

▼ Build a Classification Model with Spark with iris dataset

```

1 from pyspark.sql import SparkSession
2 from pyspark.ml.feature import StringIndexer, VectorAssembler, IndexToString
3 from pyspark.ml.classification import LogisticRegression
4 from pyspark.ml.evaluation import MulticlassClassificationEvaluator
5 import urllib.request
6
7 # Initialize Spark session
8 spark = SparkSession.builder.appName("IrisClassification").getOrCreate()
9
10 # Download the dataset
11 data_url = "https://archive.ics.uci.edu/ml/machine-learning-databases/iris/iris.data"
12 iris_path = "/tmp/iris.csv"
13 urllib.request.urlretrieve(data_url, iris_path)
14
15 # Load dataset into DataFrame
16 df = spark.read.csv(iris_path, inferSchema=True)\
17     .toDF("sepal_length", "sepal_width", "petal_length", "petal_width", "species")
18
19 # Step 1: Encode species as numerical label
20 indexer = StringIndexer(inputCol="species", outputCol="label")
21 indexer_model = indexer.fit(df)
22 df = indexer_model.transform(df)
23
24 # Step 2: Assemble features into a single vector
25 assembler = VectorAssembler(
26     inputCols=["sepal_length", "sepal_width", "petal_length", "petal_width"],
27     outputCol="features"
28 )
29 df = assembler.transform(df)
30
31 # Step 3: Train/test split
32 train, test = df.randomSplit([0.7, 0.3], seed=42)
33
34 # Step 4: Train logistic regression model
35 lr = LogisticRegression(featuresCol="features", labelCol="label")
36 model = lr.fit(train)
37
38 # Step 5: Make predictions
39 predictions = model.transform(test)
40
41 # Step 6: Convert prediction index back to species name
42 label_converter = IndexToString(
43     inputCol="prediction",
44     outputCol="predicted_species",
45     labels=indexer_model.labels
46 )
47 predictions = label_converter.transform(predictions)
48
49 # Step 7: Display predictions
50 predictions.select("features", "label", "prediction", "species", "predicted_species")\
51     .show(20, truncate=False)
52
53 # Step 8: Evaluate accuracy
54 evaluator = MulticlassClassificationEvaluator(
55     labelCol="label", predictionCol="prediction", metricName="accuracy"
56 )
57 accuracy = evaluator.evaluate(predictions)
58 print(f"Test Accuracy: {accuracy:.2f}")

```

```

+-----+-----+-----+-----+
|features|label|prediction|species|predicted_species|
+-----+-----+-----+-----+
|[4.4,3.0,1.3,0.2]|0.0|0.0|Iris-setosa|Iris-setosa|
|[4.6,3.2,1.4,0.2]|0.0|0.0|Iris-setosa|Iris-setosa|
|[4.6,3.6,1.0,0.2]|0.0|0.0|Iris-setosa|Iris-setosa|
|[4.7,3.2,1.3,0.2]|0.0|0.0|Iris-setosa|Iris-setosa|
|[4.8,3.1,1.6,0.2]|0.0|0.0|Iris-setosa|Iris-setosa|
|[4.8,3.4,1.6,0.2]|0.0|0.0|Iris-setosa|Iris-setosa|
|[4.8,3.4,1.9,0.2]|0.0|0.0|Iris-setosa|Iris-setosa|
|[4.9,3.1,1.5,0.1]|0.0|0.0|Iris-setosa|Iris-setosa|
|[4.9,3.1,1.5,0.1]|0.0|0.0|Iris-setosa|Iris-setosa|
|[5.0,2.3,3.3,1.0]|1.0|1.0|Iris-versicolor|Iris-versicolor|
|[5.0,3.0,1.6,0.2]|0.0|0.0|Iris-setosa|Iris-setosa|
|[5.0,3.4,1.6,0.4]|0.0|0.0|Iris-setosa|Iris-setosa|
|[5.0,3.5,1.3,0.3]|0.0|0.0|Iris-setosa|Iris-setosa|
|[5.0,3.5,1.6,0.6]|0.0|0.0|Iris-setosa|Iris-setosa|
|[5.1,2.5,3.0,1.1]|1.0|1.0|Iris-versicolor|Iris-versicolor|
|[5.1,3.4,1.5,0.2]|0.0|0.0|Iris-setosa|Iris-setosa|
|[5.1,3.5,1.4,0.2]|0.0|0.0|Iris-setosa|Iris-setosa|

```

[5.1,3.8,1.6,0.2]	0.0	0.0	Iris-setosa	Iris-setosa	
[5.2,3.4,1.4,0.2]	0.0	0.0	Iris-setosa	Iris-setosa	
[5.2,3.5,1.5,0.2]	0.0	0.0	Iris-setosa	Iris-setosa	

+-----+-----+-----+-----+
only showing top 20 rows

Test Accuracy: 0.98

✓ Build a Clustering Model with Spark with a dataset of your choice

```
1 # Install PySpark
2 !pip install -q pyspark
3
4 # Import libraries
5 from pyspark.sql import SparkSession
6 from pyspark.ml.clustering import KMeans
7 from pyspark.ml.feature import VectorAssembler
8 from pyspark.ml.evaluation import ClusteringEvaluator
9 import urllib.request
10
11 # Download the Iris dataset locally
12 url = "https://archive.ics.uci.edu/ml/machine-learning-databases/iris/iris.data"
13 file_path = "/content/iris.csv"
14 urllib.request.urlretrieve(url, file_path)
15
16 # Start Spark session
17 spark = SparkSession.builder.appName("IrisKMeans").getOrCreate()
18
19 # Load the local CSV into Spark
20 df = spark.read.csv(file_path, inferSchema=True)
21 columns = ["sepal_length", "sepal_width", "petal_length", "petal_width", "class"]
22 df = df.toDF(*columns)
23
24 # Drop label column for unsupervised clustering
25 df_unlabeled = df.drop("class")
26
27 # Assemble features
28 assembler = VectorAssembler(inputCols=df_unlabeled.columns, outputCol="features")
29 data = assembler.transform(df_unlabeled)
30
31 # Apply KMeans clustering
32 kmeans = KMeans(featuresCol="features", k=3, seed=42)
33 model = kmeans.fit(data)
34 predictions = model.transform(data)
35
36 # Evaluate clustering
37 evaluator = ClusteringEvaluator()
38 silhouette = evaluator.evaluate(predictions)
39 print(f"Silhouette Score = {silhouette:.4f}")
40
41 # Show cluster centers
42 print("\nCluster Centers:")
43 for i, center in enumerate(model.clusterCenters()):
44     print(f"Cluster {i+1}: {center}")
```

🔗 Silhouette Score = 0.7342

Cluster Centers:
Cluster 1: [5.006 3.418 1.464 0.244]
Cluster 2: [6.85384615 3.07692308 5.71538462 2.05384615]
Cluster 3: [5.88360656 2.74098361 4.38852459 1.43442623]

✓ Build a Recommendation Engine with Spark with a custom dataset

```
1 !pip install pyspark ipywidgets --quiet
2
3 import pandas as pd
4 from pyspark.sql import SparkSession
5 from pyspark.ml.recommendation import ALS
6 from pyspark.sql.functions import col
7 import ipywidgets as widgets
8 from IPython.display import display
9
10 # Initialize Spark
11 spark = SparkSession.builder.master("local[*]").appName("MovieRecommender").getOrCreate()
12
13 # Sample MovieLens-like data
14 data = [
15     (0, 100, 4.0),
16     (0, 101, 3.5),
```

```
17 (1, 100, 5.0),
18 (1, 102, 4.0),
19 (2, 101, 2.5),
20 (2, 103, 4.5),
21 (3, 100, 4.0),
22 (3, 102, 5.0)
23 ]
24 ratings_df = spark.createDataFrame(data, ["userId", "movieId", "rating"])
25
26 # Train ALS model
27 als = ALS(userCol="userId", itemCol="movieId", ratingCol="rating", coldStartStrategy="drop", nonnegative=True)
28 model = als.fit(ratings_df)
29
30 # Widget setup
31 user_input = widgets.BoundedIntText(value=0, min=0, max=10, description='User ID:')
32 output = widgets.Output()
33
34 def on_click_submit(b):
35     output.clear_output()
36     with output:
37         try:
38             user_id = user_input.value
39             user_df = spark.createDataFrame([(user_id,)], ["userId"])
40             recommendations = model.recommendForUserSubset(user_df, 3).collect()
41             if recommendations:
42                 recs = recommendations[0]['recommendations']
43                 for rec in recs:
44                     print(f"Movie ID: {rec['movieId']} | Predicted Rating: {rec['rating']:.2f}")
45             else:
46                 print("No recommendations found.")
47         except Exception as e:
48             print("Error:", e)
49
50 submit_button = widgets.Button(description="Get Recommendations")
51 submit_button.on_click(on_click_submit)
52
53 # Display everything
54 display(widgets.VBox([user_input, submit_button, output]))
55
```



1.6/1.6 MB 31.6 MB/s eta 0:00:00

User ID:

Get Recommendation...

Movie ID: 100 | Predicted Rating: 4.82

Movie ID: 103 | Predicted Rating: 4.27

Movie ID: 102 | Predicted Rating: 4.06