

Cleaning Data with PySpark (Dataframe Details)

Defining a schema

Creating a defined schema helps with data quality and import performance. As mentioned during the lesson, we'll create a simple schema to read in the following columns:

- Name
- Age
- City

The `Name` and `City` columns are `StringType()` and the `Age` column is an `IntegerType()`.

```
1 # Import the pyspark.sql.types library
2 from pyspark.sql.types import *
3
4 # Define a new schema using the StructType method
5 people_schema = StructType([
6     # Define a StructField for each field
7     StructField('name', StringType(), False),
8     StructField('age', IntegerType(), False),
9     StructField('city', StringType(), False)
10 ])
```

Use lazy processing

Lazy processing operations will usually return in about the same amount of time regardless of the actual quantity of data. Remember that this is due to Spark not performing any transformations until an action is requested.

For this exercise, we'll be defining a Data Frame (`aa_dfw_df`) and add a couple transformations. Note the amount of time required for the transformations to complete when defined vs when the data is actually queried. These differences may be short, but they will be noticeable. When working with a full Spark cluster with larger quantities of data the difference will be more apparent.

```

1 # Load the CSV file
2 aa_dfw_df = spark.read.format('csv').options(Header=True).load('AA_DFW_2018.csv.gz')
3
4 # Add the airport column using the F.lower() method
5 aa_dfw_df = aa_dfw_df.withColumn('airport', F.lower(aa_dfw_df['Destination Airport']))
6
7 # Drop the Destination Airport column
8 aa_dfw_df = aa_dfw_df.drop(aa_dfw_df['Destination Airport'])
9
10 # Show the DataFrame
11 aa_dfw_df.show()

```

Saving a DataZFrame in Parquet format

When working with Spark, you'll often start with CSV, JSON, or other data sources. This provides a lot of flexibility for the types of data to load, but it is not an optimal format for Spark. The

`Parquet` format is a columnar data store, allowing Spark to use *predicate pushdown*. This means Spark will only process the data necessary to complete the operations you define versus reading the entire dataset. This gives Spark more flexibility in accessing the data and often drastically improves performance on large datasets.

In this exercise, we're going to practice creating a new Parquet file and then process some data from it.

The `spark` object and the `df1` and `df2` DataFrames have been setup for you.

```

1 # View the row count of df1 and df2
2 print("df1 Count: %d" % df1.count())
3 print("df2 Count: %d" % df2.count())
4
5 # Combine the DataFrames into one
6 df3 = df1.union(df2)
7
8 # Save the df3 DataFrame in Parquet format
9 df3.write.parquet('AA_DFW_ALL.parquet', mode='overwrite')
10
11 # Read the Parquet file into a new DataFrame and run a count
12 print(spark.read.parquet('AA_DFW_ALL.parquet').count())

```

SQL and Parquest

Parquet files are perfect as a backing data store for SQL queries in Spark. While it is possible to run the same queries directly via Spark's Python functions, sometimes it's easier to run SQL queries alongside the Python options.

For this example, we're going to read in the Parquet file we created in the last exercise and register it as a SQL table. Once registered, we'll run a quick query against the table (aka, the Parquet file).

The `spark` object and the `AA_DFW_ALL.parquet` file are available for you automatically.

```
1 # Read the Parquet file into flights_df
2 flights_df = spark.read.parquet('AA_DFW_ALL.parquet')
3
4 # Register the temp table
5 flights_df.createOrReplaceTempView('flights')
6
7 # Run a SQL query of the average flight duration
8 avg_duration = spark.sql('SELECT avg(flight_duration) from flights').collect()[0]
9 print('The average flight time is: %d' % avg_duration)
```

Dog Parsing

You've done a considerable amount of cleanup on the initial dataset, but now need to analyze the data a bit deeper. There are several questions that have now come up about the type of dogs seen in an image and some details regarding the images. You realize that to answer these questions, you need to process the data into a specific type. Before you can use it, you'll need to create a schema / type to represent the dog details.

The `joined_df` DataFrame is as you last defined it, and the `pyspark.sql.types` have all been imported.

```
1 # Select the dog details and show 10 untruncated rows
2 print(joined_df.select('dog_list').show(truncate=False))
3
4 # Define a schema type for the details in the dog list
5 DogType = StructType([
6     StructField("breed", StringType(), False),
7     StructField("start_x", IntegerType(), False),
8     StructField("start_y", IntegerType(), False),
9     StructField("end_x", IntegerType(), False),
10    StructField("end_y", IntegerType(), False)
11 ])
```

Per image count

Your next task in building a data pipeline for this dataset is to create a few analysis oriented columns. You've been asked to calculate the number of dogs found in each image based on your `dog_list` column created earlier. You have also created the `DogType` which will allow better

parsing of the data within some of the data columns.

The `joined_df` is available as you last defined it, and the `DogType` structtype is defined. `pyspark.sql.functions` is available under the `F` alias.

```
1 # Create a function to return the number and type of dogs as a tuple
2 def dogParse(doglist):
3     dogs = []
4     for dog in doglist:
5         (breed, start_x, start_y, end_x, end_y) = dog.split(',')
6         dogs.append((breed, int(start_x), int(start_y), int(end_x), int(end_y)))
7     return dogs
8
9 # Create a UDF
10 udfDogParse = F.udf(dogParse, ArrayType(DogType))
11
12 # Use the UDF to list of dogs and drop the old column
13 joined_df = joined_df.withColumn('dogs', udfDogParse('dog_list')).drop('dog_list')
14
15 # Show the number of dogs in the first 10 rows
16 joined_df.select(F.size('dogs')).show(10)
```

Percentage dog pixels

The final task for parsing the dog annotation data is to determine the percentage of pixels in each image that represents a dog (or dogs). You'll need to use the various techniques you've learned in this course to help calculate this information and add it as columns for later analysis.

To calculate the percentage of pixels, first calculate the total number of pixels representing each dog then sum them for the image. You can calculate the bounding box with the formula:

$(X_{end} - X_{start}) * (Y_{end} - Y_{start})$

NOTE: You can ignore the possibility of overlapping bounding boxes in this instance.

For the percentage, calculate the total number of "dog" pixels divided by the total size of the image, multiplied by 100.

The `joined_df` DataFrame is as you last used it. `pyspark.sql.functions` is aliased to `F`.

```
1 # Define a UDF to determine the number of pixels per image
2 def dogPixelCount(doglist):
3     totalpixels = 0
4     for dog in doglist:
5         totalpixels += (dog[3] - dog[1])*(dog[4] - dog[2])
6     return totalpixels
7
```

```
8 # Define a UDF for the pixel count
9 udfDogPixelCount = F.udf(dogPixelCount, IntegerType())
10 joined_df = joined_df.withColumn('dog_pixels', udfDogPixelCount('dogs'))
11
12 # Create a column representing the percentage of pixels
13 joined_df = joined_df.withColumn('dog_percent', (joined_df.dog_pixels / (joined_df.width *
14 joined_df.height)) * 100)
15
16 # Show the first 10 annotations with more than 60% dog
17 joined_df.filter(joined_df.dog_percent > 60).show(10)
```