

Amazon SageMaker Studio

File Edit View Run Kernel Git Tabs Settings Help

Feedback

Terminal 1 C2\_W1\_Assignment.ipynb prepare\_data.py

```

1 ##### Import required modules #####
2 #####
3 #####
4 from sklearn.model_selection import train_test_split
5 from sklearn.utils import resample
6 import functools
7 import multiprocessing
8
9 from datetime import datetime
10 from time import gtime, strftime, sleep
11
12 import pandas as pd
13 import argparse
14 import subprocess
15 import sys
16 import os
17 import re
18 import collections
19 import json
20 import csv
21 import glob
22 from pathlib import Path
23 import time
24 import boto3
25
26 subprocess.check_call([sys.executable, "-m", "conda", "install", "-c", "pytorch", "pytorch==1.6.0", "-y"])
27
28 subprocess.check_call([sys.executable, "-m", "conda", "install", "-c", "conda-forge", "transformers==3.5.1", "-y"])
29 from transformers import RobertaTokenizer
30
31 subprocess.check_call([sys.executable, '-m', 'pip', 'install', 'sagemaker==2.35.0'])
32 import sagemaker
33
34 from sagemaker.session import Session
35 from sagemaker.feature_store.feature_group import FeatureGroup
36 from sagemaker.feature_store.feature_definition import (
    FeatureDefinition,
    FeatureTypeEnum,
)
37
38 #####
39 ###### Setup environmental variables #####
40
41 #####
42 #####
43 #####
44
45 region = os.environ['AWS_DEFAULT_REGION']
46 sts = boto3.Session(region_name=region).client(service_name='sts', region_name=region)
47 iam = boto3.Session(region_name=region).client(service_name='iam', region_name=region)
48 featurestore_runtime = boto3.Session(region_name=region).client(service_name='sagemaker-featurestore-runtime', region_name=region)
49 sm = boto3.Session(region_name=region).client(service_name='sagemaker', region_name=region)
50

```

1 8 1 Git: refreshing... Python

Ln 188, Col 36 Spaces: 4 prepare\_data.py

```

53 assumed_role_name = assumed_roleArn.split('/')[-2]
54 get_role_response = iam.get_role(RoleName=assumed_role_name)
55 role = get_role_response['Role']['Arn']
56 bucket = sagemaker.Session().default_bucket()
57
58 sagemaker_session = sagemaker.Session(boto_session=boto3.Session(region_name=region),
59                                         sagemaker_client=sm,
60                                         sagemaker_featurestore_runtime_client=featurestore_runtime)
61
62 # List of sentiment classes: -1 - negative; 0 - neutral; 1 - positive
63 classes = [-1, 0, 1]
64
65 # Label IDs of the target class (sentiment) setup as a dictionary
66 classes_map = {
67     -1: 0,
68     0: 1,
69     1: 2
70 }
71
72 # Tokenization model
73 PRE_TRAINED_MODEL_NAME = 'roberta-base'
74
75 # Create the tokenizer to use based on pre trained model
76 tokenizer = RobertaTokenizer.from_pretrained(PRE_TRAINED_MODEL_NAME)
77
78 #####
79 #####
80 #####
81 # Functions which can be considered as tools
82
83 def cast_object_to_string(data_frame):
84     for label in data_frame.columns:
85         if data_frame.dtypes[label] == 'object':
86             data_frame[label] = data_frame[label].astype("str").astype("string")
87     return data_frame
88
89
90 def wait_for_feature_group_creation_complete(feature_group):
91     try:
92         status = feature_group.describe().get("FeatureGroupStatus")
93         print(f'Feature Group status: {status}')
94         while status == "Creating":
95             print("Waiting for Feature Group Creation")
96             time.sleep(5)
97
98         if status != "Created":
99             print(f'Feature Group status: {status}')
100            raise RuntimeError(f"Failed to create feature group {feature_group.name}")
101        print(f"FeatureGroup {feature_group.name} successfully created.")
102    except:
103        print('No feature group created yet.')
104
105 def list_arg(raw_value):
106     """argparse type for a list of strings"""
107     return str(raw_value).split(',')
108
109 def to_sentiment(star_rating):
110     if star_rating in {1, 2}: # negative
111         return -1
112     if star_rating == 3: # neutral
113         return 0
114     if star_rating in {4, 5}: # positive
115         return 1
116
117 #####
118 #####
119 ###### Create or load Feature Group #####

```

Ln 188, Col 36 Spaces: 4 prepare\_data.py

Ln 188, Col 36 Spaces: 4 prepare\_data.py

Ln 188, Col 36 Spaces: 4 prepare\_data.py

In 188 Col 36 Spaces: 4 prepare\_data.py

```

243     default=128
244 )
245 parser.add_argument('--feature-store-offline-prefix', type=str,
246     default=None,
247 )
248 parser.add_argument('--feature-group-name', type=str,
249     default=None,
250 )
251
252     return parser.parse_args()
253
254 #####
255 ##### Processing functions #####
256 #####
257
258 def _preprocess_file(file,
259     balance_dataset,
260     max_seq_length,
261     prefix,
262     feature_group_name):
263
264     print('file {}'.format(file))
265     print('balance_dataset {}'.format(balance_dataset))
266     print('max_seq_length {}'.format(max_seq_length))
267     print('prefix {}'.format(prefix))
268     print('feature_group_name {}'.format(feature_group_name))
269
270     # Create a feature group
271     # The Feature Group that was set in the main notebook cannot be passed here - it will be used later in the notebook for other purposes
272     # you need to create a Feature Group with the same Feature Definitions within the processing job
273     feature_group = create_or_load_feature_group(prefix, feature_group_name)
274
275     filename_without_extension = Path(Path(file).stem).stem
276
277     # read file
278     df = pd.read_csv(file,
279         index_col=0)

```

Ln 188, Col 36 Spaces: 4 prepare\_data.py

```

282 df = df.dropna()
283 df = df.reset_index(drop=True)
284 print('Shape of dataframe {}'.format(df.shape))
285
286 # convert star rating into sentiment
287 df['sentiment'] = df['Rating'].apply(lambda star_rating: to_sentiment(star_rating=star_rating))
288 print('Shape of dataframe with sentiment {}'.format(df.shape))
289
290 # convert sentiment (-1, 0, 1) into Label_id (0, 1, 2)
291 df['label_id'] = df['sentiment'].apply(lambda sentiment: classes_map[sentiment])
292 print('df[label_id] after using classes_map: {}'.format(df['label_id']))
293
294 df['input_ids'] = df['Review Text'].apply(lambda review: convert_to_bert_input_ids(review, max_seq_length))
295 print('df[input_ids] after calling convert_to_bert_input_ids: {}'.format(df['input_ids']))
296
297 # convert the index into a review_id
298 df.reset_index(inplace=True)
299 df = df.rename(columns = {'index': 'review_id',
300                     'Review Text': 'review_body'})
301
302 # drop all columns except the following:
303 df = df[['review_id', 'sentiment', 'label_id', 'input_ids', 'review_body']]
304 df = df.reset_index(drop=True)
305
306 print(df.head())
307
308 print('Shape of dataframe after dropping columns {}'.format(df.shape))
309
310 # balance the dataset by sentiment down to the minority class
311 if balance_dataset:
312
313     df_unbalanced_grouped_by = df.groupbyby('sentiment')
314     df_balanced = df.unbalanced_grouped_by.apply(lambda x: x.sample(df_unbalanced_grouped_by.size().min()).reset_index(drop=True))
315     print('Shape of balanced df: {}'.format(df_balanced.shape))
316
317     print(df_balanced['sentiment'].head())
318
319     df = df_balanced
320
321     # adding timestamp as date column
322     timestamp = datetime.now().strftime("%Y-%m-%dT%H:%M:%S%z")
323     print(timestamp)
324
325     df['date'] = timestamp
326
327

```

Ln 188, Col 36 Spaces: 4 prepare\_data.py

```

328 # split dataset
329 print('Shape of dataframe before splitting {}'.format(df.shape))
330
331 print('train split percentage {}'.format(args.train_split_percentage))
332 print('validation split percentage {}'.format(args.validation_split_percentage))
333 print('test split percentage {}'.format(args.test_split_percentage))
334
335 holdout_percentage = 1.00 - args.train_split_percentage
336 print('holdout percentage {}'.format(holdout_percentage))
337 df_train, df_holdout = train_test_split(df,
338                                         test_size=holdout_percentage,
339                                         stratify=df['sentiment'])
340
341 test_holdout_percentage = args.test_split_percentage / holdout_percentage
342 print('test holdout percentage {}'.format(test_holdout_percentage))
343 df_validation, df_test = train_test_split(df_holdout,
344                                         test_size=test_holdout_percentage,
345                                         stratify=df_holdout['sentiment'])
346
347 df_train = df_train.reset_index(drop=True)
348 df_validation = df_validation.reset_index(drop=True)
349 df_test = df_test.reset_index(drop=True)
350
351 print('Shape of train dataframe {}'.format(df_train.shape))
352 print('Shape of validation dataframe {}'.format(df_validation.shape))
353 print('Shape of test dataframe {}'.format(df_test.shape))
354
355 train_data = '{}/sentiment/train'.format(args.output_data)
356 validation_data = '{}/sentiment/validation'.format(args.output_data)
357 test_data = '{}/sentiment/test'.format(args.output_data)
358
359 ## write TSV Files
360 df_train.to_csv('{}/part-{}-{}.tsv'.format(train_data, args.current_host, filename_without_extension), sep='\t', index=False)
361 df_validation.to_csv('{}/part-{}-{}.tsv'.format(validation_data, args.current_host, filename_without_extension), sep='\t', index=False)
362 df_test.to_csv('{}/part-{}-{}.tsv'.format(test_data, args.current_host, filename_without_extension), sep='\t', index=False)
363
364 # dataframe
365 df_train.head()
366 df_validation.head()

```

```
367     df_test.head()
368
369     column_names = ['review_id', 'sentiment', 'date', 'label_id', 'input_ids', 'review_body']
370
371     df_train_records = df_train[column_names]
```

Ln 188, Col 36 Spaces: 4 prepare\_data.py

```
374
375 df_validation_records = df_validation[column_names]
376 df_validation_records['split_type'] = 'validation'
377 df_validation_records.head()
378
379 df_test_records = df_test[column_names]
380 df_test_records['split_type'] = 'test'
381 df_test_records.head()
382
383 # add record to feature store
384 df_fs_train_records = cast_object_to_string(df_train_records)
385 df_fs_validation_records = cast_object_to_string(df_validation_records)
386 df_fs_test_records = cast_object_to_string(df_test_records)
387
388 print('Ingesting features...')
389 feature_group.ingest(
390     data_frame=df_fs_train_records, max_workers=3, wait=True
391 )
392 feature_group.ingest(
393     data_frame=df_fs_validation_records, max_workers=3, wait=True
394 )
395 feature_group.ingest(
396     data_frame=df_fs_test_records, max_workers=3, wait=True
397 )
398
399 offline_store_status = None
400 while offline_store_status != 'Active':
401     try:
402         offline_store_status = feature_group.describe()['OfflineStoreStatus']['Status']
403     except:
404         pass
405     print('Offline store status: {}'.format(offline_store_status))
406     sleep(15)
407 print('...features ingested!')
408
409
410 def process(args):
411     # (1) Create a local file if it's not supported by the
412     # (2) Create a local file if it's not supported by the
```

Ln 188, Col 36 Spaces: 4 prepare\_data.py

```
validation_data = '{}/sentiment/validation'.format(args.output_data)
test_data = '{}/sentiment/test'.format(args.output_data)

# partial functions allow to derive a function with some parameters to a function with fewer parameters
# and fixed values set for the more limited function
# here '_preprocess_file' will be more limited function than '_preprocess_file' with fixed values for some parameters
preprocess_file = functools.partial(_preprocess_file,
                                     balance_dataset=args.balance_dataset,
                                     max_seq_length=args.max_seq_length,
                                     prefix=args.feature_store_offline_prefix,
                                     feature_group_name=args.feature_group_name)

input_files = glob.glob('{}/*.csv'.format(args.input_data))

num_cpus = multiprocessing.cpu_count()
print('num_cpus {}'.format(num_cpus))

p = multiprocessing.Pool(num_cpus)
p.map(preprocess_file, input_files)

print('Listing contents of {}'.format(preprocessed_data))
dirs_output = os.listdir(preprocessed_data)
for file in dirs_output:
    print(file)

print('Listing contents of {}'.format(train_data))
dirs_output = os.listdir(train_data)
for file in dirs_output:
    print(file)

print('Listing contents of {}'.format(validation_data))
dirs_output = os.listdir(validation_data)
for file in dirs_output:
    print(file)

print('Listing contents of {}'.format(test_data))
dirs_output = os.listdir(test_data)
for file in dirs_output:
    print(file)

print('Complete')
#####
#####
```

Ln 188, Col 36 Spaces: 4 prepare\_data.py

```
466 if __name__ == "__main__":
467
468     args = parse_args()
469     print('Loaded argument')
470     print(args)
471
472     print('Environment var')
473     print(os.environ)
474
475     process(args)
```

