



**Title:** Report

**Name:** Umm-e-Hani , Laiba Batool , Dawood Tanvir

**Roll No.:** 21I-1715 , 21I-1781, 21I-1665

**Section:** M

# MODEL TRAINING:

We tried training the model on the whole dataset but that wasn't possible. I kept getting the error given below. Tried to train the dataset on 50% data still the same error. I tried multiple solutions but it was giving a memory error. After many tries I decided to pick up on 1 million rows and still the same error so I had to train the model on an even smaller dataset.

```
-----
Py4JJavaError                                Traceback (most recent call last)
Cell In[8], line 1
----> 1 transformed = pipeline.fit(sampled_df).transform(sampled_df)

File C:\spark\python\pyspark\ml\base.py:205, in Estimator.fit(self, dataset, params)
    203     return self.copy(params)._fit(dataset)
    204     else:
--> 205         return self._fit(dataset)
    206     else:
    207         raise TypeError(
    208             "Params must be either a param map or a list/tuple of param maps, "
    209             "but got %s." % type(params)
    210         )

File C:\spark\python\pyspark\ml\pipeline.py:134, in Pipeline._fit(self, dataset)
    132     dataset = stage.transform(dataset)
    133     else: # must be an Estimator
--> 134     model = stage.fit(dataset)
    135     transformers.append(model)
    136     if i < indexOfLastEstimator:

File C:\spark\python\pyspark\ml\base.py:205, in Estimator.fit(self, dataset, params)
    203     return self.copy(params)._fit(dataset)
    204     else:
--> 205         return self._fit(dataset)
    206     else:
    207         raise TypeError(
    208             "Params must be either a param map or a list/tuple of param maps, "
    209             "but got %s." % type(params)
    210         )

File C:\spark\python\pyspark\ml\wrapper.py:383, in JavaEstimator._fit(self, dataset)
    382     def _fit(self, dataset: DataFrame) -> JM:
--> 383         java_model = self._fit_java(dataset)
    384         model = self._create_model(java_model)
    ...
```

# KAFKA CONNECTION AND MONGODB:

This code is a Flask-based web application that uses Apache Spark and Kafka to provide product recommendations to users based on their reviews. The application provides a web form where users can input their review, and then uses Spark to recommend products to the user. The recommended products are stored in a MongoDB database, and a Kafka producer sends the recommendations to the Kafka consumer, which displays the recommended products to the user in real-time. The code uses a pre-trained machine learning model built using the ALS algorithm for collaborative filtering to recommend products to users. The application has several dependencies, including the Flask, Kafka, and PySpark libraries. It also requires a MongoDB database and a pre-trained ALS model. The code has two routes - the first route displays the web form, while the second route displays the recommended products. The first route sends a post request with the user's review, and the application uses Spark to recommend products to the user. The recommended products are stored in the MongoDB database, and the Kafka producer sends the recommendations to the Kafka consumer. The Kafka consumer then displays the recommended products to the user on the second route. Overall, this Flask-based application provides an effective way to recommend Amazon products to users based on their reviews, using powerful technologies like Spark and Kafka for large-scale data processing and real-time communication

```
1  from flask import Flask, render_template, request
2  import warnings
3  from flask import redirect
4  from pymongo import MongoClient, collection
5
6  import numpy as np
7  from kafka import KafkaProducer
8  from kafka import KafkaConsumer
9
10 warnings.simplefilter("ignore", UserWarning)
11 warnings.filterwarnings('ignore')
12 import pyspark
13 import pyspark
14 from pyspark.sql import SparkSession
15 from pyspark.sql.types import StructType, StructField, StringType, DoubleType, BooleanType, LongType, IntegerType
16 from pyspark.sql.functions import col, sum
17 from pyspark.sql.functions import mean, min, max
18 from pyspark.sql.functions import from_unixtime
19 import numpy
20 from pyspark.sql.functions import udf
21 from pyspark.sql.types import StringType
22 import pyspark.sql.functions as F
23 from pyspark.ml.feature import StringIndexer
24 import os
25 import sys
26
27 os.environ['PYSPARK_PYTHON'] = sys.executable
28 os.environ['PYSPARK_DRIVER_PYTHON'] = sys.executable
29
30 # spark = SparkSession.builder.appName("Amazon Reviews").getOrCreate()
31 # https://img.freepik.com/free-vector/musical-pentagram-sound-waves-notes-background_1017-33911.jpg?w=2000
```

```

25 import sys
26
27 os.environ['PYSPARK_PYTHON'] = sys.executable
28 os.environ['PYSPARK_DRIVER_PYTHON'] = sys.executable
29
30 # spark = SparkSession.builder.appName("Amazon Reviews").getOrCreate()
31 # https://img.freepik.com/free-vector/musical-pentagram-sound-waves-notes-background_1017-33911.jpg?w=2000
32 spark = SparkSession.builder.appName("Amazon Reviews") \
33     .config("spark.driver.memory", "4g") \
34     .config("spark.executor.memory", "4g") \
35     .config("spark.executor.heartbeatInterval", "5000s") \
36     .config("spark.network.timeout", "10000s") \
37     .getOrCreate()
38
39 mongo_client = MongoClient('mongodb://localhost:27017')
40 db = mongo_client['RECOMMENDATIONS']
41 collection = db['User_Recommendation']
42 app = Flask(__name__)
43
44 recommendations = []
45
46
47 @app.route('/', methods=['GET'])
48 def home():
49     return render_template('review.html')
50
51
52 @app.route('/', methods=['POST'])
53 def index():
54     from pyspark.ml.recommendation import ALS, ALSModel
55     models = ALSModel.load("E:\\AmazonResc")

```

```

49     return render_template('review.html')
50
51
52 @app.route('/', methods=['POST'])
53 def index():
54     from pyspark.ml.recommendation import ALS, ALSModel
55     models = ALSModel.load("E:\\AmazonResc")
56     user_id = request.form['review']
57
58     user_data = spark.createDataFrame([(user_id,)], ["Revindexed"])
59
60     userRecs = models.recommendForUserSubset(user_data, numItems=5) # Specify the number of recommendations
61     rec = userRecs.collect()
62     producer = KafkaProducer(bootstrap_servers=['localhost:9092'])
63
64
65     recommendations.append(str(rec[0][1][1][0]))
66     recommendations.append(str(rec[0][1][2][0]))
67     recommendations.append(str(rec[0][1][3][0]))
68     recommendations.append(str(rec[0][1][4][0]))
69
70     document = {
71         'user_id': user_id,
72         'recommendations': recommendations
73     }
74     #insert document in mongodb
75     collection.insert_one(document)
76     recommendations.append(user_id)
77
78
79     for recommendation in recommendations:

```

```

49     return render_template('review.html')
50
51
52 @app.route('/', methods=['POST'])
53 def index():
54     from pyspark.ml.recommendation import ALS, ALSModel
55     models = ALSModel.load("E:\\AmazonRec")
56     user_id = request.form['review']
57
58     user_data = spark.createDataFrame([(user_id,)], ["Revindexed"])
59
60     userRecs = models.recommendForUserSubset(user_data, numItems=5) # Specify the number of recommendations
61     rec = userRecs.collect()
62     producer = KafkaProducer(bootstrap_servers=['localhost:9092'])
63
64
65     recommendations.append(str(rec[0][1][1][0]))
66     recommendations.append(str(rec[0][1][2][0]))
67     recommendations.append(str(rec[0][1][3][0]))
68     recommendations.append(str(rec[0][1][4][0]))
69
70     document = {
71         'user_id': user_id,
72         'recommendations': recommendations
73     }
74     #insert document in mongodb
75     collection.insert_one(document)
76     recommendations.append(user_id)
77
78
79     for recommendation in recommendations:

```

```

76     recommendations.append(user_id)
77
78
79     for recommendation in recommendations:
80         producer.send('Rec', value=recommendation.encode())
81
82     return redirect('/result')
83
84
85 @app.route('/result')
86 def display_result():
87     consumer = KafkaConsumer(
88         'Rec',
89         bootstrap_servers='localhost:9092',
90         auto_offset_reset='earliest',
91         group_id='i23'
92     )
93     i = 0
94     rec = []
95     # Poll for new messages
96     for message in consumer:
97         value = message.value.decode()
98         rec.append(value)
99         consumer.commit_async()
100         i += 1
101         if i == 4:
102             break
103
104     return render_template("outputRec.html", a=rec[0], b=rec[1], c=rec[2], d=rec[3])
105
106

```

## Write your review

Name

Email

Your ID

Submit



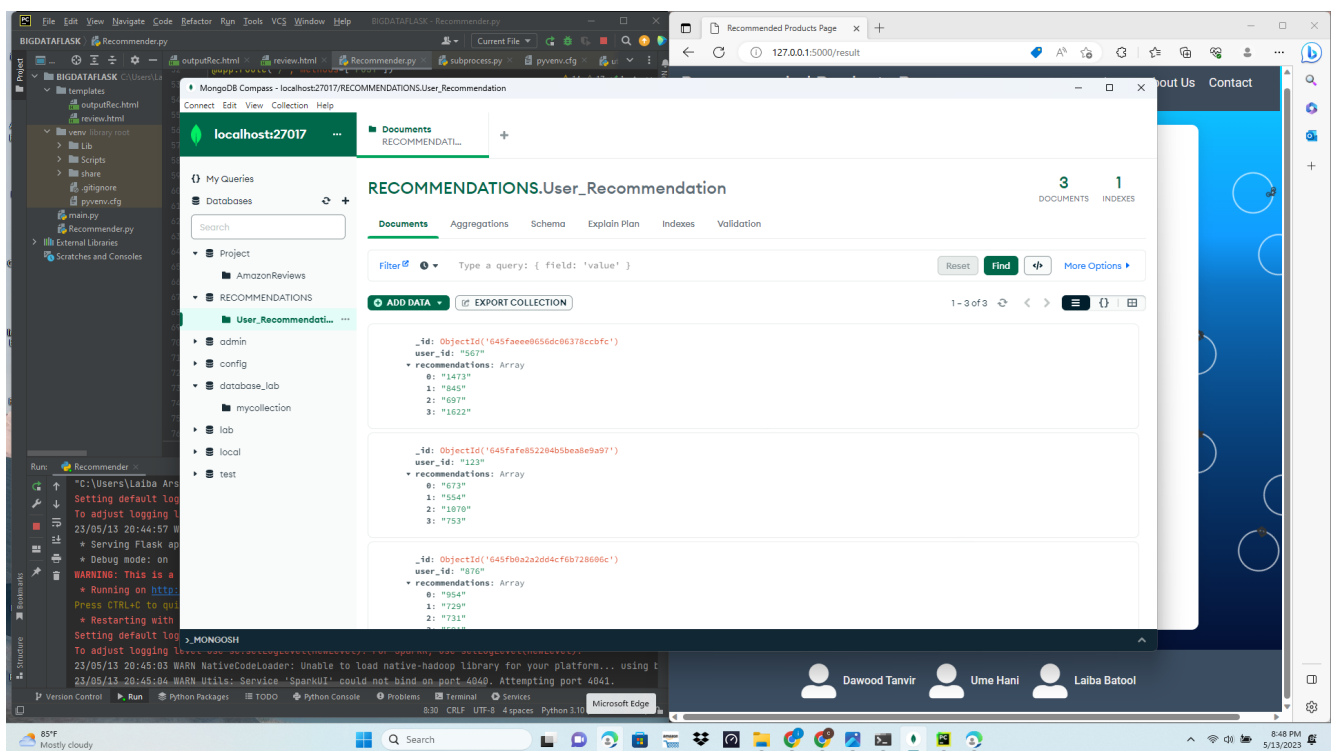
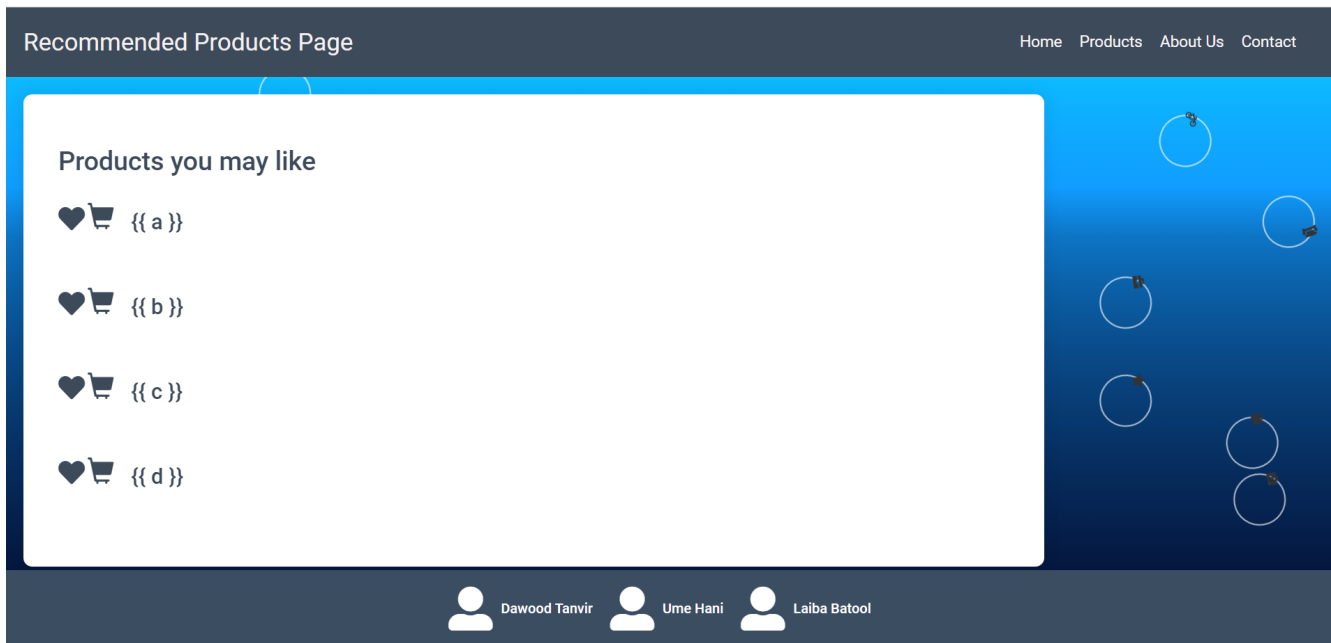
Dawood Tanvir



Ume Hani



Laiba Batool



## CONTRIBUTIONS:

UMM E HANI : MODEL TRAINING

DAWOOD & LAIBA : MONGO AND KAFKA SETUP