

Amazon SageMaker Studio    File Edit View Run Kernel Git Tabs Settings Help    Feedback

Terminal 1    C2\_W1\_Assignment.ipynb    prepare\_data.py    2 vCPU + 4 GiB Python 3 (Data Science) Share

## Feature transformation with Amazon SageMaker processing job and Feature Store

### Introduction

In this lab you will start with the raw [Women's Clothing Reviews](#) dataset and prepare it to train a BERT-based natural language processing (NLP) model. The model will be used to classify customer reviews into positive (1), neutral (0) and negative (-1) sentiment.

You will convert the original review text into machine-readable features used by BERT. To perform the required feature transformation you will configure an Amazon SageMaker processing job, which will be running a custom Python script.

### Table of Contents

- 1. Configure the SageMaker Feature Store
  - 1.1. Configure dataset
  - 1.2. Configure the SageMaker feature store
    - Exercise 1
- 2. Transform the dataset
  - Exercise 2
  - Exercise 3
- 3. Query the Feature Store
  - 3.1. Export training, validation, and test datasets from the Feature Store
    - Exercise 4
  - 3.2. Export TSV from Feature Store
  - 3.3. Check that the dataset in the Feature Store is balanced by sentiment
    - Exercise 5
    - Exercise 6
    - Exercise 7

```
[1]: # please ignore warning messages during the installation
| pip install --disable-pip-version-check -q sagemaker==2.35.0
| conda install -q -y pytorch==1.6.0 -c pytorch
| pip install --disable pip version check -q transformers==3.5.1

/opt/conda/lib/python3.7/site-packages/secretstorage/dhcrypto.py:16: CryptographyDeprecationWarning: int_from_bytes is deprecated, use int.from_bytes instead
  from cryptography.utils import int_from_bytes
/opt/conda/lib/python3.7/site-packages/secretstorage/util.py:25: CryptographyDeprecationWarning: int_from_bytes is deprecated, use int.from_bytes instead
  from cryptography.utils import int_from_bytes
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
Collecting package metadata (current_repotdata.json): ...working... done
Solving environment: ...working... done

## Package Plan ##

environment location: /opt/conda

added / updated specs:
- pytorch==1.6.0

The following packages will be downloaded:

```

package	build
ca-certificates-2021.7.5	h06a4308_1 113 KB
certifi-2021.5.30	py37h06a4308_0 139 KB
cudatoolkit-10.2.89	hfdb86e86_1 365.1 MB
ninja-1.10.2	hff7bd54_1 1.4 MB
pytorch-1.6.0	py3.7_cuda10.2.89_cudnn7.6.5_0 537.7 MB pytorch

```
Total: 904.5 MB
```

```
The following NEW packages will be INSTALLED:

```

package	build
cudatoolkit	pkgs/main/linux-64::cudatoolkit-10.2.89-hfd86e86_1
ninja	pkgs/main/linux-64::ninja-1.10.2-hff7bd54_1
pytorch	pytorch/linux-64::pytorch-1.6.0-py3.7_cuda10.2.89_cudnn7.6.5_0

```
The following packages will be UPDATED:

```

package	build
ca-certificates	conda-forge::ca-certificates-2021.5.3~ --> pkgs/main::ca-certificates-2021.7.5-h06a4308_1

```
The following packages will be SUPERSEDED by a higher-priority channel:

```

package	build
certifi	conda-forge::certifi-2021.5.30-py37h8~ --> pkgs/main::certifi-2021.5.30-py37h06a4308_0

```
Preparing transaction: ...working... done
Verifying transaction: ...working... done
Executing transaction: ...working... done
/opt/conda/lib/python3.7/site-packages/secretstorage/dhcrypto.py:16: CryptographyDeprecationWarning: int_from_bytes is deprecated, use int.from_bytes instead
  from cryptography.utils import int_from_bytes
/opt/conda/lib/python3.7/site-packages/secretstorage/util.py:25: CryptographyDeprecationWarning: int_from_bytes is deprecated, use int.from_bytes instead
  from cryptography.utils import int_from_bytes
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
```

```
[2]: import boto3
import sagemaker
import botocore

config = botocore.config.Config(user_agent_extra='dlai-pds/c2/w1')

# Low-Level service client of the boto3 session
sm = boto3.client(service_name='sagemaker',
                   config=config)

featurestore_runtime = boto3.client(service_name='sagemaker-featurestore-runtime',
                                    config=config)

sess = sagemaker.Session(sagemaker_client=sm,
                        sagemaker_featurestore_runtime_client=featurestore_runtime)
```

```
bucket = sess.default_bucket()
role = sagemaker.get_execution_role()
region = sess.boto_region_name
```

## 1. Configure the SageMaker Feature Store

### 1.1. Configure dataset

The raw dataset is in the public S3 bucket. Let's start by specifying the S3 location of it:

```
[3]: raw_input_data_s3_uri = 's3://dlai-practical-data-science/data/raw/'
print(raw_input_data_s3_uri)
```

```
s3://dlai-practical-data-science/data/raw/
```

List the files in the S3 bucket (in this case it will be just one file):

```
[4]: aws s3 ls $raw_input_data_s3_uri
2021-04-30 02:21:06    8457214 womens_clothing_ecommerce_reviews.csv
```

### 1.2. Configure the SageMaker feature store

As the result of the transformation, in addition to generating files in S3 bucket, you will also save the transformed data in the **Amazon SageMaker Feature Store** to be used by others in your organization, for example.

To configure a Feature Store you need to setup a **Feature Group**. This is the main resource containing all of the metadata related to the data stored in the Feature Store. A Feature Group should contain a list of **Feature Definitions**. A Feature Definition consists of a name and the data type. The Feature Group also contains an online store configuration and an offline store configuration controlling where the data is stored. Enabling the online store allows quick access to the latest value for a record via the [GetRecord API](#). The offline store allows storage of the data in your S3 bucket. You will be using the offline store in this lab.

Let's setup the Feature Group name and the Feature Store offline prefix in S3 bucket (you will use those later in the lab):

```
[5]: import time
timestamp = int(time.time())

feature_group_name = 'reviews-feature-group-' + str(timestamp)
feature_store_offline_prefix = 'reviews-feature-store-' + str(timestamp)

print('Feature group name: {}'.format(feature_group_name))
print('Feature store offline prefix in S3: {}'.format(feature_store_offline_prefix))

Feature group name: reviews-feature-group-1628645791
Feature store offline prefix in S3: reviews-feature-store-1628645791
```

Taking two features from the original raw dataset (`Review_Text` and `Rating`), you will transform it preparing to be used for the model training and then to be saved in the Feature Store. Here you will define the related features to be stored as a list of `FeatureDefinition`.

```
[6]: from sagemaker.feature_store.feature_definition import (
    FeatureDefinition,
    FeatureTypeEnum,
)

feature_definitions= [
    # unique ID of the review
    FeatureDefinition(feature_name='review_id', feature_type=FeatureTypeEnum.STRING),
    # ingestion timestamp
    FeatureDefinition(feature_name='date', feature_type=FeatureTypeEnum.STRING),
    # sentiment: -1 (negative), 0 (neutral) or 1 (positive). It will be found the Rating values (1, 2, 3, 4, 5)
    FeatureDefinition(feature_name='sentiment', feature_type=FeatureTypeEnum.STRING),
    # Label ID of the target class (sentiment)
    FeatureDefinition(feature_name='label_id', feature_type=FeatureTypeEnum.STRING),
    # reviews encoded with the BERT tokenizer
    FeatureDefinition(feature_name='input_ids', feature_type=FeatureTypeEnum.STRING),
    # original Review Text
    FeatureDefinition(feature_name='review_body', feature_type=FeatureTypeEnum.STRING),
    # train/validation/test Label
    FeatureDefinition(feature_name='split_type', feature_type=FeatureTypeEnum.STRING)
]
```

### Exercise 1

Create the feature group using the feature definitions defined above.

**Instructions:** Use the `FeatureGroup` function passing the defined above feature group name and the feature definitions.

```
[7]: from sagemaker.feature_store.feature_group import FeatureGroup

feature_group = FeatureGroup(
    name=..., # Feature Group name
    feature_definitions=..., # a List of Feature Definitions
    sagemaker_session=sess # SageMaker session
)

print(feature_group)

FeatureGroup(name='reviews-feature-group-1628645791', sagemaker_session=<sagemaker.session.Session object at 0x7f765a434090>, feature_definitions=[FeatureDefinition(feature_name='review_id', feature_type=<FeatureTypeEnum.STRING: 'String'>), FeatureDefinition(feature_name='date', feature_type=<FeatureTypeEnum.STRING: 'String'>), FeatureDefinition(feature_name='sentiment', feature_type=<FeatureTypeEnum.STRING: 'String'>), FeatureDefinition(feature_name='label_id', feature_type=<FeatureTypeEnum.STRING: 'String'>), FeatureDefinition(feature_name='input_ids', feature_type=<FeatureTypeEnum.STRING: 'String'>), FeatureDefinition(feature_name='review_body', feature_type=<FeatureTypeEnum.STRING: 'String'>), FeatureDefinition(feature_name='split_type', feature_type=<FeatureTypeEnum.STRING: 'String'>)])
```

You will use the defined Feature Group later in this lab, the actual creation of the Feature Group will take place in the processing job. Now let's move into the setup of the processing job to transform the dataset.

## 2. Transform the dataset

You will configure a SageMaker processing job to run a custom Python script to balance and transform the raw data into a format used by BERT model.

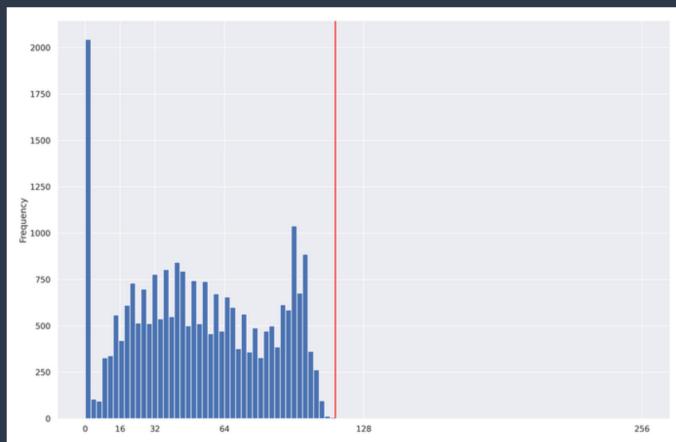
Set the transformation parameters including the instance type, instance count, and train/validation/test split percentages. For the purposes of this lab, you will use a relatively small instance type. Please refer to [this link](#) for additional instance types that may work for your use case outside of this lab.

You can also choose whether you want to balance the dataset or not. In this case, you will balance the dataset to avoid class imbalance in the target variable, `sentiment`.

Another important parameter of the model is the `max_seq_length`, which specifies the maximum length of the classified reviews for the RoBERTa model. If the sentence is shorter than the maximum length parameter, it will be padded. In another case, when the sentence is longer, it will be truncated from the right side.

Since a smaller `max_seq_length` leads to faster training and lower resource utilization, you want to find the smallest power-of-2 that captures 100% of our reviews. For this dataset, the 100th percentile is 115. However, it's best to stick with powers-of-2 when using BERT. So let's choose 128 as this is the smallest power-of-2 greater than 115. You will see below how the shorter sentences will be padded to a maximum length.

```
mean      52.512374
std       31.387048
min      1.000000
10%     10.000000
20%     22.000000
30%     32.000000
40%     41.000000
50%     51.000000
60%     61.000000
70%     73.000000
80%     88.000000
90%     97.000000
100%    115.000000
max      115.000000
```



```
[8]: processing_instance_type='ml.c5.xlarge'
processing_instance_count=1
train_split_percentage=0.90
validation_split_percentage=0.05
test_split_percentage=0.05
balance_dataset=True
max_seq_length=128
```

To balance and transform our data, you will use a scikit-learn-based processing job. This is essentially a generic Python processing job with scikit-learn pre-installed. You can specify the version of scikit-learn you wish to use. Also pass the SageMaker execution role, processing instance type and instance count.

```
[9]: from sagemaker.sklearn.processing import SKLearnProcessor
processor = SKLearnProcessor(
    framework_version='0.23-1',
    role=role,
    instance_type=processing_instance_type,
    instance_count=processing_instance_count,
    env={AWS_DEFAULT_REGION: region},
    max_runtime_in_seconds=7200
)
```

The processing job will be running the Python code from the file `src/prepare_data.py`. In the following exercise you will review the contents of the file and familiarize yourself with main parts of it.

## Exercise 2

1. Open the file `src/prepare_data.py`. Go through the comments to understand its content.
2. Find and review the `convert_to_bert_input_ids()` function, which contains the RoBERTa `tokenizer` configuration.
3. Complete method `encode_plus` of the RoBERTa `tokenizer`. Pass the `max_seq_length` as a value for the argument `max_length`. It defines a pad to a maximum length specified.
4. Save the file `src/prepare_data.py` (with the menu command File -> Save Python File).

*This cell will take approximately 1-2 minutes to run.*

```
[13]: import sys, importlib
sys.path.append('src/')

# import the 'prepare_data.py' module
import prepare_data

# reload the module if it has been previously loaded
if 'prepare_data' in sys.modules:
    importlib.reload(prepare_data)

input_ids = prepare_data.convert_to_bert_input_ids("this product is great!", max_seq_length)

updated_correctly = False

if len(input_ids) != max_seq_length:
    print('#####')
    print('Please check that the function \'convert_to_bert_input_ids\' in the file src/prepare_data.py is complete.')
    print('#####')
    raise Exception('Please check that the function \'convert_to_bert_input_ids\' in the file src/prepare_data.py is complete.')
else:
```

```
print('#####')
print('Updated correctly!')
print('#####')

updated_correctly = True
```

Review the results of tokenization for the given example ("this product is great!")

Launch the processing job with the custom script passing defined above parameters

```
[15]: from sagemaker.processing import ProcessingInput, ProcessingOutput

if (updated_correctly):

    processor.run(code='src/prepare_data.py',
                  inputs=[

                      ProcessingInput(source=raw_input_data_s3_uri,
                                      destination='/opt/ml/processing/input/data/',
                                      s3_data_distribution_type='ShardedByS3Key')
                  ],
                  outputs=[

                      ProcessingOutput(output_name='sentiment-train',
                                      source='/opt/ml/processing/output/sentiment/train',
                                      s3_upload_mode='EndOfJob'),
                      ProcessingOutput(output_name='sentiment-validation',
                                      source='/opt/ml/processing/output/sentiment/validation',
                                      s3_upload_mode='EndOfJob'),
                      ProcessingOutput(output_name='sentiment-test',
                                      source='/opt/ml/processing/output/sentiment/test',
                                      s3_upload_mode='EndOfJob')
                  ],
                  arguments=[--train-split-percentage, str(train_split_percentage),
                             '--validation-split-percentage', str(validation_split_percentage),
                             '--test-split-percentage', str(test_split_percentage),
                             '--balance-dataset', str(balance_dataset),
                             '--max-seq-length', str(max_seq_length),
                             '--feature-store-offline-prefix', str(feature_store_offline_prefix),
                             '--feature-group-name', str(feature_group_name)
                  ],
                  logs=True,
                  wait=False)

else:
    print('#####')
    print('Please update the code correctly above.')
    print('#####')
```

```

Job Name: sagemaker-skikit-learn-2021-08-11-02-29-20-395
Inputs: [{"InputName": "input-1", "AppManaged": False, "S3Input": {"S3Uri": 's3://dlai-practical-data-science/data/raw/'}, "LocalPath": '/opt/ml/input/processing/input/data/'}, {"S3DataType": 'S3Prefix', "S3InputMode": 'File', "S3DataDistributionType": 'ShardedByS3Key', "S3CompressionType": 'None'}], {"InputName": "code", "AppManaged": False, "S3Input": {"S3Uri": 's3://sagemaker-us-east-1-512205290621/sagemaker-skikit-learn-2021-08-11-02-29-20-395/input/code/prepare_data.py'}, "LocalPath": '/opt/ml/processing/input/code'}, {"S3DataType": 'S3Prefix', "S3InputMode": 'File', "S3DataDistributionType": 'FullyReplicated', "S3CompressionType": 'None'}]
Outputs: [{"OutputName": "sentiment-train", "AppManaged": False, "S3Output": {"S3Uri": 's3://sagemaker-us-east-1-512205290621/sagemaker-skikit-learn-2021-08-11-02-29-20-395/output/sentiment-train'}, "LocalPath": '/opt/ml/processing/output/sentiment-train', "S3UploadMode": 'EndOfJob'}}, {"OutputName": "sentiment-validation", "AppManaged": False, "S3Output": {"S3Uri": 's3://sagemaker-us-east-1-512205290621/sagemaker-skikit-learn-2021-08-11-02-29-20-395/output/sentiment-validation'}, "LocalPath": '/opt/ml/processing/output/sentiment-validation', "S3UploadMode": 'EndOfJob'}], {"OutputName": "sentiment-test", "AppManaged": False, "S3Output": {"S3Uri": 's3://sagemaker-us-east-1-512205290621/sagemaker-skikit-learn-2021-08-11-02-29-20-395/output/sentiment-test'}, "LocalPath": '/opt/ml/processing/output/sentiment-test', "S3UploadMode": 'EndOfJob'}]}

```

You can see the information about the processing jobs using the `describe` function. The result is in dictionary format. Let's pull the processing job name:

```
[16]: scikit_processing_job_name = processor.jobs[-1].describe()['ProcessingJobName']
       print('Processing job name: {}'.format(scikit_processing_job_name))
Processing job name: sagemaker-scikit-learn-2021-08-11-02-29-395
```

### Exercise 3

Pull the processing job status from the processing job description

**Instructions:** Print the keys of the processing job description dictionary, choose the one related to the status of the processing job and print the value of it.

```
[17]: print(processor.jobs[-1].describe().keys())
dict_keys(['ProcessingInputs', 'ProcessingOutputConfig', 'ProcessingJobName', 'ProcessingResources', 'StoppingCondition', 'AppSpecification', 'Environment', 'RoleArn', 'ProcessingJobArn', 'ProcessingJobStatus', 'LastModifiedTime', 'CreationTime', 'ResponseMetadata'])

[20]: ### BEGIN SOLUTION - DO NOT delete this comment for grading purposes
scikit_processing_job_status = processor.jobs[-1].describe()['ProcessingJobStatus'] # Replace None
### END SOLUTION - DO NOT delete this comment for grading purposes
print('Processing job status: {}'.format(scikit_processing_job_status))
```

- open the link
  - notice that you are in the section **Amazon SageMaker** -> **Processing jobs**
  - check the name of the processing job, its status and other available information

Review processing job

Wait for about 5 minutes to review the CloudWatch Logs. You may open the file `src/prepare_data.py` again and examine the outputs of the code in the CloudWatch Logs.

```
[22]: from IPython.core.display import display, HTML  
  
display(HTML('<b>Review <a target="blank" href="https://console.aws.amazon.com/cloudwatch/home?region={}#logStream:group=/aws/sagemaker/Processing' +</b>))
```

**Review CloudWatch logs after about 5 minutes**

After the completion of the processing job you can also review the output in the S3 bucket.

**Review S3 output data after the processing job has completed**

**Wait for the processing job to complete**

*This cell will take approximately 15 minutes to run.*

*Please wait until ^^ Processing Job ^^ completes above*

Inspect the transformed and balanced data in the S3 bucket.

```
[25]: processing_job_description = running_processor.describe()

output_config = processing_job_description['ProcessingOutputConfig']
for output in output_config['Outputs']:
    if output['OutputName'] == 'sentiment-train':
        processed_train_data_s3_uri = output['S3Output']['S3Uri']
    if output['OutputName'] == 'sentiment-validation':
        processed_validation_data_s3_uri = output['S3Output']['S3Uri']
    if output['OutputName'] == 'sentiment-test':
        processed_test_data_s3_uri = output['S3Output']['S3Uri']

print(processed_train_data_s3_uri)
print(processed_validation_data_s3_uri)
print(processed_test_data_s3_uri)

s3://sagemaker-us-east-1-512205290621/sagemaker-sklearn-2021-08-11-02-29-20-395/output/sentiment-train
s3://sagemaker-us-east-1-512205290621/sagemaker-sklearn-2021-08-11-02-29-20-395/output/sentiment-validation
s3://sagemaker-us-east-1-512205290621/sagemaker-sklearn-2021-08-11-02-29-20-395/output/sentiment-test
```

```
[26]: !aws s3 ls $processed_train_data_s3_uri/
        2021-08-11 02:43:15    4878637 part-algo-1-womens_clothing_ecommerce_reviews.ts
```

```
[27]: !aws s3 ls $processed_validation_data_s3_uri/
```

2021-08-11 02:45:19 203423 part-digital-womens\_clothing\_ecommerce\_reviews.ts

2021-08-11 02:43:16 273697 part-algo-1-womens\_clothing\_ecommerce\_reviews ts

**Copy the data into the folder balanced .**

```
[29]: !aws s3 cp ./processed_train_data_s3_uri/part-algo-1-womens_clothing_ecommerce_reviews.tsv ./balanced/sentiment-train/
!aws s3 cp ./processed_validation_data_s3_uri/part-algo-1-womens_clothing_ecommerce_reviews.tsv ./balanced/sentiment-validation/
!aws s3 cp ./processed_test_data_s3_uri/part-algo-1-womens_clothing_ecommerce_reviews.tsv ./balanced/sentiment-test/

download: s3://sagemaker-us-east-1-512285296621/sagemaker-skicit-learn-2021-08-11-02-29-20-395/output/sentiment-train/part-algo-1-womens_clothing_ecommerce_reviews.tsv to balanced/sentiment-train/part-algo-1-womens_clothing_ecommerce_reviews.tsv
download: s3://sagemaker-us-east-1-512285296621/sagemaker-skicit-learn-2021-08-11-02-29-20-395/output/sentiment-validation/part-algo-1-womens_clothing_ecommerce_reviews.tsv to balanced/sentiment-validation/part-algo-1-womens_clothing_ecommerce_reviews.tsv
download: s3://sagemaker-us-east-1-512285296621/sagemaker-skicit-learn-2021-08-11-02-29-20-395/output/sentiment-test/part-algo-1-womens_clothing_ecommerce_reviews.tsv to balanced/sentiment-test/part-algo-1-womens_clothing_ecommerce_reviews.tsv
```

**Review the training, validation and test data outputs.**

### 3. Query the Feature Store

In addition to transforming the data and saving in S3 bucket, the processing job populates the feature store with the transformed and balanced data. Let's query this data using Amazon Athena.

### 3.1. Export training, validation, and test datasets from the Feature Store

Here you will do the export only for the training dataset, as an example.

Use `athena_query()` function to create an Athena query for the defined above Feature Group. Then you can pull the table name of the Amazon Glue Data Catalog table which is auto-generated by Feature Store.

```
[33]: feature_store_query = feature_group.athena_query()

feature_store_table = feature_store_query.table_name

query_string = """
    SELECT date,
           review_id,
           sentiment,
           label_id,
           input_ids,
           review_body
      FROM "{}"
     WHERE split_type='train'
     LIMIT 5
""".format(feature_store_table)

print('Glue Catalog table name: {}'.format(feature_store_table))
print('Running query: {}'.format(query_string))

Glue Catalog table name: reviews-feature-group-1628645791-1628649409
Running query:
SELECT date,
       review_id,
       sentiment,
       label_id,
       input_ids,
       review_body
  FROM "reviews-feature-group-1628645791-1628649409"
 WHERE split_type='train'
 LIMIT 5
```

Configure the S3 location for the query results. This allows us to re-use the query results for future queries if the data has not changed. We can even share this S3 location between team members to improve query performance for common queries on data that does not change often.

```
[34]: output_s3_uri = 's3://{}{}/query_results/{}'.format(bucket, feature_store_offline_prefix
print(output_s3_uri)
```

Exercise 4

**Instructions:** Use `feature-store`, `query-pnp`, function passing the constructed above query string and the location of the output S3 bucket.

```
feature_store_query.run(  
    query_string=..., # query string
```

```

        output_location=... # Location of the output S3 bucket
    )

[37]: feature_store_query.run(
    ### BEGIN SOLUTION - DO NOT delete this comment for grading purposes
    query_string=query_string, # Replace None
    output_location='s3://sagemaker-us-east-1-512205290621/query_results/reviews-feature-store-1628645791/' # Replace None
    ### END SOLUTION - DO NOT delete this comment for grading purposes
)

feature_store_query.wait()

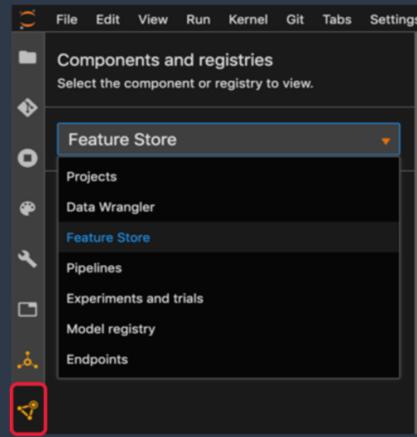
[38]: import pandas as pd
pd.set_option("max_colwidth", 100)

df_feature_store = feature_store_query.as_dataframe()
df_feature_store

```

	date	review_id	sentiment	label_id	input_ids	review_body
0	2021-08-11T02:37:14Z	12144	0	1	[0, 16587, 5, 299, 1521, 1437, 45, 98, 203, 5, 10199, 4, 939, 197, 348, 1166, 5, 299, 21, 727, 2...	Love the top design not so much the fabric. i should've read the top was 100% polyester. return...
1	2021-08-11T02:37:14Z	3419	0	1	[0, 713, 2551, 7, 422, 10, 410, 828, 650, 15, 162, 734, 55, 1026, 12365, 87, 15, 5, 1421, 11, 5...	This seemed to run a little bit small on me... more form fitting than on the model in the photo...
2	2021-08-11T02:37:14Z	8006	-1	0	[0, 100, 3584, 42, 3588, 13, 127, 11412, 4, 5, 8089, 32, 203, 55, 5676, 8, 4520, 172, 51, 2082, ...	I purchased this dress for my graduation. the colors are much more intense and bright than they ...
3	2021-08-11T02:37:14Z	13012	0	1	[0, 100, 770, 7, 657, 42, 299, 328, 15, 162, 24, 21, 7735, 352, 2233, 219, 4, 24, 21, 350, 1810, ...	I wanted to love this top! on me it was weirdly boxy. it was too wide and too short and it was ...
4	2021-08-11T02:37:14Z	19440	-1	0	[0, 713, 3588, 56, 98, 203, 801, 53, 5, 25490, 29716, 19750, 24, 4, 939, 1381, 15, 1836, 37863, ...	This dress had so much potential but the horizontal knit ruined it. i tried on size xx in grey ...

Review the Feature Store in SageMaker Studio



### 3.2. Export TSV from Feature Store

Save the output as a TSV file:

```

[39]: df_feature_store.to_csv('./feature_store_export.tsv',
                           sep='\t',
                           index=False,
                           header=True)

```

```

[40]: !head -n 5 ./feature_store_export.tsv

```

date	review_id	sentiment	label_id	input_ids	review_body
2021-08-11T02:37:14Z	12144	0	1	[0, 16587, 5, 299, 1521, 1437, 45, 98, 203, 5, 10199, 4, 939, 197, 348, 1166, 5, 299, 21, 727, 2...	Love the top design not so much the fabric. i should've read the top was 100% polyester. returning and finding a better quality top.
2021-08-11T02:37:14Z	3419	0	1	[0, 713, 2551, 7, 422, 10, 410, 828, 650, 15, 162, 734, 55, 1026, 12365, 87, 15, 5, 1421, 11, 5...	This seemed to run a little bit small on me... more form fitting than on the model in the photo. also it's thin. i'd be sold on it for sure if it were just a little longer too because i like the almost-lingerie-like look. the lace applique is pretty and if you purchase it for that then it really shows up better on the two lighter colors - it really gets lost all in black.
2021-08-11T02:37:14Z	8006	-1	0	[0, 100, 3584, 42, 3588, 13, 127, 11412, 4, 5, 8089, 32, 203, 55, 5676, 8, 4520, 172, 51, 2082, 80...	I purchased this dress for my graduation. the colors are much more intense and bright than they appear online. i'm kind of concerned that i look like a carnival tent. i like the style but the fabric feels cheap for the cost of the dress. the fabric is 100% synthetic. the fit is great. the dress is true to size.
2021-08-11T02:37:14Z	13012	0	1	[0, 100, 770, 7, 657, 42, 299, 328, 15, 162, 24, 21, 7735, 352, 2233, 219, 4, 24, 21, 350, 1810, 8, 358...	I wanted to love this top! on me it was weirdly boxy. it was too wide and too short and it was more sheer than i expected it. the fabric is lightweight but wrinkles really easy and loses shape pretty quickly. i also thought the fabric would be heavier and have more structure than it has. i love the scalloped detail at the arms but after about five minutes of wear this felt super crumpled and just kind of slopp...
2021-08-11T02:37:14Z	19440	-1	0	[0, 713, 3588, 56, 98, 203, 801, 53, 5, 25490, 29716, 19750, 24, 4, 939, 1381, 15, 1836, 37863, ...]	unfortunately it will be going back.

Upload TSV to the S3 bucket:

```

[41]: !aws s3 cp ./feature_store_export.tsv s3://#bucket/feature_store/feature_store_export.tsv
upload: ./feature_store_export.tsv to s3://sagemaker-us-east-1-512205290621/feature_store/feature_store_export.tsv

```

Check the file in the S3 bucket:

```

[42]: !aws s3 ls --recursive s3://#bucket/feature_store/feature_store_export.tsv
2021-08-11 02:51:28        4562 feature_store/feature_store_export.tsv

```

### 3.3. Check that the dataset in the Feature Store is balanced by sentiment

Now you can setup an Athena query to check that the stored dataset is balanced by the target class `sentiment`.

## Exercise 5

Write an SQL query to count the total number of the reviews per sentiment stored in the Feature Group.

**Instructions:** Pass the SQL statement of the form

```
SELECT category_column, COUNT(*) AS new_column_name
FROM table_name
GROUP BY category_column
```

into the variable `query_string_count_by_sentiment`. Here you would need to use the column `sentiment` and give a name `count_reviews` to the new column with the counts.

```
[43]: feature_store_query_2 = feature_group.athena_query()

# Replace all None
### BEGIN SOLUTION - DO NOT delete this comment for grading purposes
query_string_count_by_sentiment = """
SELECT sentiment, COUNT(*) AS count_reviews
FROM "{}"
GROUP BY sentiment
""".format(feature_store_table)
### END SOLUTION - DO NOT delete this comment for grading purposes
```

## Exercise 6

Query the feature store.

**Instructions:** Use `run` function of the Feature Store query, passing the new query string `query_string_count_by_sentiment`. The output S3 bucket will remain unchanged. You can follow the example above.

```
[44]: feature_store_query_2.run(
    ### BEGIN SOLUTION - DO NOT delete this comment for grading purposes
    query_string=query_string_count_by_sentiment, # Replace None
    output_location='s3://sagemaker-us-east-1-512205290621/query_results/reviews-feature-store-1628645791/' # Replace None
    ### END SOLUTION - DO NOT delete this comment for grading purposes
)

feature_store_query_2.wait()

df_count_by_sentiment = feature_store_query_2.as_dataframe()
df_count_by_sentiment
```

	sentiment	count_reviews
0	0	2051
1	-1	2051
2	1	2051

## Exercise 7

Visualize the result of the query in the bar plot, showing the count of the reviews by sentiment value.

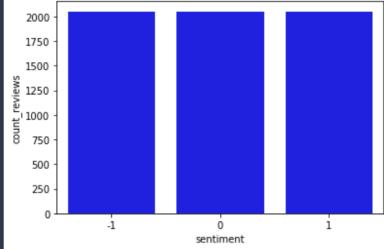
**Instructions:** Pass the resulting data frame `df_count_by_sentiment` into the `barplot` function of the `seaborn` library.

```
sns.barplot(
    data=...,
    x='...',
    y='...',
    color="blue"
)

[46]: import seaborn as sns

sns.barplot(
    ### BEGIN SOLUTION - DO NOT delete this comment for grading purposes
    data=df_count_by_sentiment, # Replace None
    x='sentiment', # Replace None
    y='count_reviews', # Replace None
    ### END SOLUTION - DO NOT delete this comment for grading purposes
    color="blue"
)
```

```
[46]: <matplotlib.axes._subplots.AxesSubplot at 0x7f75f16f4a90>
```



Upload the notebook and `prepare_data.py` file into S3 bucket for grading purposes.

**Note:** you may need to save the file before the upload.

```
[47]: aws s3 cp ./C2_W1_Assignment.ipynb s3://$bucket/C2_W1_Assignment_Learner.ipynb
aws s3 cp ./src/prepare_data.py s3://$bucket/src/C2_W1_prepare_data_Learner.py

upload: ./C2_W1_Assignment.ipynb to s3://sagemaker-us-east-1-512205290621/C2_W1_Assignment_Learner.ipynb
upload: src/prepare_data.py to s3://sagemaker-us-east-1-512205290621/src/C2_W1_prepare_data_Learner.py
```

Please go to the main lab window and click on `Submit` button (see the `Finish the lab` section of the instructions).

```
[ ]:
```



1 8 1 Git: idle Python 3 (Data Science) | Idle Kernel: Idle | Instance MEM

Mode: Command ⌘ Ln 1, Col 1 C2\_W1\_Assignment.ipynb