

```

1 from google.colab import drive
2 drive.mount('/content/drive')

→ Drive already mounted at /content/drive; to attempt to forcibly remount, call drive.mount("/content/drive", force_remount=True).

1 !tar xf /content/drive/MyDrive/Big_Data/spark-3.1.1-bin-hadoop3.2.tgz
2 !pip install -q findspark

→ tar: /content/drive/MyDrive/Big_Data/spark-3.1.1-bin-hadoop3.2.tgz: Cannot open: No such file or directory
tar: Error is not recoverable: exiting now

1 !apt-get install openjdk-8-jdk-headless -qq > /dev/null

1 import os
2 os.environ["JAVA_HOME"] = "/usr/lib/jvm/java-8-openjdk-amd64"
3 os.environ["SPARK_HOME"] = "/content/spark-3.1.1-bin-hadoop3.2"

1 import findspark
2 findspark.init()

1 from pyspark.ml.classification import MultilayerPerceptronClassifier, RandomForestClassifier, LinearSVC, OneVsRest
2 from pyspark.sql.types import LongType, StringType, DoubleType, IntegerType, ArrayType
3 from pyspark.ml.evaluation import MulticlassClassificationEvaluator, RegressionEvaluator
4 from pyspark.ml.feature import VectorAssembler, StandardScaler
5 from pyspark.ml.regression import LinearRegression
6 from pyspark.ml.linalg import Vectors, VectorUDT
7 from pyspark.sql import SparkSession, Window
8 from pyspark import SparkContext, SparkConf
9 from pyspark.ml.recommendation import ALS
10 from pyspark.sql import SparkSession, Row
11 from pyspark.ml.clustering import KMeans
12 from pyspark.sql.window import Window
13 from pyspark.ml.fpm import FPGrowth
14 from itertools import combinations
15 import pyspark.sql.functions as F
16 import matplotlib.pyplot as plt
17 from functools import reduce
18 import pandas as pd
19 import numpy as np
20 import findspark
21 import os

1 path = "/content/drive/MyDrive/Big_Data/datasets/"

```

## ✓ Câu 1: Phân cụm dữ liệu

```

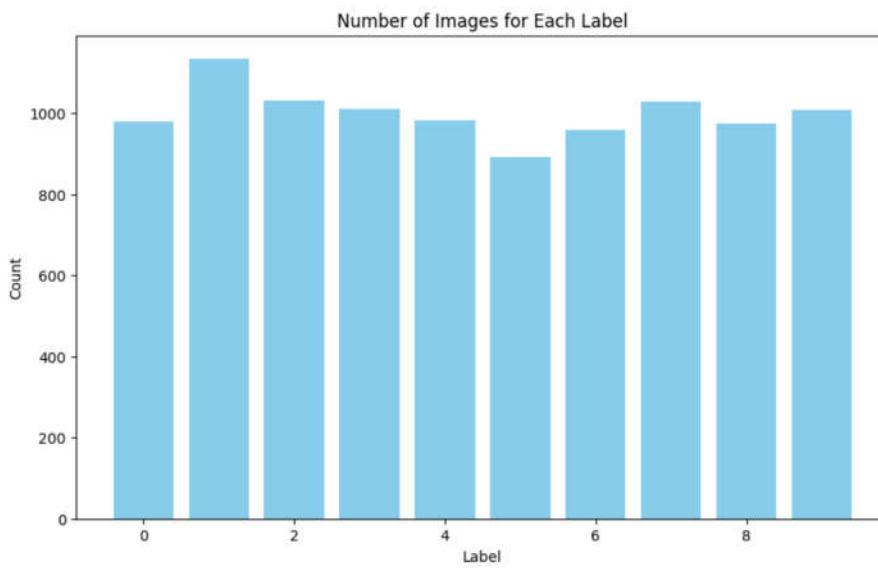
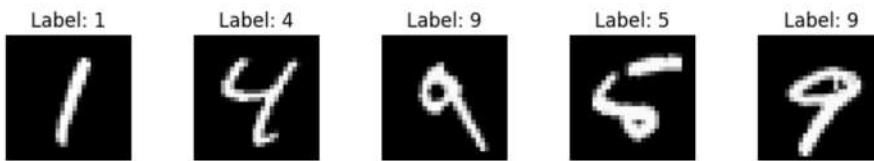
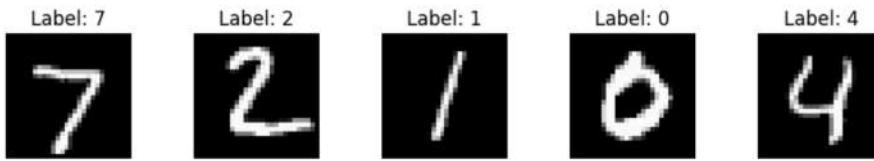
1 # Khởi tạo SparkSession
2 spark = SparkSession.builder \
3     .appName("Final") \
4     .getOrCreate()
5
6 # Đọc dữ liệu từ file CSV vào DataFrame
7 data_path_mnist_mini= "/content/drive/MyDrive/Big_Data/datasets/mnist_mini.csv"
8 df = spark.read.csv(data_path_mnist_mini, header=False, inferSchema=True)
9
10 # Xem một số dòng đầu tiên của DataFrame
11 df.show(3)
12 print("Number of rows:", df.count())
13
14 # Chọn một số dòng dữ liệu để hiển thị
15 sample_data = df.take(10)
16
17 # Hiển thị một số ảnh từ dữ liệu mẫu
18 fig, axs = plt.subplots(2, 5, figsize=(10, 4))
19 fig.subplots_adjust(hspace=0.5, wspace=0.5)
20 for i in range(10):
21     img = np.array(sample_data[i][1:]).reshape(28, 28)
22     axs[i//5, i%5].imshow(img, cmap='gray')
23     axs[i//5, i%5].set_title(f"Label: {sample_data[i][0]}")
24     axs[i//5, i%5].axis('off')
25 plt.show()
26
27 # Thống kê số lượng các nhãn
28 label_counts = df.groupBy("_c0").count().orderBy("_c0").collect()
29
30 # Trích xuất nhãn và số lượng từ kết quả
31 labels = [row["_c0"] for row in label_counts]
32 counts = [row["count"] for row in label_counts]
33
34 # Vẽ biểu đồ
35 plt.figure(figsize=(10, 6))
36 plt.bar(labels, counts, color='skyblue')
37 plt.xlabel('Label')
38 plt.ylabel('Count')
39 plt.title('Number of Images for Each Label')
40 plt.show()

```

```
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
| _c0|_c1|_c2|_c3|_c4|_c5|_c6|_c7|_c8|_c9|_c10|_c11|_c12|_c13|_c14|_c15|_c16|_c17|_c18|_c
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
| 7| 0| 0| 0| 0| 0| 0| 0| 0| 0| 0| 0| 0| 0| 0| 0| 0| 0| 0| 0| 0| 0| 0| 0| 0| 0| 0| 0| 0|
| 2| 0| 0| 0| 0| 0| 0| 0| 0| 0| 0| 0| 0| 0| 0| 0| 0| 0| 0| 0| 0| 0| 0| 0| 0| 0| 0| 0| 0|
| 1| 0| 0| 0| 0| 0| 0| 0| 0| 0| 0| 0| 0| 0| 0| 0| 0| 0| 0| 0| 0| 0| 0| 0| 0| 0| 0| 0| 0|
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
```

only showing top 3 rows

Number of rows: 10000

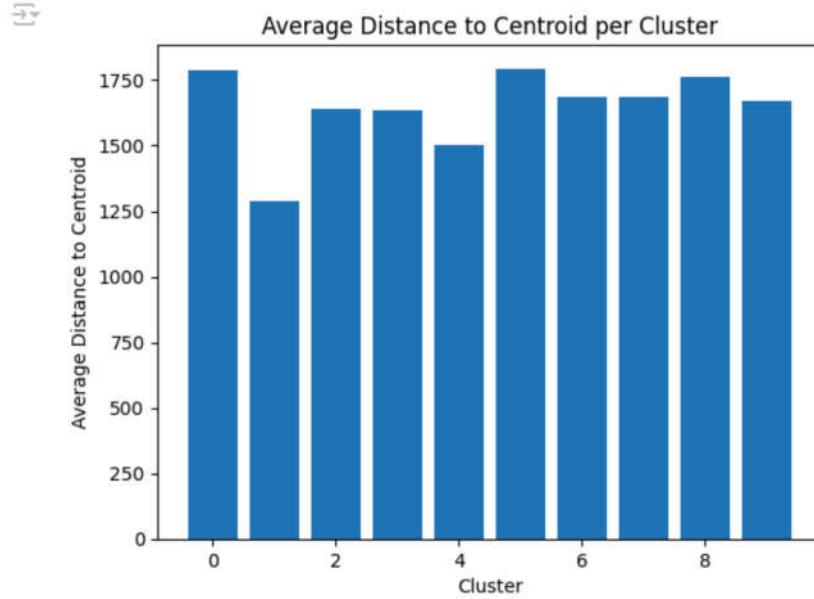


```
1 class MNISTClustering:
2     def __init__(self,spark, data_path, weighted_indices, k=10):
3         """
4             Initializes the MNISTClustering class.
5             Parameters:
6                 - data_path: str, path to the MNIST dataset.
7                 - weighted_indices: list of int, indices of rows to be weighted more heavily.
8                 - k: int, number of clusters for KMeans.
9         """
10        self.spark = spark
11        self.df = self.spark.read.csv(data_path, header=False, inferSchema=True)
12        self.weighted_indices = weighted_indices
13        self.k = k
14        self.predictions = None
15
16    def preprocess(self):
17        """
18            Preprocesses the data by assembling features and applying weights to specified indices.
19        """
20        assembler = VectorAssembler(inputCols=self.df.columns[1:], outputCol="features")
21        self.df = assembler.transform(self.df)
22        self.df = self.df.withColumn("weight", F.lit(1))
23        self.df = self.df.withColumn("index", F.monotonically_increasing_id())
24        windowSpec = Window.orderBy("index")
25        self.df = self.df.withColumn("index", F.row_number().over(windowSpec) - 1)
26        self.df = self.df.withColumn(
27            "weight",
28            F.when(F.col("index").isin(self.weighted_indices), F.col("weight") * 100).otherwise(F.col("weight")))
29
30
31    def cluster(self):
32        """
33            Performs KMeans clustering on the preprocessed data.
34        """
35        kmeans = KMeans(k=self.k, seed=42, featuresCol="features", weightCol="weight")
36        self.model = kmeans.fit(self.df)
37        self.predictions = self.model.transform(self.df)
38
39    def analyze(self):
40        """
41            Analyzes the clustering results by computing the average distance of points to their cluster centroids.
42        """
43
```

```

42     """
43     Returns:
44     - avg_distances: list of float, average distances to the centroids for each cluster.
45     """
46     centroids = self.model.clusterCenters()
47     def compute_distance(features, centroid):
48         return float(np.linalg.norm(features - centroid))
49     distances = self.predictions.rdd.map(
50         lambda row: (row['prediction'], (compute_distance(row['features'], centroids[row['prediction']]), 1)))
51     ).reduceByKey(lambda a, b: (a[0] + b[0], a[1] + b[1])).mapValues(
52         lambda x: x[0] / x[1]
53     ).collect()
54     avg_distances = [0] * self.k
55     for cluster, avg_distance in distances:
56         avg_distances[cluster] = avg_distance
57     return avg_distances
58
59 def visualize(self, avg_distances):
60     """
61     Visualizes the average distance to centroid per cluster using a bar chart.
62     Parameters:
63     - avg_distances: list of float, average distances to the centroids for each cluster.
64     """
65     plt.bar(range(self.k), avg_distances)
66     plt.xlabel('Cluster')
67     plt.ylabel('Average Distance to Centroid')
68     plt.title('Average Distance to Centroid per Cluster')
69     plt.show()
70
71 def run(self):
72     """
73     Runs the full clustering analysis: preprocessing, clustering, analysis, and visualization.
74     """
75     self.preprocess()
76     self.cluster()
77     avg_distances = self.analyze()
78     self.visualize(avg_distances)
79 if __name__ == "__main__":
80     data_path_mnist_mini = "/content/drive/MyDrive/Big_Data/datasets/mnist_mini.csv"
81     weighted_indices = [0, 1, 2, 3, 4, 7, 8, 11, 18, 61]
82     analysis = MNISTClustering(spark, data_path_mnist_mini, weighted_indices)
83     analysis.run()

```



## ▼ Câu 2: Giảm số chiều với SVD

```

1 from pyspark.ml.linalg import Vectors as MLVectors
2 from pyspark.mllib.linalg import Vectors as MLLibVectors
3 from pyspark.mllib.linalg.distributed import RowMatrix
4 import numpy as np
5 import matplotlib.pyplot as plt
6
7 class MNISTSVD:
8     def __init__(self, spark, predictions, k=3, n=100):
9         """
10         Initializes the MNISTSVD class.
11         Parameters:
12         - predictions: DataFrame, predictions from a clustering model.
13         - k: int, number of principal components to retain.
14         - n: int, number of data points to sample for visualization.
15         """
16         self.spark = spark
17         self.predictions = predictions
18         self.k = k
19         self.n = n
20
21     def reduce_dimensions(self):
22         """
23         Reduces the dimensions of the input data using SVD.
24         Returns:
25         - svd_result: DataFrame, transformed data with reduced dimensions.
26         """
27         assembled_data = self.predictions.select("features", "prediction")
28         prediction = self.predictions.select("prediction")

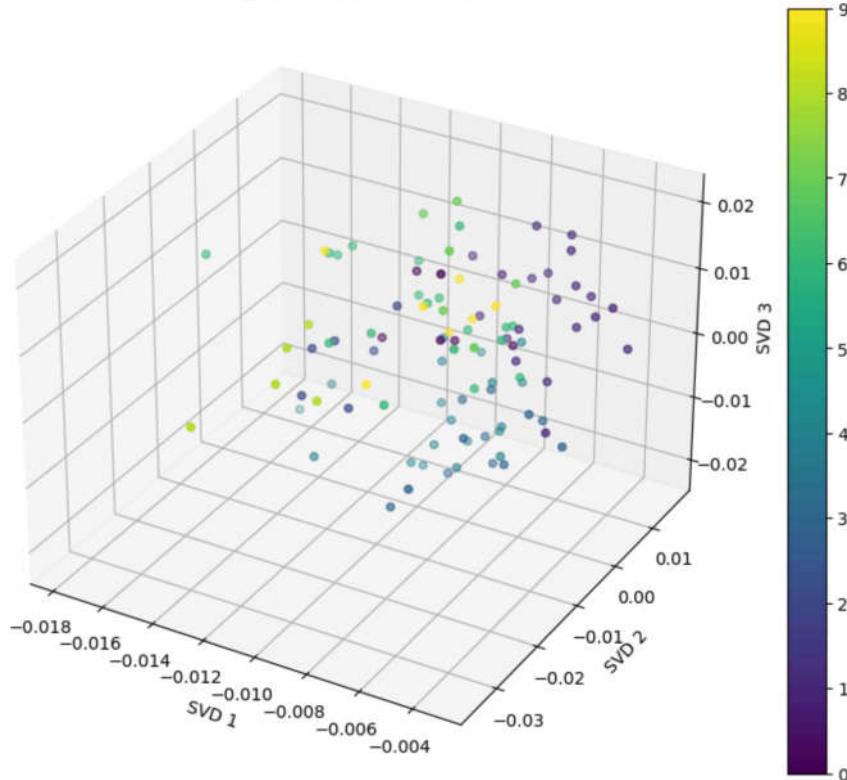
```

```

50     rows = assembled_data.rdd.map(lambda row: (MLLibVectors.dense(row["features"].toArray()), row["prediction"]))
51     mat = RowMatrix(rows.map(lambda row: row[0]))
52     svd = mat.computeSVD(self.k, computeU=True)
53     U = svd.U.rows.map(lambda row: row.tolist())
54     svd_df = spark.createDataFrame(U, ["col1", "col2", "col3"])
55     svd_df = svd_df.withColumn("index", F.monotonically_increasing_id())
56     prediction = prediction.withColumn("index", F.monotonically_increasing_id())
57     svd_result = svd_df.join(prediction, on="index").drop("index")
58
59     return svd_result
60
61 def sample_data(self, svd_result):
62     """
63         Samples data for visualization.
64         Parameters:
65             - svd_result: DataFrame, data with reduced dimensions.
66         Returns:
67             - svd_features_np: numpy array, sampled SVD features.
68             - sampled_predictions: numpy array, sampled predictions.
69     """
70     sampled_data = svd_result.sample(withReplacement=False, fraction=float(self.n) / svd_result.count(), seed=3).limit(100)
71     sampled_data_pd = sampled_data.toPandas()
72     svd_features_np = np.array(sampled_data_pd[["col1", "col2", "col3"]].values.tolist())
73     sampled_predictions = np.array(sampled_data_pd["prediction"].tolist())
74
75     return svd_features_np, sampled_predictions
76
77 def plot_3d(self, svd_features_np, sampled_predictions):
78     """
79         Plots the sampled data in 3D.
80         Parameters:
81             - svd_features_np: numpy array, sampled SVD features.
82             - sampled_predictions: numpy array, sampled predictions.
83     """
84     fig = plt.figure(figsize=(10, 8))
85     ax = fig.add_subplot(111, projection='3d')
86     sc = ax.scatter(svd_features_np[:, 0], svd_features_np[:, 1], svd_features_np[:, 2], c=sampled_predictions, cmap="viridis")
87     plt.colorbar(sc)
88     ax.set_xlabel("SVD 1")
89     ax.set_ylabel("SVD 2")
90     ax.set_zlabel("SVD 3")
91     plt.title("3D SVD of MNIST Data")
92     plt.show()
93
94 def run(self):
95     """
96         Runs the entire SVD analysis.
97     """
98     svd_result = self.reduce_dimensions()
99     svd_features_np, sampled_predictions = self.sample_data(svd_result)
100    self.plot_3d(svd_features_np, sampled_predictions)
101
102 if __name__ == "__main__":
103     predictions_df = analysis.predictions
104     svd_analysis = MNISTSVD(spark, predictions_df)
105     svd_analysis.run()
106
107

```

3D SVD of MNIST Data



## ✓ Câu 3: Khuyến nghị sản phẩm với Collaborative Filtering

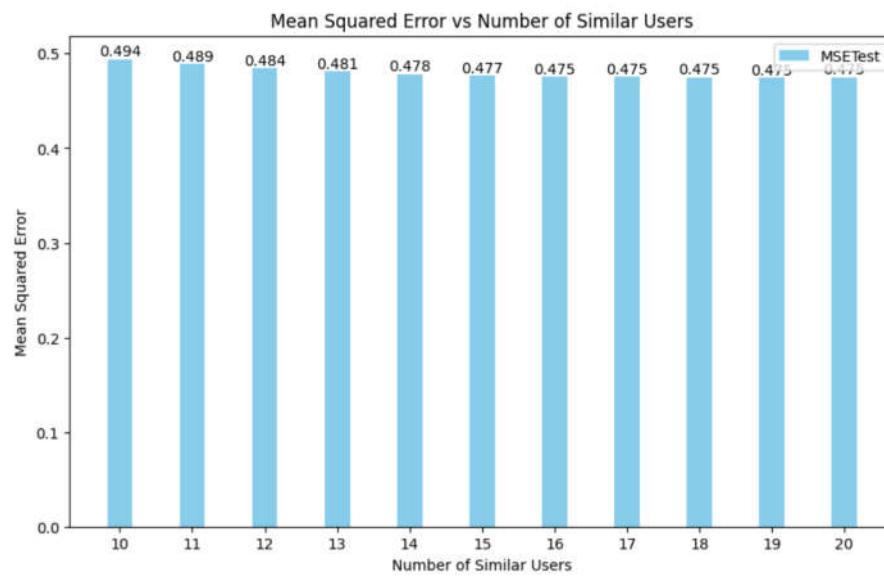
```
2     def __init__(self, filepath, rank_value=10, regParam_value=0.1, maxIter_value=10, list_num_users=list(range(10,21))):  
3         self.spark = SparkSession.builder.appName("CollaborativeFiltering").getOrCreate()  
4         self.data = self.spark.read.csv(filepath, header=True, inferSchema=True)  
5         (self.train, self.test) = self.data.randomSplit([0.7, 0.3], seed=96)  
6         self.rank_value = int(rank_value)  
7         self.regParam_value = float(regParam_value)  
8         self.maxIter_value = int(maxIter_value)  
9         self.list_num_users = list(map(lambda x: int(x), list_num_users))  
10        self.best_mse_test = float('inf')  
11        self.best_model = None  
12        self.best_params = None  
13  
14    def train_model(self, rank, regParam, maxIter):  
15        als = ALS(rank=rank,  
16                   regParam=regParam,  
17                   maxIter=maxIter,  
18                   userCol="user",  
19                   itemCol="item",  
20                   ratingCol="rating")  
21        model = als.fit(self.train)  
22        return model  
23  
24    def evaluate_model(self, model, num_recommendations):  
25        recommendations = model.recommendForAllItems(num_recommendations)  
26        exploded_recommendations = recommendations.withColumn("recommendation", F.explode("recommendations")) \  
27                           .select("item", F.col("recommendation.user").alias("user"), F.col("recommendation.rating"))  
28        avg_predictions = exploded_recommendations.groupBy("item").agg(F.avg("prediction").alias("avg_prediction"))  
29        avg_ratings = self.test.groupBy("item").agg(F.avg("rating").alias("avg_rating"))  
30        joined_df = exploded_recommendations.join(avg_predictions, "item", "left") \  
31                           .join(avg_ratings, "item", "left") \  
32                           .join(self.test, ["user", "item"], "left")  
33        joined_df = joined_df.withColumn("rating",  
34                                         F.when(F.col("rating").isNull(),  
35                                              F.when(F.col("avg_rating").isNull(), F.col("avg_prediction"))  
36                                              .otherwise(F.col("avg_rating"))  
37                                         )  
38                                         .otherwise(F.col("rating")))  
39        evaluator = RegressionEvaluator(metricName="mse", labelCol="rating", predictionCol="prediction")  
40        mse_test = evaluator.evaluate(joined_df)  
41        return mse_test  
42  
43    def run_experiment(self):  
44        results = []  
45        for numUsers in self.list_num_users:  
46            rank = self.rank_value  
47            regParam = self.regParam_value  
48            maxIter = self.maxIter_value  
49            model = self.train_model(rank, regParam, maxIter)  
50            mse_test = self.evaluate_model(model, numUsers)  
51            results.append((rank, regParam, maxIter, numUsers, mse_test))  
52            print(f'Rank: {rank}, RegParam: {regParam}, MaxIter: {maxIter}, NumUsers: {numUsers}, MSETest: {mse_test}')  
53            if mse_test < self.best_mse_test:  
54                self.best_mse_test = mse_test  
55                self.best_model = model  
56                self.best_params = {'rank': rank, 'regParam': regParam, 'maxIter': maxIter, 'numUsers': numUsers}  
57  
58        df_results = self.spark.createDataFrame(results, ['Rank', 'RegParam', 'MaxIter', 'NumUsers', 'MSETest'])  
59        return df_results  
60  
61    def plot_results(self, df_results):  
62        df_results = df_results.orderBy(F.col('NumUsers')).toPandas()  
63        plt.figure(figsize=(10, 6))  
64        num_users = df_results["NumUsers"]  
65        x = np.arange(len(num_users))  
66        bar_width = 0.35  
67        bars2 = plt.bar(x, df_results["MSETest"], bar_width, label='MSETest', color='skyblue')  
68        plt.title('Mean Squared Error vs Number of Similar Users')  
69        plt.xlabel('Number of Similar Users')  
70        plt.ylabel('Mean Squared Error')  
71        plt.xticks(x, num_users)  
72        for i in range(len(bars2)):  
73            bar = bars2[i]  
74            yval = bar.get_height()  
75            plt.text(  
76                bar.get_x() + bar.get_width() / 2.0,  
77                yval,  
78                round(yval, 3),  
79                ha='center',  
80                va='bottom'  
81            )  
82        plt.legend()  
83        plt.show()  
84  
85    def make_recommendations_user(self, user_id, num_recommendations):  
86        user_recs = self.best_model.recommendForAllUsers(num_recommendations)  
87        user_recs = user_recs.filter(user_recs.user == user_id) \  
88                           .select("user", F.explode("recommendations")) \  
89                           .alias("rec") \  
90                           .select("user", F.col("rec.item").alias("item"), F.col("rec.rating").alias("rating"))  
91        return user_recs  
92  
93    def make_recommendations_item(self, item_id, num_recommendations):  
94        item_recs = self.best_model.recommendForAllItems(num_recommendations)  
95        item_recs = item_recs.filter(item_recs.item == item_id) \  
96                           .select("item", F.explode("recommendations")) \  
97                           .alias("rec") \  
98                           .select("item", F.col("rec.user").alias("user"), F.col("rec.rating").alias("rating"))  
99        return item_recs  
100  
101   def stop_spark(self):
```

```
102     self.spark.stop()
```

```
1 filepath = '/content/drive/MyDrive/BigData/CK/datasets/ratings2k.csv'
2 cf = CollaborativeFiltering(filepath=filepath)
3 df_results = cf.run_experiment()
4
5 # rank_value = 20
6 # regParam_value = 0.5
7 # maxIter_value = 5
8 # list_num_users = list(range(20,31))
9 # cf = CollaborativeFiltering(filepath=filepath, rank_value=rank_value, regParam_value=regParam_value, maxIter_value=maxIter_value, list_
10 # df_results = cf.run_experiment()
```

```
→ Rank: 10, RegParam: 0.1, MaxIter: 10, NumUsers: 10, MSETest: 0.4940042546096179
Rank: 10, RegParam: 0.1, MaxIter: 10, NumUsers: 11, MSETest: 0.488957627048654
Rank: 10, RegParam: 0.1, MaxIter: 10, NumUsers: 12, MSETest: 0.48441952358985496
Rank: 10, RegParam: 0.1, MaxIter: 10, NumUsers: 13, MSETest: 0.4813058037642575
Rank: 10, RegParam: 0.1, MaxIter: 10, NumUsers: 14, MSETest: 0.4783468208341385
Rank: 10, RegParam: 0.1, MaxIter: 10, NumUsers: 15, MSETest: 0.4767723293577801
Rank: 10, RegParam: 0.1, MaxIter: 10, NumUsers: 16, MSETest: 0.47549411717843354
Rank: 10, RegParam: 0.1, MaxIter: 10, NumUsers: 17, MSETest: 0.4754643987717024
Rank: 10, RegParam: 0.1, MaxIter: 10, NumUsers: 18, MSETest: 0.4749176913611994
Rank: 10, RegParam: 0.1, MaxIter: 10, NumUsers: 19, MSETest: 0.47464702807694176
Rank: 10, RegParam: 0.1, MaxIter: 10, NumUsers: 20, MSETest: 0.47488509599190015
```

```
1 cf.plot_results(df_results)
```



```
1 print(f'Best params: {cf.best_params}, Best MSE: {cf.best_mse_test}')
2 res = cf.make_recommendations_user(user_id=55, num_recommendations=20)
3 res.show()
```

```
→ Best params: {'rank': 10, 'regParam': 0.1, 'maxIter': 10, 'numUsers': 19}, Best MSE: 0.47464702807694176
+---+---+-----+
|user|item| rating|
+---+---+-----+
| 55 | 352 | 4.4253345|
| 55 | 335 | 3.9472497|
| 55 | 422 | 3.6893568|
| 55 | 199 | 3.4357786|
| 55 | 456 | 3.4119596|
| 55 | 196 | 3.370279|
| 55 | 95 | 3.3507597|
| 55 | 257 | 3.3421004|
| 55 | 8 | 3.3195114|
| 55 | 439 | 3.3124905|
| 55 | 437 | 3.285025|
| 55 | 239 | 3.2773952|
| 55 | 438 | 3.2386374|
| 55 | 262 | 3.2120264|
| 55 | 251 | 3.1818209|
| 55 | 356 | 3.1670544|
| 55 | 464 | 3.1670544|
| 55 | 454 | 3.1670544|
| 55 | 12 | 3.109748|
| 55 | 15 | 3.0618691|
+---+---+-----+
```

```
1 print(f'Best params: {cf.best_params}, Best MSE: {cf.best_mse_test}')
2 res = cf.make_recommendations_item(item_id=352, num_recommendations=20)
3 res.show()
```

```
→ Best params: {'rank': 10, 'regParam': 0.1, 'maxIter': 10, 'numUsers': 19}, Best MSE: 0.47464702807694176
+---+---+-----+
|item|user| rating|
+---+---+-----+
| 352 | 1 | 4.8855176|
| 352 | 32 | 4.453047|
| 352 | 55 | 4.4253345|
| 352 | 41 | 4.3970222|
| 352 | 50 | 4.2273903|
```

```

| 352| 72| 4.015167|
| 352| 71| 3.9610286|
| 352| 69| 3.9535818|
| 352| 56| 3.929315|
| 352| 64| 3.9155207|
| 352| 46| 3.3819664|
| 352| 21| 3.3557096|
| 352| 14| 3.189015|
| 352| 35| 3.1153104|
| 352| 49| 3.0909026|
| 352| 42| 2.9488816|
| 352| 22| 2.8665767|
| 352| 9| 2.8573644|
| 352| 33| 2.8508577|
| 352| 48| 2.84253|
+---+---+-----+

```

## ✓ Câu 4: Dự đoán giá chứng khoán

```

1 class StockData:
2     def __init__(self, file_path, k=3):
3         self.file_path = file_path
4         self.spark = SparkSession.builder.appName("StockPrediction").getOrCreate()
5         self.k = k
6         self.data = None
7         self.new_data = None
8         self.df_with_vectors = None
9         self.window_spec = Window.orderBy(F.col("year").desc(), F.col("month").desc(), F.col("day").desc())
10
11    def load_and_prepare_data(self):
12        data = self.spark.read.csv(self.file_path, header=True, inferSchema=True)
13        split_date = F.split(data["Ngay"], "/")
14        data = data.withColumn("day", split_date.getItem(0).cast("int")) \
15            .withColumn("month", split_date.getItem(1).cast("int")) \
16            .withColumn("year", split_date.getItem(2).cast("int")) \
17            .select("day", "month", "year", "HVN")
18        return data
19
20    def calculate_fluctuation(self, data):
21        lead_price = F.lead(F.col("HVN"), 1).over(self.window_spec)
22        fluctuation = F.coalesce((F.col("HVN") - lead_price) / lead_price, F.lit(0.0))
23        data = data.withColumn("fluctuation", fluctuation)
24        return data
25
26    def calculate_amplitudes(self, data):
27        window_spec = Window.orderBy(F.col("year").desc(), F.col("month").desc(), F.col("day").desc())
28        win_spec_k = window_spec.rowsBetween(1, self.k)
29        new_data = data.withColumn("k_day_amplitude", F.reverse(F.collect_list("fluctuation").over(win_spec_k)))
30        new_data = new_data.withColumn("k_day_amplitude",
31            F.when(F.size(new_data["k_day_amplitude"]) < self.k,
32                F.concat(F.array_repeat(F.lit(0), self.k - F.size(new_data["k_day_amplitude"])), new_data["k"]
33                .otherwise(new_data["k_day_amplitude"])))
34
35        return new_data
36
37
38    def transform_to_vector(self, new_data):
39        def list_to_vector_udf_nested(l):
40            return Vectors.dense(l)
41
42        list_to_vector_udf_nested = F.udf(list_to_vector_udf_nested, VectorUDT())
43        df_with_vectors = new_data.withColumn('last_k_fluctuation', list_to_vector_udf_nested(new_data['k_day_amplitude'])) \
44            .select("last_k_fluctuation", "fluctuation", "month")
45        return df_with_vectors
46
47    def split_data(self, data):
48        train_data = data.filter(F.col("month") <= 6).select("last_k_fluctuation", "fluctuation")
49        test_data = data.filter(F.col("month") > 6).select("last_k_fluctuation", "fluctuation")
50        return train_data, test_data
51
52    def create_data(self):
53        self.data = self.load_and_prepare_data()
54        self.data = self.calculate_fluctuation(self.data)
55        self.new_data = self.calculate_amplitudes(self.data)
56        self.df_with_vectors = self.transform_to_vector(self.new_data)
57
58    def get_data(self):
59        if self.df_with_vectors is None:
60            self.create_data()
61        return self.df_with_vectors
62
63
64
65 # Usage
66 file_path = path+'stockHVN2022.csv'
67 stock_data = StockData(file_path)
68 stock_data.create_data()
69 data = stock_data.get_data()
70 train_data, test_data = stock_data.split_data(data)
71 print("Train data:")
72 train_data.show(5, truncate=False)
73 print("Test data:")
74 test_data.show(5, truncate=False)

→ Train data:
+-----+-----+
|last_k_fluctuation|fluctuation|
+-----+-----+

```

```
+-----+
|[0.06885245901639349, 0.012269938650306704, -0.027272727272723] | -0.034267912772585715|
|[-0.003267973856209197, 0.06885245901639349, 0.012269938650306704] | -0.027272727272723 |
|[-0.003257328990227944, -0.003267973856209197, 0.06885245901639349] | 0.012269938650306704 |
|[0.02675585284280939, -0.003257328990227944, -0.003267973856209197] | 0.06885245901639349 |
|[0.023972602739726005, 0.02675585284280939, -0.003257328990227944] | -0.003267973856209197|
+-----+
only showing top 5 rows
```

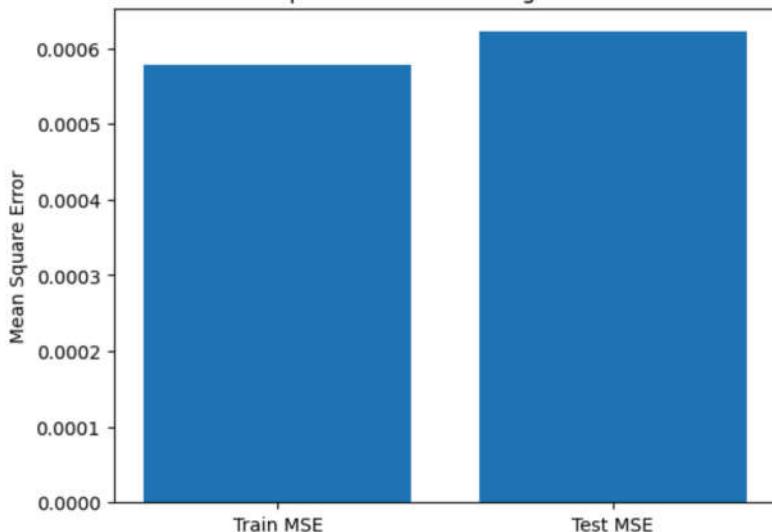
Test data:

```
+-----+
|last_k_fluctuation | fluctuation |
+-----+
|[-0.06918918918918925, 0.06968641114982595, 0.02280130293159599] | -0.01273885350318463|
|[-0.06565656565656569, -0.06918918918918925, 0.06968641114982595] | 0.02280130293159599 |
|[0.014344262295082027, -0.06565656565656569, -0.06918918918918925] | 0.06968641114982595 |
|[-0.06153846153846159, 0.014344262295082027, -0.065656565656569] | -0.06918918918918925 |
|[0.014634146341463448, -0.06153846153846159, 0.014344262295082027] | -0.065656565656569 |
+-----+
only showing top 5 rows
```

```
1 class StockModel:
2     def __init__(self, train_data, test_data):
3         self.train_data = train_data
4         self.test_data = test_data
5         self.lr_model = None
6
7     def train(self):
8         lr = LinearRegression(featuresCol="last_k_fluctuation", labelCol="fluctuation")
9         self.lr_model = lr.fit(self.train_data)
10
11    def evaluate(self):
12        train_predictions = self.lr_model.transform(self.train_data)
13        test_predictions = self.lr_model.transform(self.test_data)
14
15        evaluator = RegressionEvaluator(labelCol="fluctuation", predictionCol="prediction", metricName="mse")
16        train_mse = evaluator.evaluate(train_predictions)
17        test_mse = evaluator.evaluate(test_predictions)
18
19        return train_mse, test_mse
20
21    def plot_mse(self, train_mse, test_mse):
22        plt.bar(['Train MSE', 'Test MSE'], [train_mse, test_mse])
23        plt.ylabel('Mean Square Error')
24        plt.title('Mean Square Error on Training and Test Sets')
25        plt.show()
26
27
28 stock_model = StockModel(train_data, test_data)
29 stock_model.train()
30 train_mse, test_mse = stock_model.evaluate()
31 print("Train MSE:", train_mse)
32 print("Test MSE:", test_mse)
33 stock_model.plot_mse(train_mse, test_mse)
```

Train MSE: 0.0005788157807408016  
Test MSE: 0.0006217807601481779

Mean Square Error on Training and Test Sets



## ✓ Câu 5: Phân loại đa lớp với pyspark

```
1 class MNISTClassifier:
2     def __init__(self, spark):
3         self.spark = spark
4
5     def load_data(self, file_path):
6         """
7             Loads the MNIST dataset from a CSV file using SparkSession's read.csv() method.
8         """
9         self.data = self.spark.read.csv(file_path, header=False, inferSchema=True)
10
11    def preprocess_data(self):
12        """
13            Preprocesses the dataset by renaming the first column to 'label' and creating a feature vector using VectorAssembler.
14        """
15        col4_data = col4_data.withColumnRenamed('0', 'label')
```

```
42     self.data = self.data.withColumnRenamed('label', 'label')
43
44     assembler = VectorAssembler(inputCols=self.data.columns[1:], outputCol="features")
45     self.data = assembler.transform(self.data).select('label', 'features')
46
47
48     def split_data(self, train_ratio):
49         """
50             Splits the dataset into training and testing sets based on the provided ratio.
51         """
52         return self.data.randomSplit([train_ratio, 1 - train_ratio])
53
54
55     def train_model(self, model, train_data):
56         """
57             Trains a machine learning model on the training data using Spark ML fit() method.
58         """
59         return model.fit(train_data)
60
61
62     def evaluate_accuracy(self, model, data):
63         """
64             Evaluates the accuracy of a model on the given data using MulticlassClassificationEvaluator with accuracy metric.
65         """
66         predictions = model.transform(data)
67         evaluator = MulticlassClassificationEvaluator(metricName="accuracy", labelCol='label')
68         accuracy = evaluator.evaluate(predictions)
69         return accuracy
70
71
72     def evaluate_log_loss(self, model, data):
73         """
74             Evaluates the log loss of a model on the given data using MulticlassClassificationEvaluator with logLoss metric.
75         """
76         predictions = model.transform(data)
77         evaluator = MulticlassClassificationEvaluator(metricName="logLoss", labelCol='label')
78         log_loss = evaluator.evaluate(predictions)
79         return log_loss
80
81
82     def plot_accuracy(self, accuracies_train, accuracies_test):
83         """
84             Plots the accuracy of different models on the training and testing data using matplotlib.
85         """
86         models = ['OneVsRest Linear SVM', 'MLP', 'Random Forest']
87         x = range(len(models))
88
89         plt.bar(x, accuracies_train, width=0.4, label='Train', align='center')
90         plt.bar(x, accuracies_test, width=0.4, label='Test', align='edge')
91
92         plt.xlabel('Models')
93         plt.ylabel('Accuracy')
94         plt.title('Model Accuracy Comparison')
95         plt.xticks(x, models)
96         plt.legend()
97         plt.show()
98
99
100    def run(self, file_path, train_ratio=0.8):
101        """
102            Runs the entire MNISTClassifier pipeline, including data loading, preprocessing, training, evaluation, and plotting.
103        """
104        self.load_data(file_path)
105        self.preprocess_data()
106        train_data, test_data = self.split_data(train_ratio)
107
108        mlp = MultilayerPerceptronClassifier(layers=[784, 256, 128, 10], seed=123, featuresCol='features', labelCol='label')
109        rf = RandomForestClassifier(numTrees=100, seed=123, featuresCol='features', labelCol='label')
110        lsvc = LinearSVC(maxIter=10, regParam=0.1, featuresCol='features', labelCol='label')
111        lsvc_ovr = OneVsRest(classifier=lsvc, featuresCol='features', labelCol='label')
112        models = [lsvc_ovr, mlp, rf]
113
114        accuracies_train = []
115        log_losses_train = []
116        accuracies_test = []
117        log_losses_test = []
118
119        for model in models:
120            trained_model = self.train_model(model, train_data)
121
122            train_accuracy = self.evaluate_accuracy(trained_model, train_data)
123            test_accuracy = self.evaluate_accuracy(trained_model, test_data)
124            accuracies_train.append(train_accuracy)
125            accuracies_test.append(test_accuracy)
126
127            if model != lsvc_ovr and model != rf:
128                train_log_loss = self.evaluate_log_loss(trained_model, train_data)
129                test_log_loss = self.evaluate_log_loss(trained_model, test_data)
130                log_losses_train.append(train_log_loss)
131                log_losses_test.append(test_log_loss)
132                print(f"Train Log Loss ({model.__class__.__name__}): {train_log_loss}")
133                print(f"Test Log Loss ({model.__class__.__name__}): {test_log_loss}")
134
135            print(f"Train Accuracy ({model.__class__.__name__}): {train_accuracy}")
136            print(f"Test Accuracy ({model.__class__.__name__}): {test_accuracy}")
137
138        self.plot_accuracy(accuracies_train, accuracies_test)
139
140
141    if __name__ == "__main__":
142        data_path_mnist_mini = "/content/drive/MyDrive/Big_Data/datasets/mnist_mini.csv"
143        classifier = MNISTClassifier(spark)
144        classifier.run(data_path_mnist_mini)
```

Train Accuracy (OneVsRest): 0.877333667460218  
Test Accuracy (OneVsRest): 0.8628033680039624  
Train Log Loss (MultilayerPerceptronClassifier): 0.0999539981884838  
Test Log Loss (MultilayerPerceptronClassifier): 0.2599573596454811  
Train Accuracy (MultilayerPerceptronClassifier): 0.9700538779601554  
Test Accuracy (MultilayerPerceptronClassifier): 0.9262010896483408  
Train Accuracy (RandomForestClassifier): 0.8758300964791379  
Test Accuracy (RandomForestClassifier): 0.8553739474987617

