

SageMaker Hosting a Model

Steps to Host a Machine Learning Model on Sagemaker	1
1. Create SageMakerInvokeEndPoint policy:	1
2. Create ml_user_predict user:	1
3. Create Notebook Instance:	2
4. Create Bucket:	2
5. Download the Datasets needed for cleaning:	3
6. Upload the Datasets needed for cleaning:	3
7. Clean Datasets needed to train the model:	3
8. Train XGBoost Model & create Endpoint Configuration:	3
9. Examine your Endpoint, once the previous step is complete:	4
10. Make predictions on the XGBoost Model:	5
11. Create an IAM role for your Lambda function:	5
12. Create a Lambda function:	5
13. Create an API Gateway to the endpoint:	7
14. Deploy the API Gateway:	7
15. Test on Postman	8
16. DELETE THE ENDPOINT to prevent charges	9
17. Recreate the endpoint for the future	9
Adding Dependencies (Libraries) To Your Endpoint	10
18. Adding Dependencies to your Lambda Function:	10

Shout Out

The steps and example notebooks are pieces of and inspired by this Udemy course here: <https://www.udemy.com/course/practical-aws-sagemaker-6-real-world-case-studies/>. This course is worth the money if you're familiar with modeling and are interested in learning how to host on AWS.

Steps to Host a Machine Learning Model on Sagemaker

1. Create **SageMakerInvokeEndPoint** policy:¹
 - a. AWS Console > IAM > Policies > Create Policy.
 - b. Service = SageMaker
 - c. Actions = Read
 - d. Resources = All Resources
 - e. Name = "SageMakerInvokeEndPoint"

¹ Section 1: 9. Lab - Configure IAM Users, Setup Command Line Interface (CLI)

- f. No tags
- g. Create Policy

2. Create **ml_user_predict** user:²

- a. AWS Console > IAM > Users > Add User > User name = “ml_user_predict”
- b. Access Type = Programmatic Access
- c. Attach existing policies directly > SageMakerInvokeEndPoint > Check box
- d. Attach existing policies directly > AmazonS3ReadOnlyAccess > Check box
- e. No tags.
- f. Create user
- g. Download .csv > Open .csv on your local machine using any text editor or IDE.

3. Create **Notebook Instance**:³

- a. AWS Console > SageMaker
- b. Notebook > Notebook Instances > Create Notebook Instance
- c. Notebook Instance Name = “xgboost-hosted-model”
- d. Permissions and encryption > IAM Role > Drop down and select **Create A New Role**

Permissions and encryption

IAM role
Notebook instances require permissions to call other services including SageMaker and S3. Choose a role or let us create a role with the [AmazonSageMakerFullAccess](#) IAM policy attached.

AmazonSageMakerServiceCatalogProductsUseRole ▲

Create a new role

Enter a custom IAM role ARN Create a new role

Use existing role

Create an IAM role ✕

Passing an IAM role gives Amazon SageMaker permission to perform actions in other AWS services on your behalf. Creating a role here will grant permissions described by the [AmazonSageMakerFullAccess](#) IAM policy to the role you create.

The IAM role you create will provide access to:

✓ S3 buckets you specify - optional

☒ Any S3 bucket
Allow users that have access to your notebook instance access to any bucket and its contents in your account.

☐ Specific S3 buckets
Example: bucket-name-1, bucket-name-2
Comma delimited. ARNs, ** and */ are not supported.

☐ None

✓ Any S3 bucket with "sagemaker" in the name

✓ Any S3 object with "sagemaker" in the name

✓ Any S3 object with the tag "sagemaker" and value "true" [See Object tagging](#)

✓ S3 bucket with a Bucket Policy allowing access to SageMaker [See S3 bucket policies](#)

Cancel Create role

- e. Create role

² Section 1: 9. Lab - Configure IAM Users, Setup Command Line Interface (CLI)

³ Section 2: 16. Lab - S3 Bucket Setup

- f. Create Notebook Instance

4. Create **Bucket**:⁴

- a. AWS Console > S3 > Buckets > Create Bucket
- b. Bucket name = “xgboost-hosted-model-bucket”
- c. Create Bucket

5. Download the **Datasets** needed for cleaning:⁵

- a. Go to the competition
 - <https://www.kaggle.com/c/bike-sharing-demand/data>
- b. Login > Rules > Accept Competition
- c. Data > Download All

6. Upload the **Datasets** needed for cleaning:

- a. Unzip the two datasets onto your local machine:
 - You should see both a test.csv and train.csv
- b. AWS Console > Amazon SageMaker > Notebook > Notebook Instances
- c. Click “Open Jupyter” for the “xgboost-hosted-model” instance
- d. Upload the two files test.csv and train.csv



e.

7. Clean **Datasets** needed to train the model:⁶

- a. AWS Console > Amazon SageMaker > Notebook > Notebook Instances
- b. Click “Open Jupyter” for the “xgboost-hosted-model” instance
- c. Upload the **bikerental_data_preparation_rev1.ipynb** notebook
 - Udemy course resource directory:
/AmazonSageMakerCourse-master/xgboost/BikeSharingRegression/bikerental_data_preparation_rev1.ipynb
 - [Google Drive Link to notebook](#)
- d. Open the notebook
- e. Cell > Run All
- f. Close the notebook tab

⁴ Section 2: 16. Lab - S3 Bucket Setup

⁵ Section 7: 61. Lab - Data Preparation Bike Rental Regression

⁶ Section 7: 61. Lab - Data Preparation Bike Rental Regression

8. Train **XGBoost Model** & create **Endpoint Configuration**:⁷

- AWS Console > Amazon SageMaker > Notebook > Notebook Instances
- Click “Open Jupyter” for the “xgboost-hosted-model” instance
- Upload the **xgboost_cloud_training_template.ipynb** notebook
 - From the Udemy course, confirm that this is NOT the iris classification notebook of the same name.
 - Udemy course resource directory:
/AmazonSageMakerCourse-master/xgboost/BikeSharingRegression/sdk1.7/xgboost_cloud_training_template.ipynb
 - [Google Drive Link to notebook](#)
- Open the notebook
- Cell > Run All
- Close the notebook tab
- NOTE:**

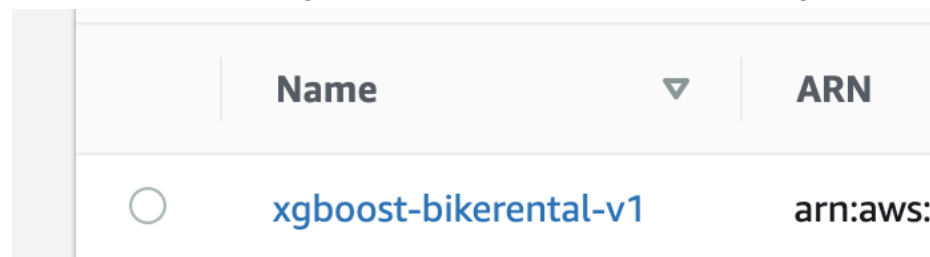
- If you run this notebook more than once, this cell:

```
In [35]: # Ref: http://sagemaker.readthedocs.io/en/latest/estimators.html
predictor = estimator.deploy(initial_instance_count=1,
                             instance_type='ml.m4.xlarge',
                             endpoint_name = 'xgboost-bikerental-v1')
```

- Will give the error:

```
ClientError: An error occurred (ValidationException) when calling the CreateEndpointConfig operation: Cannot create a
lready existing endpoint configuration "arn:aws:sagemaker:us-east-1:237397516361:endpoint-config/xgboost-bikerental-v
1".
```

- This is because you’ve already made the Endpoint Configuration, even though the Endpoint itself isn’t yet running. To see this go to: AWS Console > Amazon SageMaker > Inference > Endpoint Configuration



	Name	ARN
<input type="radio"/>	xgboost-bikerental-v1	arn:aws:

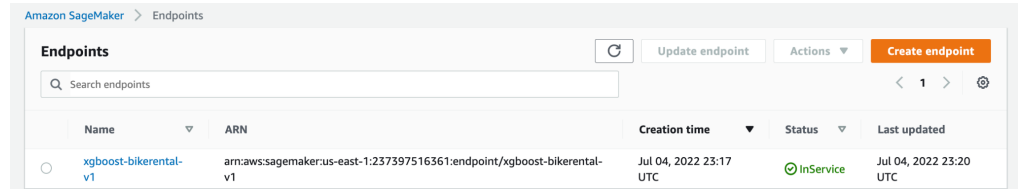
- To actually run the endpoint, you need to create a new endpoint using that configuration. To do this go to: AWS Console > Amazon SageMaker > Inference > Endpoints.
- Create Endpoint > Endpoint name = “xgboost-bikerental-v1” (Or, I use the convention of the same name as the endpoint configuration.)
- Use An Existing Configuration
- Check the box to the “xgboost-bikerental-v1” configuration > Select endpoint configuration”
- Create Endpoint

⁷ Section 7: 65. Lab - How to train using SageMaker’s built-in XGBoost Algorithm

- Never run that cell again. You can now delete the cell. Continue running the remainder of the notebook or just don't run this notebook again. It's served its purpose.

9. Examine your **Endpoint**, once the previous step is complete:

- AWS Console > Amazon SageMaker > Inference > Endpoints



- Leave the endpoint running for now but delete it when you're done with this project (which is the final step of this tutorial).

10. Make predictions on the **XGBoost Model**:⁸

- NOTE** - You don't actually need to run this step to make a hosted model but it's useful to know how to hit a deployed endpoint using a notebook.
- AWS Console > Amazon SageMaker > Notebook > Notebook Instances
- Click "Open Jupyter" for the "xgboost-hosted-model" instance
- Upload the **xgboost_cloud_prediction_template.ipynb** notebook
 - From the Udemy course, confirm that this is NOT the iris classification notebook of the same name.
 - Udemy course resource directory:
/AmazonSageMakerCourse-master/xgboost/BikeSharingRegression/sdk1.7/xgboost_cloud_prediction_template.ipynb
 - [Google Drive Link to notebook](#)
- Open the notebook tab
- Cell > Run All
- Close the notebook tab

11. Create an IAM role for your **Lambda function**:⁹

- AWS Console > IAM > Roles > Create Role.
- AWS Service > Lambda > Next.
- Assign the permission "SageMakerInvokeEndPoint" > Check the box
- Assign the permission "AWSLambdaBasicExecutionRole" > Check the box.
- Role name = "lambda_sagemaker_invoke_endpoint"

12. Create a **Lambda function**:¹⁰

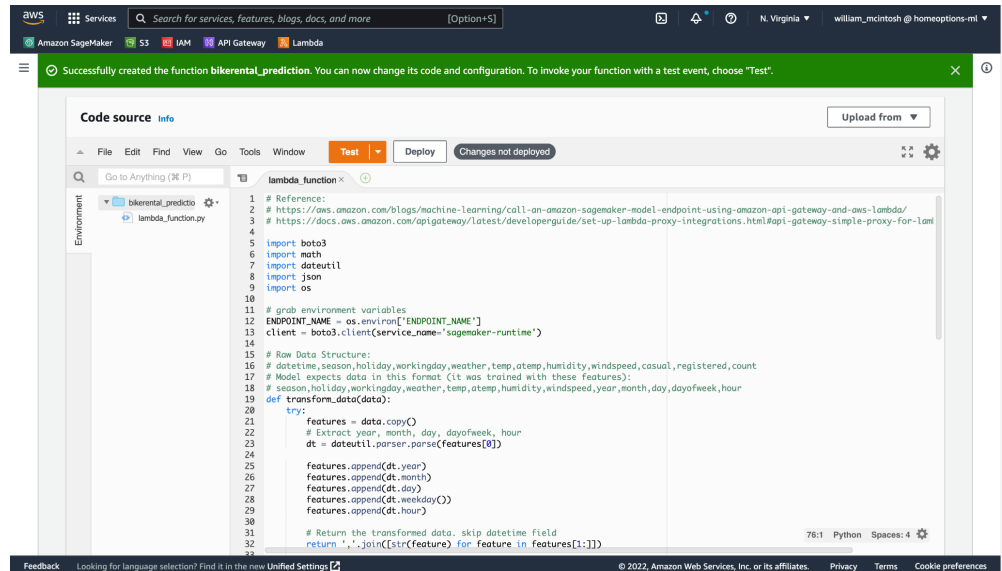
- AWS Console > Lambda > Create Function

⁸ Section 7: 67. Lab - How to run predictions against an existing endpoint

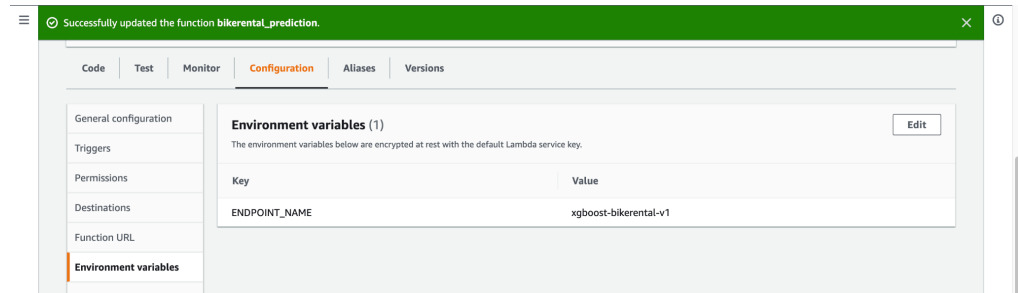
⁹ Section 8: 85. Lab - Microservice - Lambda to Endpoint

¹⁰ Section 8: 85. Lab - Microservice - Lambda to Endpoint

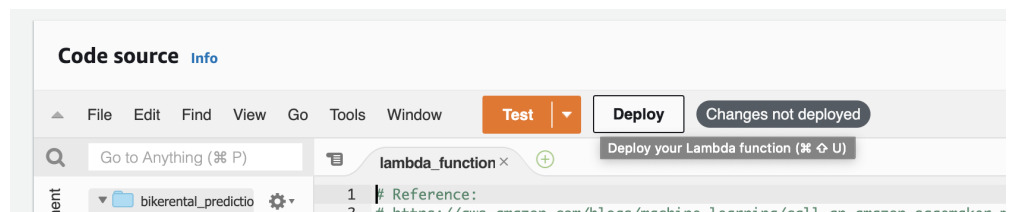
- b. Author From Scratch
- c. Function name = “bikerental_prediction”
- d. Runtime = Python3.9 (or latest supported)
- e. Change default execution role > Use existing role > lambda_sagemaker_invoke_endpoint
- f. Create function
- g. Scroll down > Copy and paste the from this file:
 - [Google Drive Link to Code](#)
 - Should look like this:



- h. Configuration tab > Environment Variables > Edit > Add Environment Variable
- i. Key = “ENDPOINT_NAME”
- j. Value = “xgboost-bikerental-v1”
- k. Save
 - Should look like this:



- l. Click **DEPLOY**



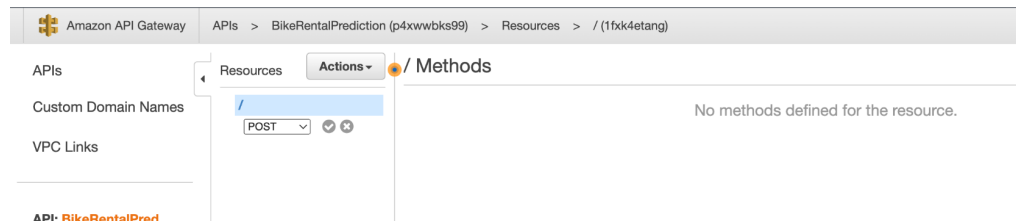
■ 

m.

13. Create an **API Gateway** to the endpoint:¹¹

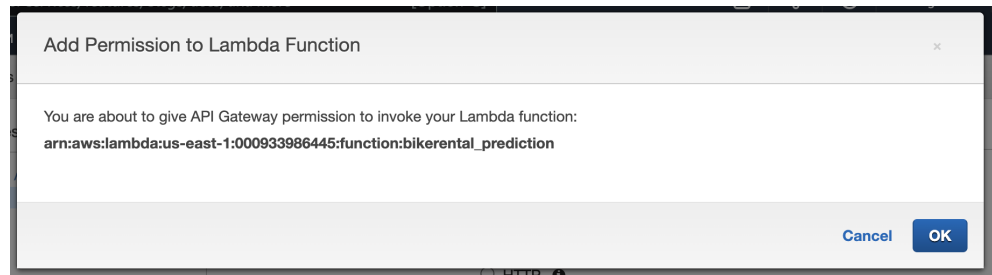
- AWS Console > API > APIs > “Create API”
- Locate “Rest API” > Build
- Click REST (non private)
- Click New API
- API name = “BikeRentalPrediction”
- Create API
- Drop down “Actions” > Create Method > POST

■ Should look like this:



- Integration Type = Lambda Function
- Lambda Function = “bikerental_prediction” (From above)

■ Add permission

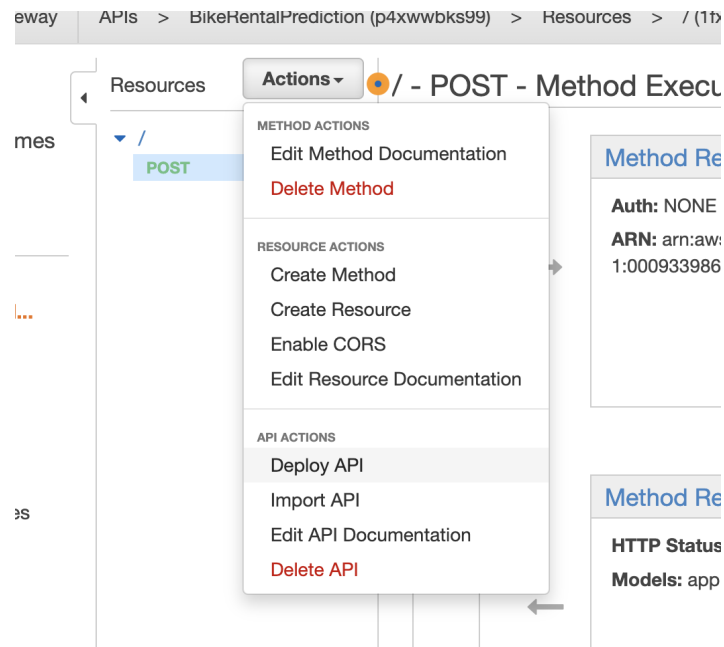


14. Deploy the API Gateway:

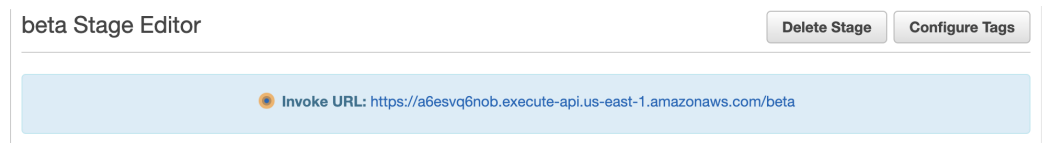
- Continuing from steps above,
- Drop down “Actions” > Deploy API

■ Should look like this:

¹¹ Section 8: 87. Lab - API Gateway, Lambda, Endpoint



- c. Deployment Stage > New Stage
- d. Stage name = "beta"
- e. Deploy
- f. Invoke URL is the url we need for testing



15. Test on Postman

- a. URL = the url from the step above
- b. POST
- c. Body > Raw
- d. Copy and paste the code below into the body

```
{
  "instances": [
    {
      "features": [
        "2012-12-19 17:00:00",
        4,
        0,
        1,
        1,
        16.4,
        20.455,

```

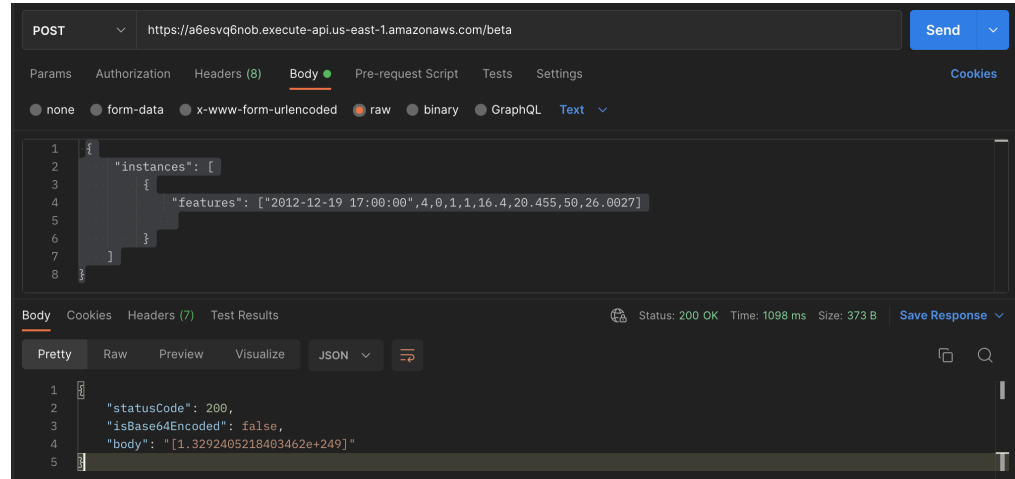


```

    50,
    26.0027
  ]
}
]
}

```

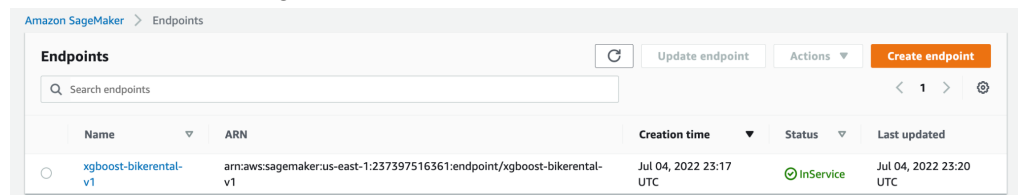
e. Send



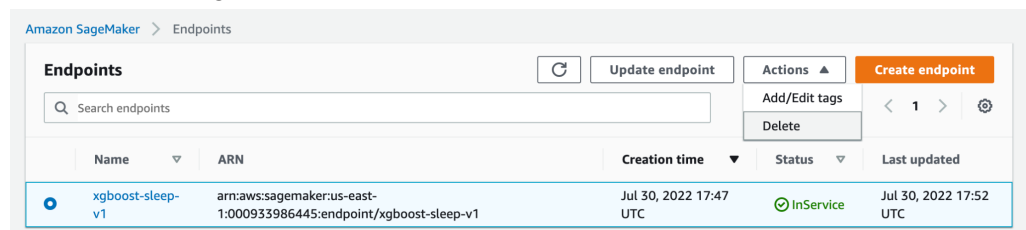
f.

16. DELETE THE ENDPOINT to prevent charges

a. AWS Console > Amazon SageMaker > Inference > Endpoints



b. Check the box on the right > Drop down Actions > Delete > Confirm delete



17. Recreate the endpoint for the future

a. AWS Console > Amazon SageMaker > Inference > Endpoints

b. Create endpoint

- c. To actually run the endpoint, you need to create a new endpoint using that configuration. To do this go to: AWS Console > Amazon SageMaker > Inference > Endpoints.
- d. Create Endpoint > Endpoint name = "xgboost-bikerental-v1" (Or, I use the convention of the same name as the endpoint configuration.)
- e. Use An Existing Configuration
- f. Check the box to the "xgboost-bikerental-v1" configuration > Select endpoint configuration"
- g. Create Endpoint
- h. **NOTE!!!** :
 - Don't forget to **delete** the endpoint when you're done!

Adding Dependencies (Libraries) To Your Endpoint

18. Adding Dependencies to your Lambda Function¹²:

- a. In order to add dependencies to your lambda function (you know, the thing at the top of your python code where you write "import requests") you'll need to upload a zip file into your lambda **which could reset all your code** if you don't copy your existing code when explained below.
- b. Open a command prompt and create a my-sourcecode-function project directory. For example, on macOS:
 - `mkdir my-sourcecode-function`
- c. Navigate to the my-sourcecode-function project directory.
 - `cd my-sourcecode-function`
- d. Copy the contents of the following sample Python code and save it in a new file named `lambda_function.py` or copy your existing `lambda_function.py` code:
 - ```
import requests
def lambda_handler(event, context):
 response = requests.get("https://www.example.com/")
 print(response.text)
 return response.text
```
- e. Your directory structure should look like this:
  - ```
my-sourcecode-function$
| lambda_function.py
```
- f. Install the requests library to a new package directory.
 - `pip install --target ./package requests`
 - `pip install --target ./package pandas`
- g. Create a deployment package with the installed library at the root.

¹² <https://docs.aws.amazon.com/lambda/latest/dg/python-package.html>

- `cd package`
`zip -r ../my-deployment-package.zip .`

h. This generates a `my-deployment-package.zip` file in your project directory. The command produces the following output:

- `adding: chardet/ (stored 0%)`
`adding: chardet/enums.py (deflated 58%)`
`...`

i. Add the `lambda_function.py` file to the root of the zip file.

- `cd ..`
`zip -g my-deployment-package.zip lambda_function.py`

j. Upload the zip file to the lambda function