

```

In [1]: import pyspark
import math
import numpy as np
from pyspark.sql import SparkSession
from pyspark.sql.types import StructType, StringType, LongType, IntegerType
from pyspark.ml.feature import VectorAssembler, StandardScaler
from pyspark.ml.feature import PCA
import matplotlib.pyplot as plt

# Feature columns
feature_columns = df_raw.columns
col_list = list(df_raw.columns)
feature_columns = list(set(col_list) - set(['c0', 'Polygon_ID', 'acq_date']))

# Assemble data
assembler = VectorAssembler(inputCols = feature_columns, outputCol = "features")
assembled_data = assembler.transform(df_raw)
# scaler = StandardScaler(inputCols="features", outputCol="scaled_features", withStd=True)
# scaler_model = scaler.fit(assembled_data)
scaled_data = scaler_model.transform(assembled_data)

# Let us perform a hyperparameter sweep
# using 'features'
def variance_search():
    k_values = range(13, 20) # Example: Try K from 2 to 10
    for k_val in k_values:
        # Fit the K-Means model to the data
        model = PCA(k=k_val, inputCol = "scaled_features", outputCol = "pcaFeatures")
        result = model.transform(scaled_data)
        explained_variance = model.explainedVariance
        print("Explained Variance: ", sum(explained_variance), "K: Value: ", k_val)
    # Return the best model
    return model

# Evaluate the clustering using the ClusteringEvaluator
silhouette_score = evaluator.evaluate(result)

# Print the WCSS (Within-Cluster Sum of Squares)
print("Silhouette Score: (silhouette_score)")
wcss = model.summary.trainingCost
print("Within cluster sum of squares (wcss):")

# Get the cluster sizes and centers
cluster_sizes = clustered_data.groupBy("prediction").count()
cluster_sizes.show()

# Return the clustered data
return clustered_data, silhouette_score, wcss

# Define a function to perform visualization
def visualize_clusters_2D(clustered_data, feature_1, feature_2, num_cluster_centers):
    # Requires:
    # clustered_data: dataframe returned from Kmeans fit
    # num_cluster_centers: integer number of clusters
    # feature_1: string with identifiers for first column name
    # feature_2: string with identifiers for second column name
    # Convert the DataFrame to a Pandas DataFrame for visualization
    pandas_df = clustered_data.select(feature_1, feature_2, "prediction").toPandas()
    # Extract the cluster assignments
    cluster_assignments = pandas_df["prediction"]
    # Extract the indices for the cluster center
    pl = col_list.index(feature_1)
    pr = col_list.index(feature_2)
    # Create a scatter plot for each cluster
    for cluster_id in range(num_cluster_centers):
        # Generate random color
        random_color = generate_random_color()
        cluster_data = pandas_df[pandas_df["prediction"] == cluster_id]
        plt.scatter(cluster_data[feature_1], cluster_data[feature_2], color=random_color,
                    label=cluster_id)
    # Set labels and title
    plt.xlabel(feature_1)
    plt.ylabel(feature_2)
    plt.title("K-Means Clustering")
    # Save figure to folder location
    mypath = "/storage/home/sxs6549/work/Project/graphs/"
    if not os.path.isdir(mypath):
        os.makedirs(mypath)
    # fig_name = "fig_cluster_2d_1" + feature_1
    path = "/storage/home/sxs6549/work/Project/graphs/figure_fig_1" + feature_1
    plt.savefig(path)
    # Show the legend
    plt.legend()
    # Display the plot
    plt.show()

# Generate random color
def generate_random_color():
    r = random.random()
    g = random.random()
    b = random.random()
    return (r, g, b)

# Define a function to perform visualization
def visualize_clusters_2D(clustered_data, feature_1, feature_2, num_cluster_centers):
    # Requires:
    # clustered_data: dataframe returned from Kmeans fit
    # num_cluster_centers: integer number of clusters
    # feature_1: string with identifiers for first column name
    # feature_2: string with identifiers for second column name
    # Convert the DataFrame to a Pandas DataFrame for visualization
    pandas_df = clustered_data.select(feature_1, feature_2, "prediction").toPandas()
    # Extract the cluster assignments
    cluster_assignments = pandas_df["prediction"]
    # Extract the indices for the cluster center
    pl = col_list.index(feature_1)
    pr = col_list.index(feature_2)
    # Create a scatter plot for each cluster
    for cluster_id in range(num_cluster_centers):
        # Generate random color
        random_color = generate_random_color()
        cluster_data = pandas_df[pandas_df["prediction"] == cluster_id]
        plt.scatter(cluster_data[feature_1], cluster_data[feature_2], color=random_color,
                    label=cluster_id)
    # Set labels and title
    plt.xlabel(feature_1)
    plt.ylabel(feature_2)
    plt.title("K-Means Clustering")
    # Save figure to folder location
    mypath = "/storage/home/sxs6549/work/Project/graphs/"
    if not os.path.isdir(mypath):
        os.makedirs(mypath)
    # fig_name = "fig_cluster_2d_1" + feature_1
    path = "/storage/home/sxs6549/work/Project/graphs/figure_fig_1" + feature_1
    plt.savefig(path)
    # Show the legend
    plt.legend()
    # Display the plot
    plt.show()

# Generate random color
def generate_random_color():
    r = random.random()
    g = random.random()
    b = random.random()
    return (r, g, b)

# Define a function to perform visualization
def visualize_clusters_2D(clustered_data, feature_1, feature_2, num_cluster_centers):
    # Requires:
    # clustered_data: dataframe returned from Kmeans fit
    # num_cluster_centers: integer number of clusters
    # feature_1: string with identifiers for first column name
    # feature_2: string with identifiers for second column name
    # Convert the DataFrame to a Pandas DataFrame for visualization
    pandas_df = clustered_data.select(feature_1, feature_2, "prediction").toPandas()
    # Extract the cluster assignments
    cluster_assignments = pandas_df["prediction"]
    # Extract the indices for the cluster center
    pl = col_list.index(feature_1)
    pr = col_list.index(feature_2)
    # Create a scatter plot for each cluster
    for cluster_id in range(num_cluster_centers):
        # Generate random color
        random_color = generate_random_color()
        cluster_data = pandas_df[pandas_df["prediction"] == cluster_id]
        plt.scatter(cluster_data[feature_1], cluster_data[feature_2], color=random_color,
                    label=cluster_id)
    # Set labels and title
    plt.xlabel(feature_1)
    plt.ylabel(feature_2)
    plt.title("K-Means Clustering")
    # Save figure to folder location
    mypath = "/storage/home/sxs6549/work/Project/graphs/"
    if not os.path.isdir(mypath):
        os.makedirs(mypath)
    # fig_name = "fig_cluster_2d_1" + feature_1
    path = "/storage/home/sxs6549/work/Project/graphs/figure_fig_1" + feature_1
    plt.savefig(path)
    # Show the legend
    plt.legend()
    # Display the plot
    plt.show()

# Generate random color
def generate_random_color():
    r = random.random()
    g = random.random()
    b = random.random()
    return (r, g, b)

# Define a function to perform visualization
def visualize_clusters_2D(clustered_data, feature_1, feature_2, num_cluster_centers):
    # Requires:
    # clustered_data: dataframe returned from Kmeans fit
    # num_cluster_centers: integer number of clusters
    # feature_1: string with identifiers for first column name
    # feature_2: string with identifiers for second column name
    # Convert the DataFrame to a Pandas DataFrame for visualization
    pandas_df = clustered_data.select(feature_1, feature_2, "prediction").toPandas()
    # Extract the cluster assignments
    cluster_assignments = pandas_df["prediction"]
    # Extract the indices for the cluster center
    pl = col_list.index(feature_1)
    pr = col_list.index(feature_2)
    # Create a scatter plot for each cluster
    for cluster_id in range(num_cluster_centers):
        # Generate random color
        random_color = generate_random_color()
        cluster_data = pandas_df[pandas_df["prediction"] == cluster_id]
        plt.scatter(cluster_data[feature_1], cluster_data[feature_2], color=random_color,
                    label=cluster_id)
    # Set labels and title
    plt.xlabel(feature_1)
    plt.ylabel(feature_2)
    plt.title("K-Means Clustering")
    # Save figure to folder location
    mypath = "/storage/home/sxs6549/work/Project/graphs/"
    if not os.path.isdir(mypath):
        os.makedirs(mypath)
    # fig_name = "fig_cluster_2d_1" + feature_1
    path = "/storage/home/sxs6549/work/Project/graphs/figure_fig_1" + feature_1
    plt.savefig(path)
    # Show the legend
    plt.legend()
    # Display the plot
    plt.show()

# Generate random color
def generate_random_color():
    r = random.random()
    g = random.random()
    b = random.random()
    return (r, g, b)

# Define a function to perform visualization
def visualize_clusters_2D(clustered_data, feature_1, feature_2, num_cluster_centers):
    # Requires:
    # clustered_data: dataframe returned from Kmeans fit
    # num_cluster_centers: integer number of clusters
    # feature_1: string with identifiers for first column name
    # feature_2: string with identifiers for second column name
    # Convert the DataFrame to a Pandas DataFrame for visualization
    pandas_df = clustered_data.select(feature_1, feature_2, "prediction").toPandas()
    # Extract the cluster assignments
    cluster_assignments = pandas_df["prediction"]
    # Extract the indices for the cluster center
    pl = col_list.index(feature_1)
    pr = col_list.index(feature_2)
    # Create a scatter plot for each cluster
    for cluster_id in range(num_cluster_centers):
        # Generate random color
        random_color = generate_random_color()
        cluster_data = pandas_df[pandas_df["prediction"] == cluster_id]
        plt.scatter(cluster_data[feature_1], cluster_data[feature_2], color=random_color,
                    label=cluster_id)
    # Set labels and title
    plt.xlabel(feature_1)
    plt.ylabel(feature_2)
    plt.title("K-Means Clustering")
    # Save figure to folder location
    mypath = "/storage/home/sxs6549/work/Project/graphs/"
    if not os.path.isdir(mypath):
        os.makedirs(mypath)
    # fig_name = "fig_cluster_2d_1" + feature_1
    path = "/storage/home/sxs6549/work/Project/graphs/figure_fig_1" + feature_1
    plt.savefig(path)
    # Show the legend
    plt.legend()
    # Display the plot
    plt.show()

# Generate random color
def generate_random_color():
    r = random.random()
    g = random.random()
    b = random.random()
    return (r, g, b)

# Define a function to perform visualization
def visualize_clusters_2D(clustered_data, feature_1, feature_2, num_cluster_centers):
    # Requires:
    # clustered_data: dataframe returned from Kmeans fit
    # num_cluster_centers: integer number of clusters
    # feature_1: string with identifiers for first column name
    # feature_2: string with identifiers for second column name
    # Convert the DataFrame to a Pandas DataFrame for visualization
    pandas_df = clustered_data.select(feature_1, feature_2, "prediction").toPandas()
    # Extract the cluster assignments
    cluster_assignments = pandas_df["prediction"]
    # Extract the indices for the cluster center
    pl = col_list.index(feature_1)
    pr = col_list.index(feature_2)
    # Create a scatter plot for each cluster
    for cluster_id in range(num_cluster_centers):
        # Generate random color
        random_color = generate_random_color()
        cluster_data = pandas_df[pandas_df["prediction"] == cluster_id]
        plt.scatter(cluster_data[feature_1], cluster_data[feature_2], color=random_color,
                    label=cluster_id)
    # Set labels and title
    plt.xlabel(feature_1)
    plt.ylabel(feature_2)
    plt.title("K-Means Clustering")
    # Save figure to folder location
    mypath = "/storage/home/sxs6549/work/Project/graphs/"
    if not os.path.isdir(mypath):
        os.makedirs(mypath)
    # fig_name = "fig_cluster_2d_1" + feature_1
    path = "/storage/home/sxs6549/work/Project/graphs/figure_fig_1" + feature_1
    plt.savefig(path)
    # Show the legend
    plt.legend()
    # Display the plot
    plt.show()

# Generate random color
def generate_random_color():
    r = random.random()
    g = random.random()
    b = random.random()
    return (r, g, b)

# Define a function to perform visualization
def visualize_clusters_2D(clustered_data, feature_1, feature_2, num_cluster_centers):
    # Requires:
    # clustered_data: dataframe returned from Kmeans fit
    # num_cluster_centers: integer number of clusters
    # feature_1: string with identifiers for first column name
    # feature_2: string with identifiers for second column name
    # Convert the DataFrame to a Pandas DataFrame for visualization
    pandas_df = clustered_data.select(feature_1, feature_2, "prediction").toPandas()
    # Extract the cluster assignments
    cluster_assignments = pandas_df["prediction"]
    # Extract the indices for the cluster center
    pl = col_list.index(feature_1)
    pr = col_list.index(feature_2)
    # Create a scatter plot for each cluster
    for cluster_id in range(num_cluster_centers):
        # Generate random color
        random_color = generate_random_color()
        cluster_data = pandas_df[pandas_df["prediction"] == cluster_id]
        plt.scatter(cluster_data[feature_1], cluster_data[feature_2], color=random_color,
                    label=cluster_id)
    # Set labels and title
    plt.xlabel(feature_1)
    plt.ylabel(feature_2)
    plt.title("K-Means Clustering")
    # Save figure to folder location
    mypath = "/storage/home/sxs6549/work/Project/graphs/"
    if not os.path.isdir(mypath):
        os.makedirs(mypath)
    # fig_name = "fig_cluster_2d_1" + feature_1
    path = "/storage/home/sxs6549/work/Project/graphs/figure_fig_1" + feature_1
    plt.savefig(path)
    # Show the legend
    plt.legend()
    # Display the plot
    plt.show()

# Generate random color
def generate_random_color():
    r = random.random()
    g = random.random()
    b = random.random()
    return (r, g, b)

# Define a function to perform visualization
def visualize_clusters_2D(clustered_data, feature_1, feature_2, num_cluster_centers):
    # Requires:
    # clustered_data: dataframe returned from Kmeans fit
    # num_cluster_centers: integer number of clusters
    # feature_1: string with identifiers for first column name
    # feature_2: string with identifiers for second column name
    # Convert the DataFrame to a Pandas DataFrame for visualization
    pandas_df = clustered_data.select(feature_1, feature_2, "prediction").toPandas()
    # Extract the cluster assignments
    cluster_assignments = pandas_df["prediction"]
    # Extract the indices for the cluster center
    pl = col_list.index(feature_1)
    pr = col_list.index(feature_2)
    # Create a scatter plot for each cluster
    for cluster_id in range(num_cluster_centers):
        # Generate random color
        random_color = generate_random_color()
        cluster_data = pandas_df[pandas_df["prediction"] == cluster_id]
        plt.scatter(cluster_data[feature_1], cluster_data[feature_2], color=random_color,
                    label=cluster_id)
    # Set labels and title
    plt.xlabel(feature_1)
    plt.ylabel(feature_2)
    plt.title("K-Means Clustering")
    # Save figure to folder location
    mypath = "/storage/home/sxs6549/work/Project/graphs/"
    if not os.path.isdir(mypath):
        os.makedirs(mypath)
    # fig_name = "fig_cluster_2d_1" + feature_1
    path = "/storage/home/sxs6549/work/Project/graphs/figure_fig_1" + feature_1
    plt.savefig(path)
    # Show the legend
    plt.legend()
    # Display the plot
    plt.show()

# Generate random color
def generate_random_color():
    r = random.random()
    g = random.random()
    b = random.random()
    return (r, g, b)

# Define a function to perform visualization
def visualize_clusters_2D(clustered_data
```



```
from IPython.core.display import HTML
display(HTML("estyle{pre { white-space: pre !important; }<style>"))
```

```
ss = SparkSession.builder.config("spark.driver.memory", "16g").appName("ProjectTree1")
ss = SparkSession.builder.config("spark.driver.memory", "5g").master("local").appName("PCAExample1").getOrCreate()
```

```
ss.sparkContext.setCheckpointDir("/storage/home/sxs6549/work/Project/scratch")
```

```
df_raw = ss.read.csv("wildfiredb.csv", header=True, inferSchema=True)
df_raw = spark.read.csv("wildfire100.csv", header=True, inferSchema=True)
column_names = df_raw.columns

df_raw = df_raw.drop("acq_date")
df_raw = df_raw.dropna()
```

CPU times: user 43.9 ms, sys: 12.3 ms, total: 56.2 ms
Wall time: 2min 40s

```
df_raw_trial = ss.read.csv("fire_small.csv", header=True, inferSchema=True)
df_raw = spark.read.csv("wildfire100.csv", header=True, inferSchema=True)
column_names = df_raw_trial.columns

df_raw = df_raw.drop("acq_date")
df_raw_trial = df_raw_trial.dropna()
```

CPU times: user 25.6 ms, sys: 3.78 ms, total: 29.4 ms
Wall time: 11.5 s

```
col_list = list(df_raw.columns)
col_list
```

```
col_list_new = list(set(col_list) - set(['_c0', 'Polygon_ID', 'acq_date', 'frp']))
col_list_new
```

```
feature_columns = df_raw.columns
col_list = list(df_raw.columns)
feature_inputs = list(set(col_list) - set(['_c0', 'Polygon_ID', 'acq_date', 'frp']))

assembler_tree = VectorAssembler(inputCols=feature_inputs, outputCol="features")
assembled_data_tree = assembler_tree.transform(df_raw)
scaler = StandardScaler(inputCols="features", outputCol="scaled_features", withStd=True)
scaler_model = scaler.fit(assembled_data)

scaled_data = scaler_model.transform(assembled_data)
```

```
pca_tree = PCA(k=36, inputCol="features", outputCol="pcaFeatures")
model_tree = pca_tree.fit(assembled_data_tree)
result_tree = model_tree.transform(assembled_data_tree)
```

```
from pyspark.ml import Pipeline
from pyspark.ml.regression import DecisionTreeRegressor
from pyspark.ml.feature import VectorAssembler, StandardScaler
from pyspark.ml.evaluation import RegressionEvaluator
from pyspark.mllib.util import MLUtils

from decision_tree_plot.decision_tree_parser import decision_tree_parse
from decision_tree_plot.decision_tree_plot import plot_trees
# Split the data into training and test sets (20% held out for testing)
(trainingData, testingData) = result_tree.randomSplit([0.8, 0.2], seed=1237)

#Code cell for Part 7
## Initialize a Pandas DataFrame to store evaluation results of all combination of hyperparameters
hyperparams_eval_df = pd.DataFrame(columns = ['max_depth', 'minInstancesPerNode', 'feature_importance'])
index = 0
# Initialize lowest_error
lowest_testing_rmse = 100000
# Set up the possible hyperparameter values to be evaluated
max_depth_list = [2, 3, 4, 5, 6, 7, 8, 9, 10]
minInstancesPerNode_list = [2, 3, 4, 5, 6, 7]
max_depth_list = [2]
minInstancesPerNode_list = [9]
#labelIndexer = StringIndexer(inputCol="class", outputCol="indexedLabel").fit(data2)
#feature_inputs = list(set(col_list) - set(['_c0', 'Polygon_ID', 'acq_date', 'frp']))
#assembler = VectorAssembler(inputCols=feature_inputs, outputCol="features")
#labelConverter = IndexToString(inputCol="prediction", outputCol="predictedClass", model_path="/storage/home/sxs6549/work/Project/fire_DTmodel_v1a")

dataFrame['_c0', 'Polygon_ID', 'acq_date', 'frp', 'double', 'Neighbour', 'int', 'Neighl
```

```
##time
for max_depth in max_depth_list:
    for minInsPN in minInstancesPerNode_list:
        trainingData.persist()
        testingData.persist()

        seed = 37
        # Construct a DT model using a set of hyper-parameter values and training data
        dt = DecisionTreeClassifier(labelCol="indexedLabel", featuresCol="features", maxDepth=max_depth, minInstancesPerNode=minInsPN, seed=seed)
        # Pipeline stages=[labelIndexer, assembler, dt, predictionConverter]
        pipeline = Pipeline(stages=[labelIndexer, assembler, dt, predictionConverter])
        model = dt.fit(trainingData)
        training_predictions = model.transform(trainingData)
        testing_predictions = model.transform(testingData)
        #evaluator = MulticlassClassificationEvaluator(labelCol="indexedLabel", predictionCol="prediction", metricName="areaUnderROC")
        evaluator = RegressionEvaluator(labelCol="frp", predictionCol="prediction", metricName="rmse")
        testing_rmse = evaluator.evaluate(testing_predictions)
        # We use 0 as default value of the 'best Model' column in the Pandas DataFrame
        # The best model will have a value 1000
        hyperparams_eval_df.loc[index] = [max_depth, minInsPN, training_rmse, testing_rmse]
        index = index + 1
        if testing_rmse < lowest_testing_rmse:
            best_max_depth = max_depth
            best_minInsPN = minInsPN
            best_index = index - 1
            best_parameters = training_rmse
            best_DTmodel = model
            best_tree = decision_tree_parse(best_DTmodel, ss, model_path)
            column = dict([(str(idx), i) for idx, i in enumerate(feature_inputs)])
            column = dict([(str(idx), i) for idx, i in enumerate(feature_inputs)])
            print('The best max depth is ', best_max_depth, ', best minInstancesPerNode = ', best_minInsPN, ', testing_rmse = ', lowest_testing_rmse)
            column = dict([(str(idx), i) for idx, i in enumerate(feature_inputs)])

The best max depth is 8, best minInstancesPerNode = 4, testing_rmse = 51.34257362
CPU times: user 1.41 s, sys: 313 ms, total: 1.72 s
Wall time: 4min 12s
```

```
#training_predictions.show(100)
```

```
#Code cell for Part 7
best_model_path_part7="/storage/home/sxs6549/work/Project/fire_DT_HPT_cluster"
```

```
#Code cell for Part 7
best_tree=decision_tree_parse(best_DTmodel, ss, best_model_path_part7)
column = dict([(str(idx), i) for idx, i in enumerate(feature_inputs)])
plot_trees(best_tree, column = column, output_path = "/storage/home/sxs6549/work/Project/fire_DT_HPT_cluster_v1a")
```

```
#Code cell for Part 7
# Store the Testing RMS in the DataFrame
hyperparams_eval_df.loc[best_index]=best_max_depth, best_minInsPN, best_parameters_testing_rmse
output_path = "/storage/home/sxs6549/work/Project/fire_HPT_cluster.csv"
hyperparams_eval_df.to_csv(output_path)
```