer(instance_count=1,
 instance_type='ml.m4.xlarge',
 strategy='MultiRecord',
 assemble_with='Line',
 output_path=batch_output)

xgb_transformer.transform(data_batch_input,
 data_type='S3Prefix',
 content_type='text/csv',
 split_type='Line')
xgb_transformer.wait(logs=False)

sagemaker.config.INFO - Not applying Sdk defaults from location: /etc/sdg/sagemaker/config.yaml
sagemaker.config.INFO - Not applying Sdk defaults from location: /home/ec2-user/.config/sagemaker/config.yaml
sagemaker.config.INFO - Not applying Sdk defaults from location: /etc/sdg/sagemaker/config.yaml
sagemaker.config.INFO - Not applying Sdk defaults from location: /home/ec2-user/.config/sagemaker/config.yaml
INFO:sagemaker:Creating training-job with name: sagemaker-xgboost-2026-01-21-04-32-30-374
2026-01-21 04:32:31 Starting - Starting the training job...
2026-01-21 04:32:45 Starting - Preparing the instances for training...
2026-01-21 04:33:08 Downloading - Downloading input data...
2026-01-21 04:33:23 Downloading - Downloading the training image.....
2026-01-21 04:34:14 Training - Training job has been completed. Training in progress....
2026-01-21 04:34:24 Uploading - Uploading generated training model...
2026-01-21 04:34:37 Completed - Training job completed
INFO:sagemaker:Creating model with name: sagemaker-xgboost-2026-01-21-04-34-42-118
INFO:sagemaker:Creating transform job with name: sagemaker-xgboost-2026-01-21-04-34-42-667
.....
CPU times: user 1.72 s, sys: 130 ms, total: 1.85 s
Wall time: 7min 41s

```


```

Step 1: Getting model statistics

Before you tune the model, re-familiarize yourself with the current model's metrics.

The setup performed a batch prediction, so you must read in the results from Amazon Simple Storage Service (Amazon S3).

File Edit View Run Kernel Git Tabs Settings Help

3_7-machinelearning.ipynb

Plot the confusion matrix and the receiver operating characteristic (ROC) curve for the original model.

```
[5]: plot_confusion_matrix(test_labels, target_predicted_binary)
```

Confusion Matrix

Predicted Class	0	1
True Class	7.00	3.00
0	2.00	19.00

```
[6]: plot_roc(test_labels, target_predicted_binary)
```

Sensitivity or TPR: 50.476159476159485
Specificity or FPR: 99.52384052384051
Precision: 66.36363636363636
Negative Predictive Value: 77.77777777777798
False Positive Rate: 30.0%
False Negative Rate: 49.523840523840524%
False Discovery Rate: 11.636363636363635%
Accuracy: 83.870907742915498
User-defined metrics initialized (additional servers needed)

Mode: Command Ln 8, Col 17 3_7-machinelearning.ipynb 1

File Edit View Run Kernel Git Tabs Settings Help

3_7-machinelearning.ipynb

```
[12]: %%time
batch_output = "s3://{}//{}//batch-out/".format(bucket,prefix)
batch_input = "s3://{}//{}//batch-in/{}".format(bucket,prefix,batch_x_file)

xgb_transformer = best_algo_model.transformer(instance_count=1,
                                             instance_type='ml.m4.xlarge',
                                             strategy="MultiRecord",
                                             assemble_with='Line',
                                             output_path=batch_output)

xgb_transformer.transform(data=batch_input,
                         data_type='S3Prefix',
                         content_type='text/csv',
                         split_type='Line')
xgb_transformer.wait(logs=False)

INFO:sagemaker:Creating model with name: sagemaker-xgboost-2026-01-21-04-52-14-892
INFO:sagemaker:Creating transform job with name: sagemaker-xgboost-2026-01-21-04-52-15-421
...
CPU times: user 670 ms, sys: 54.1 ms, total: 725 ms
Wall time: 6min 26s
```

Get the predicted target and the test labels of the model.

```
[13]: s3 = boto3.client('s3')
obj = s3.get_object(Bucket=bucket, Key="{}//batch-out/{}".format(prefix,'batch-in.csv.out'))
best_target_predicted = pd.read_csv(io.BytesIO(obj['Body']),read(),names=['class'])

def binary_convert(x):
    threshold = 0.5
    if x > threshold:
        return 1
    else:
        return 0

best_target_predicted_binary = best_target_predicted[['class']].apply(binary_convert)
test_labels = test.iloc[:,0]
```

Plot a confusion matrix for your best_target_predicted and test_labels.

Mode: Command Ln 8, Col 17 3_7-machinelearning.ipynb 1

