```
1 !pip install pyspark
```

```
Collecting pyspark
  Downloading pyspark-3.5.0.tar.gz (316.9 MB)
                                                316.9/316.9 MB 4.4 MB/s eta 0:0
  Preparing metadata (setup.py) ... done
Requirement already satisfied: py4j==0.10.9.7 in /usr/local/lib/python3.10/di
Building wheels for collected packages: pyspark
  Building wheel for pyspark (setup.py) ... done
  Created wheel for pyspark: filename=pyspark-3.5.0-py2.py3-none-any.whl size
  Stored in directory: /root/.cache/pip/wheels/41/4e/10/c2cf2467f71c678cfc8a6
Successfully built pyspark
Installing collected packages: pyspark
Successfully installed pyspark-3.5.0
```

```
1 !pip install sparkxgb
```

```
Collecting sparkxgb
  Downloading sparkxgb-0.1.tar.gz (3.6 kB)
  Preparing metadata (setup.py) ... done
Collecting pyspark==3.1.1 (from sparkxgb)
  Downloading pyspark-3.1.1.tar.gz (212.3 MB)
                                                212.3/212.3 MB 5.8 MB/s eta 0:0
  Preparing metadata (setup.py) ... done
Collecting py4j==0.10.9 (from pyspark==3.1.1->sparkxgb)
  Downloading py4j-0.10.9-py2.py3-none-any.whl (198 kB)
                                                198.6/198.6 kB 26.3 MB/s eta 0:
Building wheels for collected packages: sparkxgb, pyspark
  Building wheel for sparkxgb (setup.py) ... done
  Created wheel for sparkxgb: filename=sparkxgb-0.1-py3-none-any.whl size=562
  Stored in directory: /root/.cache/pip/wheels/b7/0c/a1/786408e13056fabeb8a72
  Building wheel for pyspark (setup.py) ... done
  Created wheel for pyspark: filename=pyspark-3.1.1-py2.py3-none-any.whl size
  Stored in directory: /root/.cache/pip/wheels/a0/3f/72/8efd988f9ae041f051c75
Successfully built sparkxgb pyspark
Installing collected packages: py4j, pyspark, sparkxgb
  Attempting uninstall: py4j
    Found existing installation: py4j 0.10.9.7
    Uninstalling py4j-0.10.9.7:
      Successfully uninstalled py4j-0.10.9.7
  Attempting uninstall: pyspark
    Found existing installation: pyspark 3.5.0
    Uninstalling pyspark-3.5.0:
      Successfully uninstalled pyspark-3.5.0
Successfully installed py4j-0.10.9 pyspark-3.1.1 sparkxgb-0.1
```

```
1 from pyspark.sql import SparkSession
```

```python
1 from pyspark.sql import SparkSession
2 from pyspark.sql.types import FloatType
3 from pyspark.sql.functions import col, isnan
4 from pyspark.ml.feature import VectorAssembler, PolynomialExpansion
5 from pyspark.ml.regression import RandomForestRegressor, LinearRegression, GBTR
6 from pyspark.ml import Pipeline
7 from pyspark.ml.evaluation import RegressionEvaluator
8 from sparkxgb import XGBoostRegressor
9
10 # Initialize Spark session
11 spark = SparkSession.builder.appName("FlightDelayPred").getOrCreate()
12
13 # Load the data
14 data_path = "/content/drive/MyDrive/Datasets/1987.csv"  # Replace with your fil
15 df = spark.read.csv(data_path, header=True, inferSchema=True)
16
17 # Data type casting and handling missing values
18 df = df.withColumn("DepTime", col("DepTime").cast(FloatType()))
19 df = df.withColumn("Distance", col("Distance").cast(FloatType()))
20 df = df.withColumn("CRSDepTime", col("CRSDepTime").cast(FloatType()))
21 df = df.withColumn("Month", col("Month").cast(FloatType()))
22 df = df.withColumn("ArrDelay", col("ArrDelay").cast(FloatType()))
23 df = df.withColumn("DepDelay", col("DepDelay").cast(FloatType()))
24 df = df.withColumn("CRSElapsedTime", col("CRSElapsedTime").cast(FloatType()))
25 df = df.withColumn("DayofMonth", col("DayofMonth").cast(FloatType()))
26 df = df.withColumn("FlightNum", col("FlightNum").cast(FloatType()))
27
28 # Remove rows with null or NaN values in target column
29 df = df.filter(df.ArrDelay.isNotNull() & (~isnan(df.ArrDelay)))
30
31 # Selecting the features and target variable
32 features = ['Month', 'CRSDepTime', 'DepTime', 'Distance', 'DepDelay', 'CRSElaps
33 target = 'ArrDelay'
34
35 # VectorAssembler to combine feature columns into a single vector column
36 assembler = VectorAssembler(inputCols=features, outputCol="features", handleInv
37
38 # Polynomial Expansion for degree 5
39 polyExpansion = PolynomialExpansion(degree=5, inputCol="features", outputCol="p
40
41 # Define Linear Regression model for polynomial regression
42 poly_lr = LinearRegression(featuresCol="polyFeatures", labelCol=target)
43
44 # Update pipeline for Polynomial Regression
45 poly_pipeline = Pipeline(stages=[assembler, polyExpansion, poly_lr])
46
```

```python
47  # Define the models
48  rf = RandomForestRegressor(featuresCol="features", labelCol=target)
49  lr = LinearRegression(featuresCol="features", labelCol=target)
50  gbt = GBTRegressor(featuresCol="features", labelCol=target)
51
52  # Pipelines for the models
53  rf_pipeline = Pipeline(stages=[assembler, rf])
54  lr_pipeline = Pipeline(stages=[assembler, lr])
55  gbt_pipeline = Pipeline(stages=[assembler, gbt])
56
57  # Evaluate each model using a pipeline
58  models = [rf_pipeline, lr_pipeline, gbt_pipeline, poly_pipeline]
59
60  # Update the models list to include the XGBoost pipeline
61  models = [rf_pipeline, lr_pipeline, gbt_pipeline, poly_pipeline]
62
63  for model in models:
64      # Train the model
65      trained_model = model.fit(df)
66
67      # Make predictions
68      predictions = trained_model.transform(df)
69
70      # Evaluate the model for RMSE
71      rmse_evaluator = RegressionEvaluator(labelCol=target, predictionCol="predic
72      rmse = rmse_evaluator.evaluate(predictions)
73
74      # Evaluate the model for R2
75      r2_evaluator = RegressionEvaluator(labelCol=target, predictionCol="predicti
76      r2 = r2_evaluator.evaluate(predictions)
77
78      # Print the RMSE and R2
79      if isinstance(model.getStages()[-1], LinearRegression) and len(model.getSta
80          model_name = "Polynomial Linear Regression"
81      else:
82          model_name = model.getStages()[-1].__class__.__name__
83      print(f"{model_name} - Root Mean Squared Error (RMSE): {rmse}, R2: {r2}")
84
85  # Stop the Spark session
86  spark.stop()
```

RandomForestRegressor — Root Mean Squared Error (RMSE): 16.45410836337029, R2
LinearRegression — Root Mean Squared Error (RMSE): 15.514098363890556, R2: 0.
GBTRegressor — Root Mean Squared Error (RMSE): 14.231853301671865, R2: 0.6959
Polynomial Linear Regression — Root Mean Squared Error (RMSE): 13.17351924772

1