


```

from pyspark.ml.feature import VectorAssembler
from pyspark.ml.regression import LinearRegression

# Define features and label
assembler = VectorAssembler(inputCols=['total_cases', 'new_cases', 'population'], outputCol="features")
data = assembler.transform(kenya_df).select('features', 'total_deaths')

# Split data
train_data, test_data = data.randomSplit([0.8, 0.2], seed=42)

# Define and train the model
lr = LinearRegression(featuresCol='features', labelCol='total_deaths')
lr_model = lr.fit(train_data)

```

This section of the script focuses on selecting features, assembling them, splitting the dataset, and training the model on the training set.

Step 5: Visualize the Model

The model's predictions were visualized using Matplotlib:

```

import matplotlib.pyplot as plt

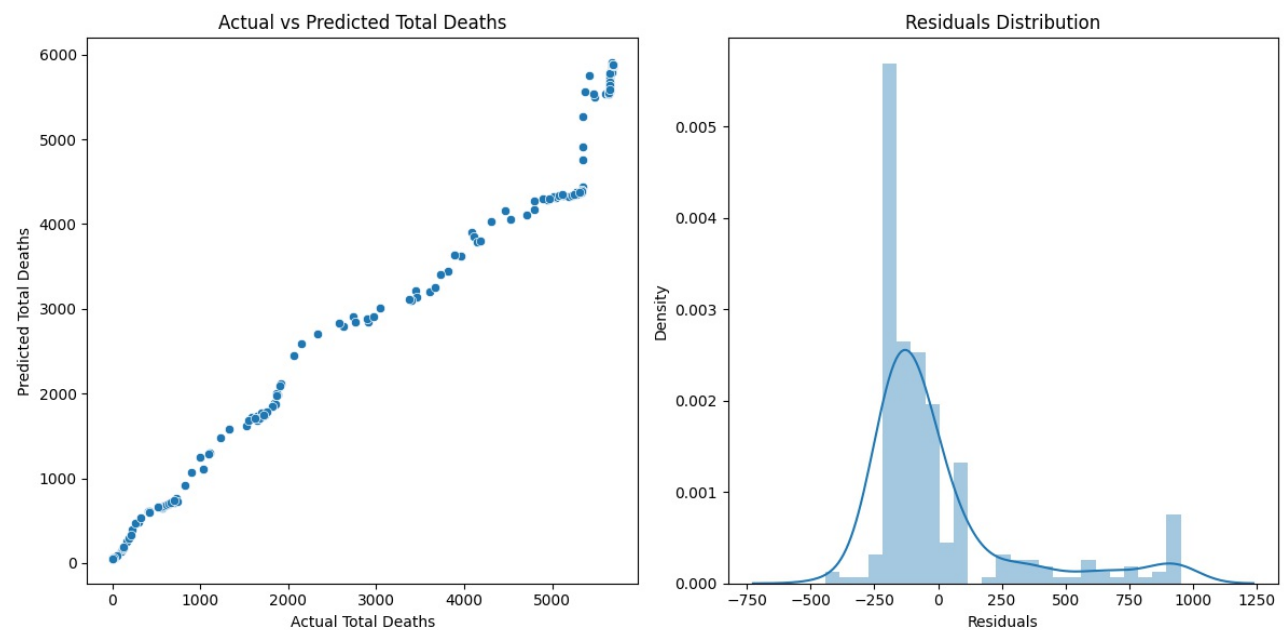
# Make predictions
predictions = lr_model.transform(test_data)

# Convert to Pandas DataFrame for visualization
pandas_df = predictions.select("total_deaths", "prediction").toPandas()

# Visualization
plt.scatter(pandas_df['total_deaths'], pandas_df['prediction'])
plt.xlabel('Actual Total Deaths')
plt.ylabel('Predicted Total Deaths')
plt.title('Actual vs Predicted Total Deaths')
plt.show()

```

The scatter plot illustrates the relationship between actual and predicted total deaths, providing a visual assessment of the model's performance.



Step 6: Test the Model

The model's performance was evaluated using the Root Mean Squared Error (RMSE) metric:

```

from pyspark.ml.evaluation import RegressionEvaluator

# Evaluate the model
evaluator = RegressionEvaluator(labelCol="total_deaths", predictionCol="prediction", metricName="rmse")
rmse = evaluator.evaluate(predictions)
print(f"Root Mean Squared Error (RMSE) on test data: {rmse}")

```

The RMSE value offers a quantitative measure of the model's accuracy, indicating how closely the predicted values match the actual death counts.

Actual vs Predicted Total Deaths

