

Project Deliverable 3

Team-6

Team Members:

1. Aniket Ghumed
2. Anantnaval Gaikward
3. Rimjhim Jain
4. Saumit Chinchkandi
5. Sudhanshu Dalvi

PHASE 3 Modeling, Evaluation, Tuning

The screenshot shows the 'Create notebook instance' page in the Amazon SageMaker console. The page is titled 'Create notebook instance' and includes a brief description of SageMaker notebooks. The main configuration section, 'Notebook instance settings', contains the following fields:

- Notebook instance name:** A text input field with the value 'amazon-data-train'.
- Notebook instance type:** A dropdown menu with the value 'ml.t3.medium'.
- Elastic Inference:** A dropdown menu with the value 'none'.
- Platform identifier:** A dropdown menu with the value 'Amazon Linux 2, Jupyter Lab 3'.

Below these settings are expandable sections for 'Additional configuration', 'Permissions and encryption', 'Network - optional', 'Git repositories - optional', and 'Tags - optional'. The 'Permissions and encryption' section is currently expanded, showing options for IAM role, root access, and encryption key.

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.

Role
LabRole

☐ Create role using the role creation wizard

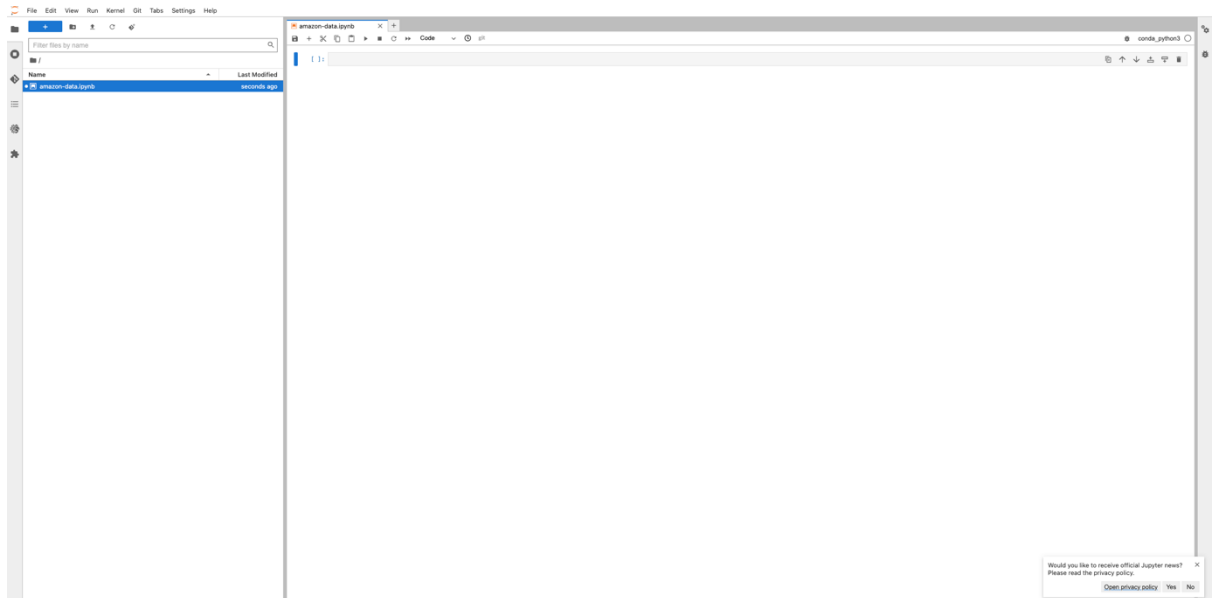
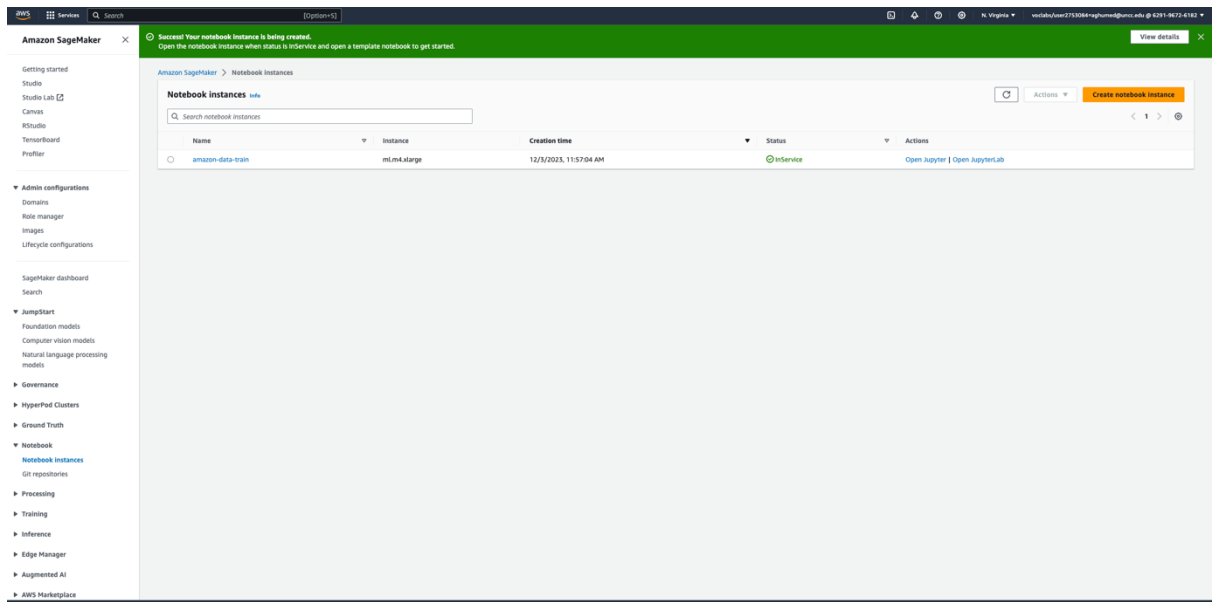
Root access - optional

☒ Enable - Give users root access to the notebook

☐ Disable - Don't give users root access to the notebook

Encryption key - optional
Encrypt your notebook data. Choose an existing KMS key or enter a key's ARN.

Key
No Custom Encryption



```
File Edit View Run Kernel Git Tabs Settings Help
amazon-data.ipynb
Name Last Modified
amazon-data.ipynb 4 minute ago

[1]: import pandas as pd
import boto3
import sagemaker
from sagemaker import get_execution_role

def load_data_from_s3(bucket, file_key):
    """Load data from S3 bucket into a pandas DataFrame."""
    sagemaker_session = sagemaker.Session()
    s3_client = sagemaker_session.boto_session.client('s3')
    response = s3_client.get_object(Bucket=bucket, Key=file_key)
    return pd.read_csv(response['Body'])

[2]: # Define your bucket and file paths
bucket = 'amazon-product-data-bucket'
categories_key = 'amazon_categories/amazon_categories.csv'
products_key = 'amazon_products/amazon_products.csv'

# Load data from S3
categories_df = load_data_from_s3(bucket, categories_key)
products_df = load_data_from_s3(bucket, products_key)

sagemaker.config INFO - Not applying SDK defaults from location: /etc/xdg/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/xdg/sagemaker/config.yaml
sagemaker.config INFO - Not applying SDK defaults from location: /home/ec2-user/.config/sagemaker/config.yaml

[3]:
```

```
File Edit View Run Kernel Git Tabs Settings Help
amazon-data.ipynb
[1]: import pandas as pd
import boto3
import sagemaker
from sagemaker import get_execution_role

def load_data_from_s3(bucket, file_key):
    """Load data from S3 bucket into a pandas DataFrame."""
    sagemaker_session = sagemaker.Session()
    s3_client = sagemaker_session.boto_session.client('s3')
    response = s3_client.get_object(Bucket=bucket, Key=file_key)
    return pd.read_csv(response['Body'])

[2]: # Define your bucket and file paths
bucket = 'amazon-product-data-bucket'
categories_key = 'amazon_categories/amazon_categories.csv'
products_key = 'amazon_products/amazon_products.csv'

# Load data from S3
categories_df = load_data_from_s3(bucket, categories_key)
products_df = load_data_from_s3(bucket, products_key)

sagemaker.config INFO - Not applying SDK defaults from location: /etc/xdg/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/xdg/sagemaker/config.yaml
sagemaker.config INFO - Not applying SDK defaults from location: /home/ec2-user/.config/sagemaker/config.yaml

[7]: # If there's no common key, you can concatenate them:
data = pd.concat([products_df, categories_df], axis=1)

[8]: # Display the first few rows of the dataframe
print("First 10 rows")
data.head(10)

[9]:
```

	asin	title	imgURL	productURL	stars	reviews	price	listPrice	category_id	isBestSeller	boughtIn	lastMonth	id	category_name
0	B014TMV5YE	Sion Softside Expandable Roller Luggage, Black...	https://m.media-amazon.com/images/I/81SdLQKYY...	https://www.amazon.com/dp/B014TMV5YE	4.5	0	139.99	0.00	104	False	2000	1.0		Beading & Jewelry Making
1	B07DGLCQXV	Luggage Sets Expandable PC+ABS Durable Suitcas...	https://m.media-amazon.com/images/I/81bQim7vR...	https://www.amazon.com/dp/B07DGLCQXV	4.5	0	189.99	209.99	104	False	1000	2.0		Fabric Decorating
2	B07XSCCYVG	Platinum Elite Softside Expandable Checked Luga...	https://m.media-amazon.com/images/I/7EAS5SzrB...	https://www.amazon.com/dp/B07XSCCYVG	4.6	0	365.49	429.99	104	False	300	3.0		Knitting & Crochet Supplies
3	B08MVFQJMJ	Freeform Hardside Expandable with Double Spinn...	https://m.media-amazon.com/images/I/6P1kSNYLQyL...	https://www.amazon.com/dp/B08MVFQJMJ	4.6	0	291.59	354.37	104	False	400	4.0		Printmaking Supplies
4	B01D,JKZBA	Winfield 2 Hardside Expandable Luggage with Sp...	https://m.media-amazon.com/images/I/6fTNJoaZCP9...	https://www.amazon.com/dp/B01D,JKZBA	4.5	0	174.99	309.99	104	False	400	5.0		Scrapbooking & Stamping Supplies
5	B07XSCD2R4	Maxlite 5 Softside Expandable Luggage with 4 S...	https://m.media-amazon.com/images/I/61...B8uSS...	https://www.amazon.com/dp/B07XSCD2R4	4.5	0	144.49	0.00	104	False	500	6.0		Sewing Products
6	B07MXF4G8K	Hard Shell Carry on Luggage Airline Approved, ...	https://m.media-amazon.com/images/I/7TCgHlYhA...	https://www.amazon.com/dp/B07MXF4G8K	4.5	0	169.99	0.00	104	False	400	7.0		Craft & Hobby Fabric
7	B07H5HSVCZ	Maxporter II 30" Hardside Spinner Trunk Luggag...	https://m.media-amazon.com/images/I/8f3h+YH0K...	https://www.amazon.com/dp/B07H5HSVCZ	4.5	0	299.99	0.00	104	False	100	8.0		Needlework Supplies
8	B08BXCBNMQ	Omni 2 Hardside Expandable Luggage with Spinn...	https://m.media-amazon.com/images/I/6KQW04m9S...	https://www.amazon.com/dp/B08BXCBNMQ	4.5	0	112.63	137.04	104	False	500	9.0		Arts, Crafts & Sewing Storage
9	B089K4XTS	Luggage Sets Expandable Lightweight Suitcases ...	https://m.media-amazon.com/images/I/81dsv5GCL...	https://www.amazon.com/dp/B089K4XTS	4.4	0	209.99	0.00	104	False	200	10.0		Painting, Drawing & Art Supplies

```
File Edit View Run Kernel Git Tabs Settings Help
amazon-data.ipynb conda_python3
8 BOB8XBCNMq Cmri 2 Hardside Expandable Luggage with Spinn... https://m.media-amazon.com/images/I/91eOWP4myS... https://www.amazon.com/dp/BOB8XBCNMq 4.5 0 112.63 137.04 104 False 500 9.0 Arts, Crafts & Sewing Storage
9 BOB8K4XTS Luggage Sets Expandable Lightweight Suitcases ... https://m.media-amazon.com/images/I/81dev5GdCL... https://www.amazon.com/dp/BOB8K4XTS 4.4 0 209.99 0.00 104 False 200 10.0 Painting, Drawing & Art Supplies

[9]: # Get the shape of the data
shape = data.shape

# Printing the shape
print("Shape of the DataFrame: Rows={}, Columns={}".format(shape[0], shape[1]))
Shape of the DataFrame: Rows=1426337, Columns=13

[10]: # Numerical Statistical Analysis
data.describe()

[10]:
      stars      reviews      price      listPrice  category_id  boughtInLastMonth      id
count  1.426337e+06  1.426337e+06  1.426337e+06  1.426337e+06  1.426337e+06  1.426337e+06  248.000000
mean    3.999512e+00  1.807508e+02  4.337540e+01  1.244916e+01  1.237409e+02  1.419823e+02  133.875000
std     1.344292e+00  1.761453e+03  1.302893e+02  4.611198e+01  7.311273e+01  8.362720e+02  77.132441
min     0.000000e+00  0.000000e+00  0.000000e+00  0.000000e+00  1.000000e+00  0.000000e+00  1.000000
25%     4.100000e+00  0.000000e+00  1.199000e+01  0.000000e+00  6.500000e+01  0.000000e+00  67.500000
50%     4.400000e+00  0.000000e+00  1.995000e+01  0.000000e+00  1.200000e+02  0.000000e+00  137.500000
75%     4.600000e+00  0.000000e+00  3.599000e+01  0.000000e+00  1.760000e+02  5.000000e+01  199.250000
max     5.000000e+00  3.465630e+05  1.973181e+04  9.999900e+02  2.700000e+02  1.000000e+05  270.000000

[11]: ## Get the information
print("Info of the data")
data.info()

Info of the data
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1426337 entries, 0 to 1426336
Data columns (total 13 columns):
#   Column      Non-Null Count  Dtype
---  ---
0   asin        1426337 non-null  object
1   title       1426336 non-null  object
2   imgURL     1426337 non-null  object
3   productURL 1426337 non-null  object
4   stars       1426337 non-null  float64
5   reviews    1426337 non-null  int64
6   price       1426337 non-null  float64
7   listPrice  1426337 non-null  float64
8   category_id 1426337 non-null  int64
9   isBestSeller 1426337 non-null  bool
10  boughtInLastMonth 1426337 non-null  int64
11  id          248 non-null      float64
12  category_name 248 non-null      object
dtypes: bool(1), float64(4), int64(3), object(5)
memory usage: 131.9+ MB

Simple 0 1 Fully initialized conda_python3 | Idle Mode: Command Ln 3, Col 12 amazon-data.ipynb 1
```

```
File Edit View Run Kernel Git Tabs Settings Help
amazon-data.ipynb conda_python3
7 listPrice 1426337 non-null float64
8 category_id 1426337 non-null int64
9 isBestSeller 1426337 non-null bool
10 boughtInLastMonth 1426337 non-null int64
11 id 248 non-null float64
12 category_name 248 non-null object
dtypes: bool(1), float64(4), int64(3), object(5)
memory usage: 131.9+ MB

[12]: #Unique values in each column
unique_values = data.nunique()
unique_values

[12]:
asin      1426337
title     1385438
imgURL    1372162
productURL 1426337
stars      41
reviews   11861
price      29961
listPrice  14518
category_id 248
isBestSeller 2
boughtInLastMonth 30
id         248
category_name 248
dtype: int64

[13]: # Get the percentage of non unique values
# Step 1: Calculate the total number of rows
total_rows = len(data)

# Step 2: Calculate the number of rows with unique values for all columns
unique_rows = len(data.drop_duplicates())

# Step 3: Calculate the number of rows with at least one duplicate value in any column
non_unique_rows = total_rows - unique_rows

# Step 4: Calculate the percentage of non-unique rows
percentage_non_unique = (non_unique_rows / total_rows) * 100

# Print the percentage of non-unique values
print("Percentage of Non-Unique Values: {:.2f}%".format(percentage_non_unique))
Percentage of Non-Unique Values: 0.00%

[14]: # Get the last few rows of the data
data.tail(10)

[14]:
      asin      title      imgURL      productURL  stars  reviews  price  listPrice  category_id  isBestSeller  boughtInLastMonth  id  category_name
1426327  B072BYSF77  Multifunctional Sports Stretchable Seamless Ca... https://m.media-amazon.com/images/I/71UTVt9hRV... https://www.amazon.com/dp/B072BYSF77  4.5  0  16.99  0.00  112  False  0  NaN  NaN
1426328  B0CDDP193JG  Manana sara Bonito hat Cotton Soft top Embroid... https://m.media-amazon.com/images/I/71Dh3FfGZ... https://www.amazon.com/dp/B0CDDP193JG  0.0  0  16.99  0.00  112  False  0  NaN  NaN
1426329  B09KMM9WYF  2 Pairs Butterfly Sunglasses Butterfly Rimless... https://m.media-amazon.com/images/I/81sOW6LgR... https://www.amazon.com/dp/B09KMM9WYF  4.5  0  9.49  10.99  112  False  0  NaN  NaN
1426330  B07N4J92M  Leather Ratchet Belts for Men Automatic Buckle... https://m.media-amazon.com/images/I/51uTK9DgA... https://www.amazon.com/dp/B07N4J92M  4.6  0  35.98  0.00  112  False  0  NaN  NaN
1426331  B00BWBURJ2  1178 Genuine Leather Mens Slim Key Case Wallet... https://m.media-amazon.com/images/I/51VARC6eU... https://www.amazon.com/dp/B00BWBURJ2  3.9  0  19.99  22.99  112  False  0  NaN  NaN
1426332  B00R3LKCO  American Flag Patriotic USA Classic 5 Panel Me... https://m.media-amazon.com/images/I/71PDJf6AA... https://www.amazon.com/dp/B00R3LKCO  4.2  0  14.95  0.00  112  False  0  NaN  NaN
1426333  B098BQ7ZQ3  Men's Baseball Cap - H2O-DRI Line Up Curved Br... https://m.media-amazon.com/images/I/8127ycev4... https://www.amazon.com/dp/B098BQ7ZQ3  4.4  0  33.99  0.00  112  False  0  NaN  NaN
1426334  B07XIMVNT1  (4 Pack) Adjustable Eyeglasses and Sunglasses... https://m.media-amazon.com/images/I/61vYWT58J... https://www.amazon.com/dp/B07XIMVNT1  3.6  0  8.54  0.00  112  False  0  NaN  NaN
1426335  B06XLBGBV9  Ax2002 Aviator Sunglasses https://m.media-amazon.com/images/I/51yJD4f1x... https://www.amazon.com/dp/B06XLBGBV9  4.5  0  64.36  67.39  112  False  0  NaN  NaN
1426336  B07GH67QC8  In Hoc Signo Vinces Knights Templar Masonic Em... https://m.media-amazon.com/images/I/91K1ZKQIO... https://www.amazon.com/dp/B07GH67QC8  4.9  0  18.79  0.00  112  False  0  NaN  NaN

Simple 0 1 Fully initialized conda_python3 | Idle Mode: Command Ln 2, Col 14 amazon-data.ipynb 1
```

```
File Edit View Run Kernel Git Tabs Settings Help
amazon-data.ipynb
1426333 B09B8BQ7ZQ3 Men's Baseball Cap - H2O-DRI Line Up Curved Br... https://m.media-amazon.com/images/I/812Tycexs4... https://www.amazon.com/dp/B09B8BQ7ZQ3 4.4 0 33.99 0.00 112 False 0 NaN NaN
1426334 B07X1MVTN1 (4 Pack) Adjustable Eyeglasses and Sunglasses ... https://m.media-amazon.com/images/I/61vvYWTS9J... https://www.amazon.com/dp/B07X1MVTN1 3.6 0 8.54 0.00 112 False 0 NaN NaN
1426335 B08XLBGBV9 Ax2002 Aviator Sunglasses https://m.media-amazon.com/images/I/51+yJD4fT... https://www.amazon.com/dp/B08XLBGBV9 4.5 0 54.36 57.39 112 False 0 NaN NaN
1426336 B07GH67QCB In Hoc Signo Vinces Knights Templar Masonic Em... https://m.media-amazon.com/images/I/91K2K2QIOE... https://www.amazon.com/dp/B07GH67QCB 4.9 0 18.79 0.00 112 False 0 NaN NaN

[15]: # Get the total number of products
total_products = len(data)
total_products

# Get data types of each column
data_types = data.dtypes
data_types

[15]: asin          object
      title       object
      imgURL      object
      productURL  object
      stars       float64
      reviews    int64
      price       float64
      listPrice   float64
      category_id int64
      isBestSeller bool
      boughtInLastMonth int64
      id          float64
      category_name object
      dtype: object

[16]: # Average price and rating by category
avg_price_by_category = data.groupby('category_name')['price'].mean()
avg_price_by_category

[16]: category_name
Abrasive & Finishing Products    149.99
Accessories & Supplies           22.79
Additive Manufacturing Products   165.99
Arts & Crafts Supplies           144.79
Arts, Crafts & Sewing Storage     112.63
...
Women's Watches                  297.49
Xbox 360 Games, Consoles & Accessories 165.99
Xbox One Games, Consoles & Accessories 399.00
Xbox Series X & S Consoles, Games & Accessories 179.99
eBook Readers & Accessories       154.99
Name: price, Length: 248, dtype: float64

[17]: # Average price and rating by category
avg_rating_by_category = data.groupby('category_name')['stars'].mean()
avg_rating_by_category

[17]: category_name
Abrasive & Finishing Products    4.4
Accessories & Supplies           4.0
Additive Manufacturing Products   4.3
Arts & Crafts Supplies           4.7
Arts, Crafts & Sewing Storage     4.5
...
Women's Watches                  4.7
Xbox 360 Games, Consoles & Accessories 4.4
Xbox One Games, Consoles & Accessories 4.9
Xbox Series X & S Consoles, Games & Accessories 5.0
eBook Readers & Accessories       4.6
Name: stars, Length: 248, dtype: float64

Simple Fully initialized conda_python3 | ide Mode: Command Ln 3, Col 23 amazon-data.ipynb 1
```

```
File Edit View Run Kernel Git Tabs Settings Help
amazon-data.ipynb
Xbox one games, consoles & accessories 399.00
Xbox Series X & S Consoles, Games & Accessories 179.99
eBook Readers & Accessories 154.99
Name: price, Length: 248, dtype: float64

[17]: # Average price and rating by category
avg_rating_by_category = data.groupby('category_name')['stars'].mean()
avg_rating_by_category

[17]: category_name
Abrasive & Finishing Products    4.4
Accessories & Supplies           4.0
Additive Manufacturing Products   4.3
Arts & Crafts Supplies           4.7
Arts, Crafts & Sewing Storage     4.5
...
Women's Watches                  4.7
Xbox 360 Games, Consoles & Accessories 4.4
Xbox One Games, Consoles & Accessories 4.9
Xbox Series X & S Consoles, Games & Accessories 5.0
eBook Readers & Accessories       4.6
Name: stars, Length: 248, dtype: float64

[18]: # Top categories by the number of products
top_categories = data['category_name'].value_counts().head(10)
top_categories

[18]: category_name
Reading & Jewelry Making    1
Small Animal Supplies       1
Heating, Cooling & Air Quality 1
Kids' Home Store            1
Home Storage & Organization 1
Wall Art                   1
Vacuum Cleaners & Floor Care 1
Ironing Products            1
Party Supplies              1
Pet Bird Supplies           1
Name: count, dtype: int64

[19]: # Best-selling products
best_selling_products = data[data['isBestSeller']]
best_selling_products

[19]:   asin          title                                     imgURL      productURL  stars  reviews  price  listPrice  category_id  isBestSeller  boughtInLastMonth  id  category_name
924  B00W66LQFO  Men's Eversoft Cotton Stay Tucked Crew T-Shirt https://m.media-amazon.com/images/I/513raGQKW... https://www.amazon.com/dp/B00W66LQFO 4.6 0 18.48 26.00 110 True 10000 NaN NaN
925  B0C4RMF5PZ  Official Renaissance World Tour Merch Disco Co... https://m.media-amazon.com/images/I/71uht5KQ... https://www.amazon.com/dp/B0C4RMF5PZ 4.8 0 40.00 0.00 110 True 4000 NaN NaN
933  B077ZMKVWM  Men's Crew T-Shirts, Multipack, Style G1100 https://m.media-amazon.com/images/I/61Z5AAGPW... https://www.amazon.com/dp/B077ZMKVWM 4.6 0 18.99 0.00 110 True 8000 NaN NaN
938  B07PHZVWX1  Men's Coolzone Boxer Briefs, Moisture Wicking ... https://m.media-amazon.com/images/I/81HfvXX0... https://www.amazon.com/dp/B07PHZVWX1 4.6 0 19.59 24.49 110 True 5000 NaN NaN
944  B01D2QRTIE  Men's Multi-Pack Mesh Ventilating Comfort Fir ... https://m.media-amazon.com/images/I/61tpudqR... https://www.amazon.com/dp/B01D2QRTIE 4.6 0 14.99 0.00 110 True 6000 NaN NaN
...
1415016 B0007LCLEP Pyle 2Way Custom Component Speaker System-6.5"... https://m.media-amazon.com/images/I/812gPY56m... https://www.amazon.com/dp/B0007LCLEP 4.0 0 44.99 47.99 26 True 0 NaN NaN
1415072 B0045502B2 BOSS Audio Systems R1002 Riot Series Car Stere... https://m.media-amazon.com/images/I/610Vz6PUF... https://www.amazon.com/dp/B0045502B2 4.1 0 35.41 56.00 26 True 0 NaN NaN
1416949 B08ZVZVN8N Pickleball Paddles, USAFA Approved Fiberglass ... https://m.media-amazon.com/images/I/71WQNbP6G... https://www.amazon.com/dp/B08ZVZVN8N 4.8 0 35.99 0.00 198 True 300 NaN NaN
1416996 B0D0HR1NWE Skechers Men's Afterburn M. Fit https://m.media-amazon.com/images/I/81viYs9J... https://www.amazon.com/dp/B0D0HR1NWE 4.4 0 40.00 74.00 198 True 200 NaN NaN
1417015 B083LL43QF FULLSOLT 3 Pack Leggings for Women Non See Thr... https://m.media-amazon.com/images/I/71cYXDO8... https://www.amazon.com/dp/B083LL43QF 4.3 0 20.39 29.99 198 True 100 NaN NaN

8520 rows x 13 columns

Simple Fully initialized conda_python3 | ide Mode: Command Ln 3, Col 22 amazon-data.ipynb 1
```

```
File Edit View Run Kernel Git Tabs Settings Help
amazon-data.ipynb
conda_python3

[19]: # Best-selling products
best_selling_products = data[data['isBestSeller']]
best_selling_products

[19]:
```

	asin	title	imgUrl	productURL	stars	reviews	price	listPrice	category_id	isBestSeller	boughtInLastMonth	id	category_name
924	B00W66LQFO	Men's Eversoft Cotton Stay Tucked Crew T-Shirt	https://m.media-amazon.com/images/I/513raGQXW...	https://www.amazon.com/dp/B00W66LQFO	4.6	0	18.48	26.00	110	True	10000	NaN	NaN
925	B0C4RMF5PZ	Official Renaissance World Tour Merch Disco Co...	https://m.media-amazon.com/images/I/71Uxhf5KQ...	https://www.amazon.com/dp/B0C4RMF5PZ	4.8	0	40.00	0.00	110	True	4000	NaN	NaN
938	B077H2VWX1	Men's Crew T-Shirts, Multipack, Style G1100	https://m.media-amazon.com/images/I/6125AAGPW...	https://www.amazon.com/dp/B077H2VWX1	4.6	0	18.99	0.00	110	True	8000	NaN	NaN
938	B077H2VWX1	Men's Coolzone Boxer Briefs, Moisture Wicking ...	https://m.media-amazon.com/images/I/81HfVvXX0...	https://www.amazon.com/dp/B077H2VWX1	4.6	0	19.59	24.49	110	True	5000	NaN	NaN
944	B01D2GRTIE	Men's Multi-Pack Mesh Ventilating Comfort Fit ...	https://m.media-amazon.com/images/I/61kxpdqRp...	https://www.amazon.com/dp/B01D2GRTIE	4.6	0	14.99	0.00	110	True	6000	NaN	NaN
...
1415016	B0007CLPE	Pyle 2Way Custom Component Speaker System-6.5"	https://m.media-amazon.com/images/I/812gPY5Bm...	https://www.amazon.com/dp/B0007CLPE	4.0	0	44.99	47.99	26	True	0	NaN	NaN
1415072	B0045502B2	BOSS Audio Systems R1002 Riot Series Car Stere...	https://m.media-amazon.com/images/I/61QVz6PURL...	https://www.amazon.com/dp/B0045502B2	4.1	0	35.41	56.00	26	True	0	NaN	NaN
1416949	B082V2V2BN	Pickleball Paddles, USAPA Approved Fiberglass ...	https://m.media-amazon.com/images/I/71WQbP6dW...	https://www.amazon.com/dp/B082V2V2BN	4.8	0	35.99	0.00	198	True	300	NaN	NaN
1416990	B0DHRN8WE	Slachters Men's Afterburn M. Fit	https://m.media-amazon.com/images/I/81vYv9dYt...	https://www.amazon.com/dp/B0DHRN8WE	4.4	0	40.00	74.00	198	True	200	NaN	NaN
1417015	B083LL43QF	FULLSOFT 3 Pack Leggings for Women Non See Thr...	https://m.media-amazon.com/images/I/71cYXOD8...	https://www.amazon.com/dp/B083LL43QF	4.3	0	20.39	29.99	198	True	100	NaN	NaN

```
8520 rows x 13 columns

[20]: # Compare the ratings
avg_rating_best_sellers = best_selling_products['stars'].mean()
avg_rating_best_sellers

[20]: 4.494837558685447

[21]: # Compare prices
avg_price_best_sellers = best_selling_products['price'].mean()
avg_price_best_sellers

[21]: 29.721798122865728

[22]: # Distribution of prices
price_distribution = data['price'].describe()
price_distribution

[22]: count    1.426337e+06
      mean    4.337548e+01
      std     1.382893e+02
      min      0.000000e+00
      25%     1.199000e+01
      50%     1.995000e+01
      75%     3.599000e+01
      max     1.973181e+04
      Name: price, dtype: float64

Simple Fully initialized conda_python3 | ide Mode: Command Ln 3, Col 19 amazon-data.ipynb 1
```

```
File Edit View Run Kernel Git Tabs Settings Help
amazon-data.ipynb
conda_python3

8520 rows x 13 columns

[20]: # Compare the ratings
avg_rating_best_sellers = best_selling_products['stars'].mean()
avg_rating_best_sellers

[20]: 4.494837558685447

[21]: # Compare prices
avg_price_best_sellers = best_selling_products['price'].mean()
avg_price_best_sellers

[21]: 29.721798122865728

[22]: # Distribution of prices
price_distribution = data['price'].describe()
price_distribution

[22]: count    1.426337e+06
      mean    4.337548e+01
      std     1.382893e+02
      min      0.000000e+00
      25%     1.199000e+01
      50%     1.995000e+01
      75%     3.599000e+01
      max     1.973181e+04
      Name: price, dtype: float64

[23]: # Average rating and review count
avg_rating = data['stars'].mean()
avg_review_count = data['reviews'].mean()

[23]: 4.494837558685447 0.03718638544799832

[24]: # Correlation between ratings and reviews
correlation = data['stars'].corr(data['reviews'])
correlation

[24]: 0.83718638544799832

[25]: # Number of missing values
missing_id_count = data['id'].isnull().sum()
missing_category_count = data['category_name'].isnull().sum()

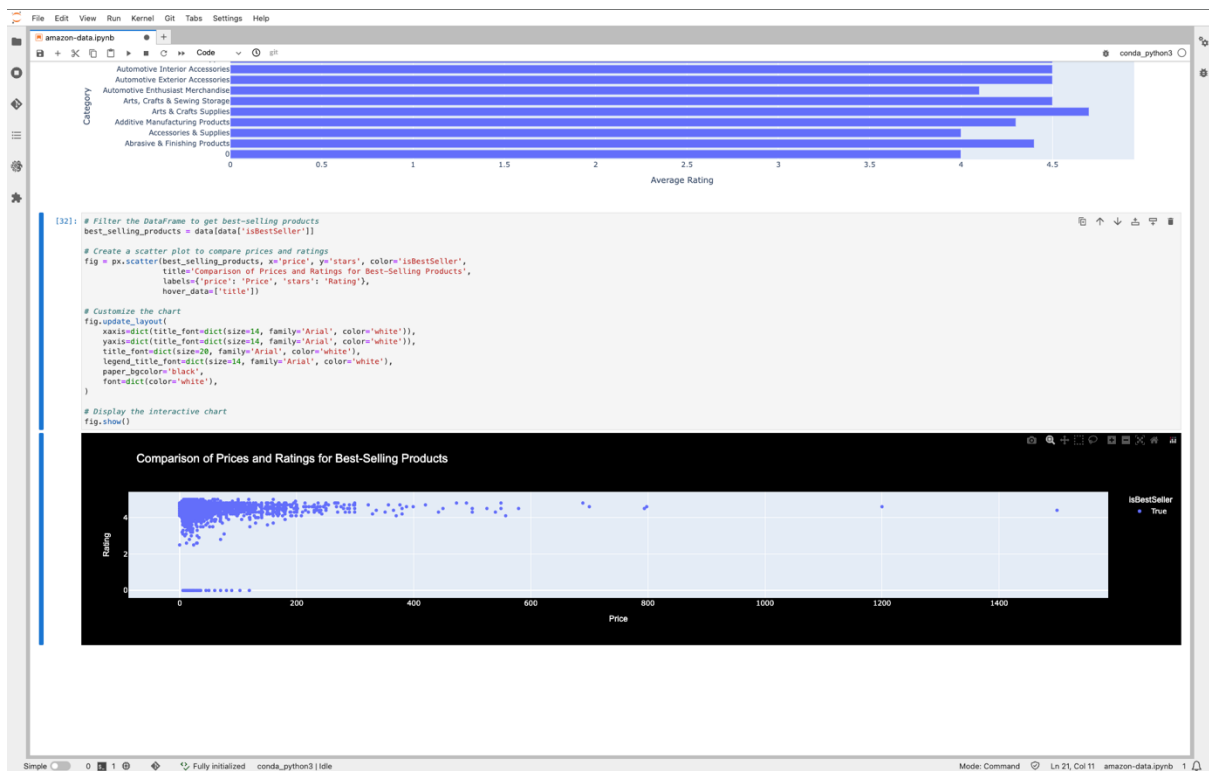
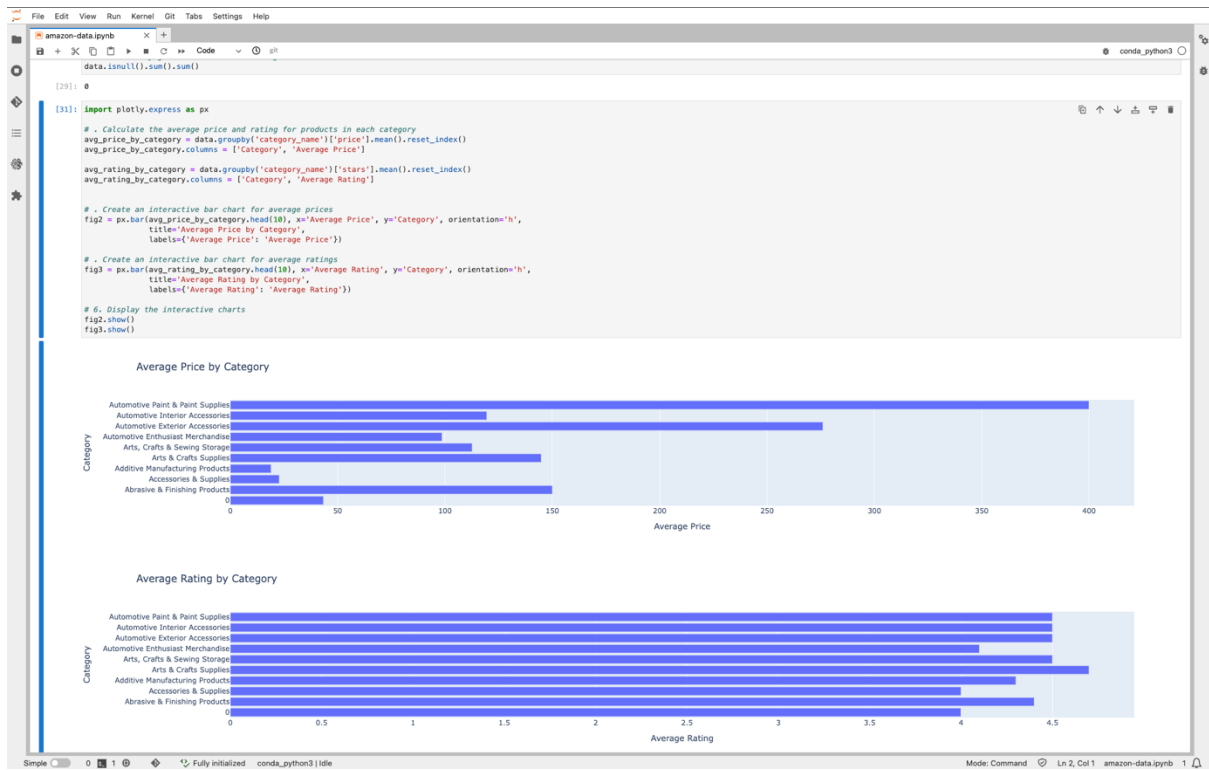
[25]: 0 1426889
     0

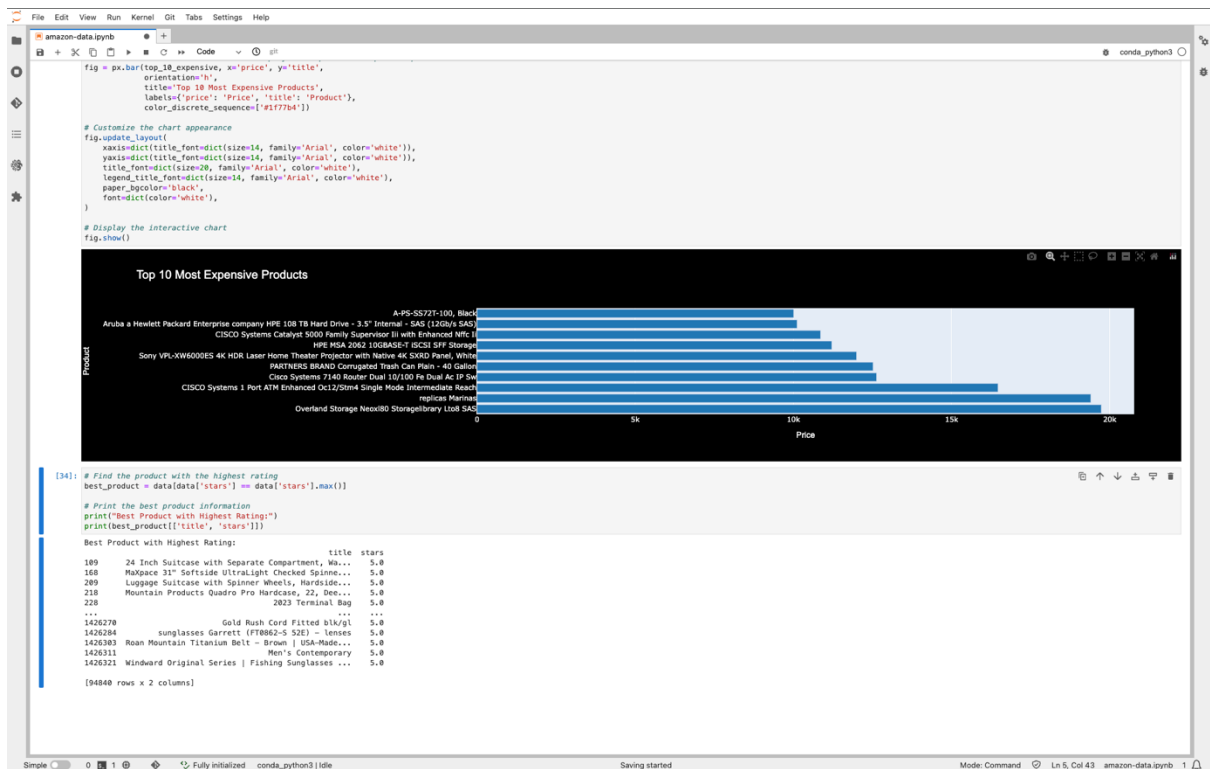
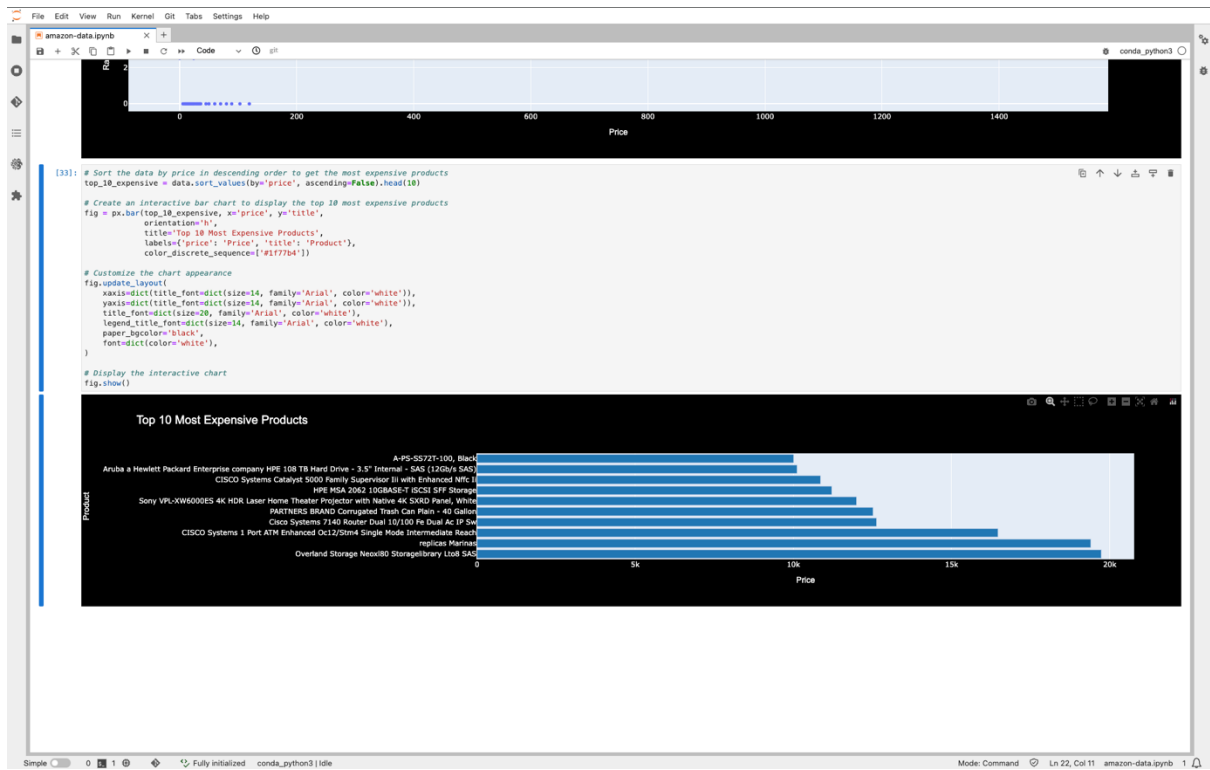
[26]: # Get the total of Missing values
data.isnull().sum()

[26]: asin      0
      title    1
      imgUrl   0
      productURL 0
      stars     0
      reviews  0
      price     0
      listPrice 0
      category_id 0
      isBestSeller 0
      boughtInLastMonth 0
      id      1426889
      category_name 1426889
      dtype: int64

[28]: # Fill Missing values with '0'
data.fillna(0, inplace=True)

Simple Fully initialized conda_python3 | ide Mode: Command Ln 1, Col 31 amazon-data.ipynb 1
```





File Edit View Run Kernel Git Tabs Settings Help

amazon-data.ipynb

train.py

amazon-data.ipynb

train.py

training_data.csv

amazon-data.ipynb

train.py

```
model = LogisticRegression()

# Train the model on the training data
model.fit(X_train, y_train)

# Make predictions on the testing data
predictions = model.predict(X_test)

# Evaluate the model
print(classification_report(y_test, predictions))

precision    recall  f1-score   support

   False    0.99    1.00    1.00   283544
    True    0.33    0.01    0.01    1724

 accuracy    0.66    0.50    0.51   285268
  macro avg   0.66    0.50    0.51   285268
 weighted avg   0.99    0.99    0.99   285268

[4]: data.to_csv("training_data.csv")

[ ]: import sagemaker
from sagemaker.sklearn.estimator import SKLearn

role = sagemaker.get_execution_role()
sagemaker_session = sagemaker.Session()

# Create a SageMaker Estimator
estimator = SKLearn(
    entry_point="train.py",
    role=role,
    instance_count=1,
    instance_type="ml.m4.xlarge",
    framework_version="0.23-1",
    py_version="py3",
    hyperparameters={"train": "s3://amazon-product-data-bucket/training_job/"})

# Start the training job
estimator.fit({"train": "s3://amazon-product-data-bucket/training_job/"})
```

Simple Fully initialized conda_python3 | idle Saving started Mode: Edit Ln 18, Col 25 amazon-data.ipynb 1

RWS Services Search [Options=5]

Global voclabi/user2753084+aghumed@uncc.edu @ 6291-9672-6182

Upload succeeded

View details below.

Upload: status

Close

The information below will no longer be available after you navigate away from this page.

Summary

Destination

s3://amazon-product-data-bucket/training_job/

Succeeded

1 file, 376.5 MB (100.00%)

Failed

0 files, 0 B (0%)

Files and folders

Configuration

Files and folders (1 Total, 376.5 MB)

Find by name

training_data.csv

-

text/csv

376.5 MB

Succeeded

-

CloudShell Feedback

© 2021, Amazon Web Services, Inc. or its affiliates. Privacy Terms Cookie preferences

File Edit View Run Kernel Git Tabs Settings Help

amazon-data.ipynb

train.py

training_data.csv

Filter files by name

Name

Last Modified

amazon-data.ipynb

seconds ago

train.py

5 minutes ago

training_data.csv

4 minutes ago

amazon-data.ipynb

train.py

conda_python3

```

# Train the model on the training data
model.fit(X_train, y_train)

# Make predictions on the testing data
predictions = model.predict(X_test)

# Evaluate the model
print(classification_report(y_test, predictions))

precision    recall  f1-score   support

   False    0.99    1.00    1.00    283544
    True    0.33    0.83    0.43     1724

 accuracy    0.66    0.50    0.51    285268
 macro avg   0.99    0.99    0.99    285268
weighted avg   0.99    0.99    0.99    285268

[4]: data.to_csv("training_data.csv")

[5]: import sagemaker
from sagemaker.sklearn.estimator import SKLearn

role = sagemaker.get_execution_role()
sagemaker_session = sagemaker.Session()

# Create a SageMaker Estimator
estimator = SKLearn(
    entry_point="train.py",
    role=role,
    instance_count=1,
    instance_type="ml.m4.xlarge",
    framework_version="0.23-1",
    py_version="py3",
    hyperparameters={"train": "s3://amazon-product-data-bucket/training_job/"})

# Start the training job
estimator.fit({"train": "s3://amazon-product-data-bucket/training_job/"})

sagemaker.config INFO - Not applying SDK defaults from location: /etc/xdg/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/xdg/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/xdg/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-scikit-learn-2023-12-03-20-24-53-363
Using provided s3_resource
2023-12-03 20:24:53 Starting - Starting the training job..

```

Simple

Fully initialized

conda_python3 | Busy

Mode: Command

Ln 20, Col 1

amazon-data.ipynb

BWS

Services

Search

[Options=5]

N. Virginia

vocalab/user2733084+aghumed@uncc.edu @ 6291-9672-6182

Amazon SageMaker

Getting started

Studio

Studio Lab

Canvas

RSudio

TensorBoard

Profiler

▼ Admin configurations

Domains

Role manager

Images

Lifecycle configurations

SageMaker dashboard

Search

▼ JumpStart

Foundation models

Computer vision models

Natural language processing models

► Governance

► HyperPod Clusters

► Ground Truth

► Notebook

► Processing

► Training

► Inference

► Edge Manager

► Augmented AI

► AWS Marketplace

Amazon SageMaker

Training jobs

Training jobs

Search training jobs

Actions

Create training job

1

Name	Creation time	Duration	Job status	Warm pool status	Time left
sagemaker-scikit-learn-2023-12-03-20-24-53-363	12/3/2023, 3:24:53 PM	-	InProgress	-	-

CloudShell

Feedback

© 2023, Amazon Web Services, Inc. or its affiliates. Privacy Terms Cookie preferences

Amazon SageMaker

Getting started
Studio
Studio Lab
Canvas
RStudio
TensorBoard
Profiler

▼ Admin configurations
Domains
Role manager
Images
Lifecycle configurations

SageMaker dashboard
Search

▼ JumpStart
Foundation models
Computer vision models
Natural language processing models

► Governance

► HyperPod Clusters

► Ground Truth

► Notebook

► Processing

► Training

► Inference

► Edge Manager

► Augmented AI

► AWS Marketplace

Amazon SageMaker > Training jobs > sagemaker-scikit-learn-2023-12-03-20-24-53-363

sagemaker-scikit-learn-2023-12-03-20-24-53-363

Clone Create model package Stop Create model

Job settings

Job name
sagemaker-scikit-learn-2023-12-03-20-24-53-363

ARN
arn:aws:sagemaker:us-east-1:629196726182:training-job/sagemaker-scikit-learn-2023-12-03-20-24-53-363

Status
In progress
- Starting
View history

Creation time
Dec 05, 2023 20:24 UTC

Last modified time
Dec 05, 2023 20:24 UTC

SageMaker metrics time series
Disabled

Training time (seconds)
-

Billable time (seconds)
-

Managed spot training savings
-

Tuning job source/parent
-

IAM role ARN
arn:aws:iam::629196726182:role/LabRole

Algorithm

Algorithm ARN
-

Additional volume size (GiB)
30

Maximum wait time for managed spot training(s)
-

Volume encryption key
-

Training image
685311688578.dkr.ecr.us-east-1.amazonaws.com/sagemaker-scikit-learn:0.23-1-cpu-py3

Maximum runtime (s)
86400

Managed spot training
Disabled

Input mode
File

Instance group

Instance type

Instance count

Keep alive period

-

ml.m4.xlarge

1

-

Input data configuration: train

Channel name

Input mode

Data source

S3 data type

train

-

S3

S3Prefix

Content type

Instance group

S3 data distribution type

-

-

FullyReplicated

Compression type

URI

None

s3://amazon-product-data-bucket/training_job/

CloudShell Feedback

Amazon S3 > Buckets > amazon-product-data-bucket > training_job/

training_job/

Copy S3 URI

Objects Properties

Objects (3)

Objects are the fundamental entities stored in Amazon S3. You can use [Amazon S3 Inventory](#) to get a list of all objects in your bucket. For others to access your objects, you'll need to explicitly grant them permissions. [Learn more](#)

Copy S3 URI Copy URL Download Open Delete Actions Create folder Upload

Find objects by prefix

<input type="checkbox"/>	Name	Type	Last modified	Size	Storage class
<input type="checkbox"/>	requirements.txt	txt	December 3, 2023, 15:41:29 (UTC-05:00)	94.0 B	Standard
<input type="checkbox"/>	train.py	py	December 3, 2023, 15:32:37 (UTC-05:00)	1.3 KB	Standard
<input type="checkbox"/>	training_data.csv	csv	December 3, 2023, 15:23:10 (UTC-05:00)	376.5 MB	Standard

CloudShell Feedback

Amazon SageMaker

sagemaker-sckit-learn-2023-12-03-22-22-32-965

Job settings

Job name: sagemaker-sckit-learn-2023-12-03-22-22-32-965
Status: Completed
ARN: arn:aws:sagemaker-us-east-1:629196726182:training-job/sagemaker-sckit-learn-2023-12-03-22-22-32-965

Algorithm

Algorithm ARN	Additional volume size (GB)	Maximum wait time for managed spot training(s)	Volume encryption key
-	30	-	-
Training image	Maximum runtime (s)	Managed spot training	
663315688378.dkr.ecr.us-east-1.amazonaws.com/sagemaker-sckit-learn-0.23-1-cpu-py3	86400	Disabled	
Input mode			
File			

Instance group	Instance type	Instance count	Keep alive period
-	ml.m4.xlarge	1	-

Input data configuration: train

Channel name	Input mode	Data source	S3 data type
train	-	S3	S3Prefix
	Content type	Instance group	S3 data distribution type
	-	-	FullyReplicated
	Compression type		URI
None			s3://amazon-product-data-bucket/training_job/

amazon-data.ipynb

```
SM_RESOURCE_CONFIG={"current_group_name":"homogeneousCluster","current_host":"algo-1","current_instance_type":"ml.m4.xlarge","hosts":["algo-1"],"instance_groups":[{"hosts":["algo-1"],"instance_group_name":"homogeneousCluster","instance_type":"ml.m4.xlarge"},"network_interface_name":"eth0"}
SM_INPUT_DATA_CONFIG={"train":{"RecordWrapperType":"None","S3DistributionType":"FullyReplicated","TrainingInputMode":"File"},"input_dir":"/opt/ml/input","is_master":true,"job_name":"sagemaker-sckit-learn-2023-12-03-22-22-32-965","log_level":20,"master_hostname":"algo-1","model_dir":"/opt/ml/model","module_dir":"/opt/ml/output/data","output_dir":"/opt/ml/output","output_intermediate_dir":"/opt/ml/output/intermediate","resource_config":{"current_group_name":"homogeneousCluster","current_host":"algo-1","current_instance_type":"ml.m4.xlarge"},"user_entry_point":"train.py"}
SM_USER_ARGS={"train":"/opt/ml/output/data","s3://amazon-product-data-bucket/training_job/"}
SM_OUTPUT_INTERMEDIATE_DIR="/opt/ml/output/intermediate"
SM_CHANNEL_TRAIN="/opt/ml/input/data/train"
SM_HP_TRADEOFFS="/amazon-product-data-bucket/training_job/PYTHONPATH/opt/ml/code/miniconda3/bin/miniconda3/lib/python3.7/zip/miniconda3/lib/python3.7/miniconda3/lib/python3.7/lib-dynload/miniconda3/lib/python3.7/site-packages"
Invoking script with the following command:
/miniconda3/bin/python train.py --train s3://amazon-product-data-bucket/training_job/
Requirement already satisfied: fspec in /miniconda3/lib/python3.7/site-packages (2023.1.0)
WARNING: Running pip as the 'root' user can result in broken permissions and conflicting behaviour with the system package manager. It is recommended to use a virtual environment instead: https://pip.pypa.io/warnings/venv
[notice] A new release of pip is available: 23.0 -> 23.3.1
[notice] To update, run: pip install --upgrade pip
sys:1: DtypeWarning: Columns (13) have mixed types.Specify dtype option on import or set low_memory=False.
2023-12-03 22:26:59 Uploading - Uploading generated training model
precision recall f1-score support
False 0.99 1.00 1.00 283544
True 0.33 0.01 0.01 1724
accuracy 0.99 0.99 285268
macro avg 0.66 0.50 0.51 285268
weighted avg 0.99 0.99 0.99 285268
2023-12-03 22:26:54,239 sagemaker-containers INFO Reporting training SUCCESS
2023-12-03 22:27:36 Completed - Training job completed
Training seconds: 170
Billable seconds: 170
```