

```
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 #!wget -q http://archive.apache.org/dist/spark/spark-3.1.1/spark-3.1.1-bin-hadoop3.2.tgz
2 !tar xF /content/drive/MyDrive/spark-3.1.1-bin-hadoop3.2.tgz
3 !pip install -q findspark

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 import pyspark as spark
2 print(spark.__version__)

    3.1.1

1 from pyspark import SparkContext,SparkConf
2 from pyspark.sql import SparkSession
3 from pyspark.sql.types import StringType
4 from pyspark.ml.fpm import FPGrowth
5 import pyspark.sql.functions as F
6 import matplotlib.pyplot as plt
7 import findspark
8 import os
9 import pyspark.sql.functions as F
10 from pyspark.sql.window import Window
11 from itertools import combinations
12 from functools import reduce

1 conf = SparkConf().setAppName("gk")
2 sc = SparkContext.getOrCreate(conf=conf)

1 def f1(path):
2     rdd = sc.textFile(path)
3     header = rdd.first()
4     rdd = rdd.filter(lambda line: line != header)
5
6     items_rdd = rdd.map(lambda line: line.split(','))
7     unique_items = items_rdd.map(lambda x: x[2]).distinct().sortBy(lambda x: x.lower())
8
9     output_path = "f1"
10    hadoop_conf = sc._jsc.hadoopConfiguration()
11    fs = sc._jvm.org.apache.hadoop.fs.FileSystem.get(hadoop_conf)
12    if fs.exists(sc._jvm.org.apache.hadoop.fs.Path(output_path)):
13        fs.delete(sc._jvm.org.apache.hadoop.fs.Path(output_path), True)
14
15    unique_items.coalesce(1).saveAsTextFile(output_path)
16
17    # In ra 10 mục hàng đầu tiên
18    first_10_names = unique_items.take(10)
19    for name in first_10_names:
20        print(name)
21
22    # In ra 10 mục hàng cuối cùng
23    last_10_names = unique_items.sortBy(lambda x: x.lower(), ascending=False).take(10)[::-1]
24    print("-----")
25    for name in last_10_names:
26        print(name)

1 def f2(path):
2     rdd = sc.textFile(path)
3     header = rdd.first()
4     rdd = rdd.filter(lambda line: line != header)
5
6     items_rdd = rdd.map(lambda line: line.split(','))
7     item_counts = items_rdd.map(lambda x: (x[2], 1)).reduceByKey(lambda a, b: a + b)
8
9     # Sắp xếp các mục hàng theo số lần xuất hiện giảm dần
10    sorted_items = item_counts.sortBy(lambda x: x[1], ascending=False)
11
12    output_path = "f2"
13    hadoop_conf = sc._jsc.hadoopConfiguration()
14    fs = sc._jvm.org.apache.hadoop.fs.FileSystem.get(hadoop_conf)
15    if fs.exists(sc._jvm.org.apache.hadoop.fs.Path(output_path)):
16        fs.delete(sc._jvm.org.apache.hadoop.fs.Path(output_path), True)
17
18    sorted_items.coalesce(1).saveAsTextFile(output_path)
19
20    top_items = sorted_items.take(10)
21    for item in top_items:
22        print(item)
23
24    # Lấy top 100 mục hàng được mua nhiều nhất
25    top_100_items = sorted_items.take(100)
26    item_names = [item[0] for item in top_100_items]
27    item_counts = [item[1] for item in top_100_items]
28
29    plt.figure(figsize=(15, 6))
30    plt.bar(item_names, item_counts)
31    plt.xlabel('Món hàng')
32    plt.ylabel('Số lần mua')
33    plt.title('Top 100 món hàng được mua nhiều nhất')
34    plt.xticks(rotation=90)
35    plt.show()

1 def f3(path):
2     rdd = sc.textFile(path)
3     header = rdd.first()
4     rdd = rdd.filter(lambda line: line != header)
5
6     data_rdd = rdd.map(lambda line: line.split(','))
7     user_counts = data_rdd.map(lambda x: ((x[0], x[1]), 1)).reduceByKey(lambda a, b: a) \
8     .map(lambda x: (x[0][0], 1)).reduceByKey(lambda a, b: a + b)
9     sorted_users = user_counts.sortBy(lambda x: x[1], ascending=False)
10
11    output_path = "f3"
12    hadoop_conf = sc._jsc.hadoopConfiguration()
13    fs = sc._jvm.org.apache.hadoop.fs.FileSystem.get(hadoop_conf)
14    if fs.exists(sc._jvm.org.apache.hadoop.fs.Path(output_path)):
15        fs.delete(sc._jvm.org.apache.hadoop.fs.Path(output_path), True)
16
17    sorted_users.coalesce(1).saveAsTextFile(output_path)
18
19    top_items = sorted_users.take(10)
20    for item in top_items:
21        print(item)
22
23    top_users = sorted_users.take(100)
24    user_names = [user[0] for user in top_users]
25    basket_counts = [user[1] for user in top_users]
26
27    plt.figure(figsize=(20, 6))
28    plt.bar(user_names, basket_counts)
29    plt.xlabel('Người dùng')
30    plt.ylabel('Số lượng giỏ hàng')
31    plt.title('Top 100 người dùng mua nhiều giỏ hàng nhất')
32    plt.xticks(rotation=90)
33    plt.show()
```

```
1 def f4(path):
2     rdd = sc.textFile(path)
3     header = rdd.first()
4     rdd = rdd.filter(lambda line: line != header)
5
6     data_rdd = rdd.map(lambda line: line.split(','))
7     user_item_counts = data_rdd.map(lambda x: (x[0], x[2])).distinct().map(lambda x: (x[0], 1)).reduceByKey(lambda a, b: a + b)
8     user_most_items = user_item_counts.sortBy(lambda x: x[1], ascending=False).first()
9
10    print("Người dùng mua nhiều món hàng phân biệt nhất:")
11    print("Mã người dùng:", user_most_items[0])
12    print("Số lượng món hàng:", user_most_items[1])
13
14    item_user_counts = data_rdd.map(lambda x: (x[2], x[0])).distinct().map(lambda x: (x[2], 1)).reduceByKey(lambda a, b: a + b)
15    item_most_users = item_user_counts.sortBy(lambda x: x[1], ascending=False).first()
16
17    print("Món hàng được mua bởi nhiều người dùng nhất:")
18    print("Tên món hàng:", item_most_users[0])
19    print("Số lượng người mua:", item_most_users[1])
20
21    output_path = "f4"
22    if sc._jvm.org.apache.hadoop.fs.FileSystem.get(sc._jsc.hadoopConfiguration()).exists(sc._jvm.org.apache.hadoop.fs.Path(output_path)):
23        sc._jvm.org.apache.hadoop.fs.FileSystem.get(sc._jsc.hadoopConfiguration()).delete(sc._jvm.org.apache.hadoop.fs.Path(output_path), True)
24    output_rdd = sc.parallelize(user_most_items + item_most_users)
25    output_rdd.coalesce(1).saveAsTextFile(output_path)
```

1 path = "/content/drive/MyDrive/baskets.csv"

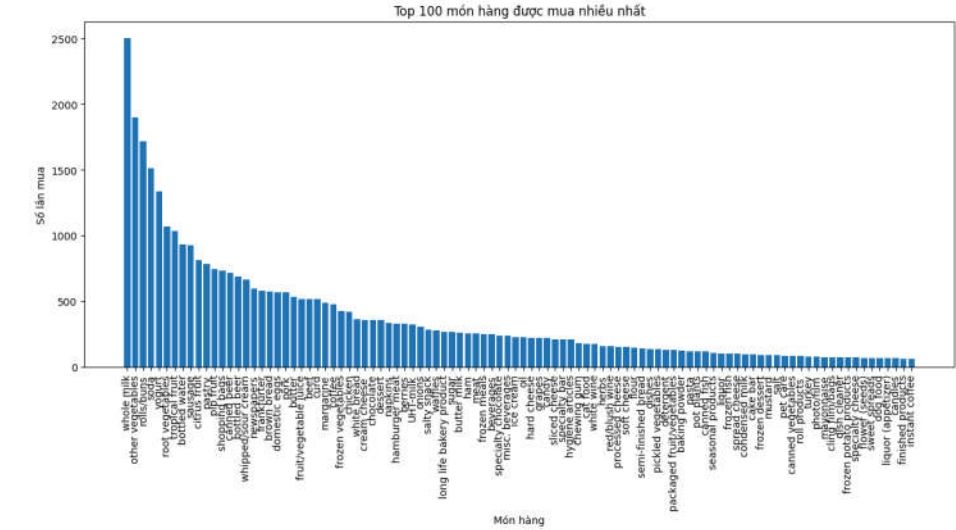
1 f1(path)

abrasive cleaner
artif. sweetener
baby cosmetics
bags
baking powder
bathroom cleaner
beef
berries
beverages
bottled beer

UHT-milk
vinegar
waffles
whipped/sour cream
whisky
white bread
white wine
whole milk
yogurt
zwieback

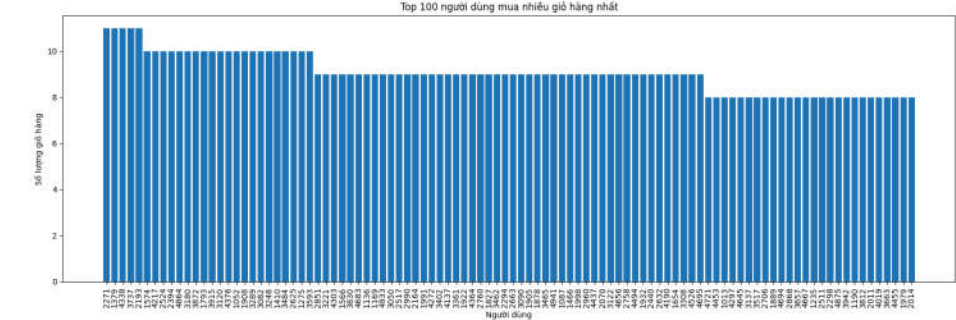
1 f2(path)

('whole milk', 2502)
('other vegetables', 1898)
('rolls/buns', 1716)
('soda', 1514)
('yogurt', 1334)
('root vegetables', 1071)
('tropical fruit', 1032)
('bottled water', 933)
('sausage', 924)
('citrus fruit', 812)



1 f3(path)

('2271', 11)
('1379', 11)
('4338', 11)
('3737', 11)
('2193', 11)
('1574', 10)
('4217', 10)
('2524', 10)
('2394', 10)
('4864', 10)



1 f4(path)

Người dùng mua nhiều món hàng phân biệt nhất:
Mã người dùng: 2051
Số lượng món hàng: 26
Món hàng được mua bởi nhiều người dùng nhất:
Tên món hàng: whole milk
Số lượng người mua: 1786

```
1 sc.stop()

1 !sudo rm -r /content/baskets

1 spark = SparkSession.builder.getOrCreate()
2
3 df = spark.read.csv(path, header=True)
4 df_basket = df.groupBy("Member_number", "Date", "year", "month", "day").agg(F.collect_set("itemDescription").alias("Basket"))
5
6 df_basket = df_basket.withColumn("year", F.col("year").cast("integer"))
7 df_basket = df_basket.withColumn("month", F.col("month").cast("integer"))
8 df_basket = df_basket.withColumn("day", F.col("day").cast("integer"))
9
10 df_basket = df_basket.orderBy(["year", "month", "day"], ascending=True)
11 df_basket.show(truncate=False)
12
13 df_Basket = df_basket.select("Basket")
14 basket_string = df_Basket.select(F.concat_ws(" ", "Basket").alias("BasketString"))
15 basket_string.coalesce(1).write.csv("baskets")

+-----+-----+-----+-----+-----+
|Member_number|Date      |year|month|day|Basket|
+-----+-----+-----+-----+-----+
|1789          |01/01/2014|2014|1     |1  |[candies, hamburger meat]|
|2789          |01/01/2014|2014|1     |1  |[yogurt, frozen vegetables]|
|2943          |01/01/2014|2014|1     |1  |[whole milk, flower (seeds)]|
|4942          |01/01/2014|2014|1     |1  |[butter, frozen vegetables]|
|1922          |01/01/2014|2014|1     |1  |[tropical fruit, other vegetables]|
|2237          |01/01/2014|2014|1     |1  |[Instant food products, bottled water]|
|1249          |01/01/2014|2014|1     |1  |[citrus fruit, coffee]|
|3797          |01/01/2014|2014|1     |1  |[whole milk, waffles]|
|1440          |01/01/2014|2014|1     |1  |[yogurt, other vegetables]|
|4260          |01/01/2014|2014|1     |1  |[soda, brown bread]|
|2727          |01/01/2014|2014|1     |1  |[hamburger meat, frozen potato products]|
|1381          |01/01/2014|2014|1     |1  |[curd, soda]|
|1659          |01/01/2014|2014|1     |1  |[specialty chocolate, frozen vegetables]|
|2542          |01/01/2014|2014|1     |1  |[bottled water, sliced cheese]|
|3956          |01/01/2014|2014|1     |1  |[yogurt, shopping bags, waffles, chocolate]|
|2610          |01/01/2014|2014|1     |1  |[domestic eggs, bottled beer, hamburger meat]|
|3681          |01/01/2014|2014|1     |1  |[dishes, onions, whipped/sour cream]|
|2974          |01/01/2014|2014|1     |1  |[bottled water, berries, whipped/sour cream]|
|2226          |01/01/2014|2014|1     |1  |[sausage, bottled water]|
|2351          |01/01/2014|2014|1     |1  |[shopping bags, cleaner]|
+-----+-----+-----+-----+-----+
only showing top 20 rows

1 # Tính số lượng giỏ hàng theo ngày
2 df_count = df_basket.groupBy("Date", "year", "month", "day").count()
3 df_count = df_count.orderBy(["year", "month", "day"], ascending=True)
4
5 pandas_df = df_count.toPandas()
6 plt.figure(figsize=(250, 20))
7 plt.plot(pandas_df["Date"], pandas_df["count"])
8 plt.xlabel("Date")
9 plt.ylabel("Number of Baskets")
10 plt.title("Number of Baskets Purchased per Day")
11 plt.xticks(rotation=30)
12 plt.show()



1 class PCY:
2     def __init__(self, basket_file, s=0.3, c=0.5):
3         self.basket_file = basket_file
4         self.s = s
5         self.c = c
6         self.spark = SparkSession.builder.appName("PCY").getOrCreate()
7         self.basket = self.create_basket()
8         self.hash_table_size = 10 # Kích thước bảng băm
9         self.hash_table = [0] * self.hash_table_size # vector có giá trị 0 hoặc 1
10        self.frequency_df = self.frequency()
11
12    def create_basket(self):
13        # Đọc file và tạo baskets
14        df = self.spark.read.csv(self.basket_file).withColumn("STT", F.monotonically_increasing_id())
15        df = df.withColumn("Basket", F.split(F.col("_c0"), ", "))
16        df = df.select("STT", "Basket")
17        return df
18
19    def frequency(self):
20        # Tìm danh sách frequency
21        windowSpec = Window.orderBy("item")
22        items_df = self.basket.select("Basket").withColumn("item", F.explode("Basket")).groupBy("item").count()
23        items_df = items_df.withColumn("id", F.row_number().over(windowSpec))
24        items_df = items_df.select("id", "item", "count")
25        items_df = items_df.filter(F.col("count") > 1).orderBy(F.lower(F.col("item")))
26        return items_df
27
28    def hash_function(self, pair, item_ids):
29        # Tính giá trị hash cho từng pair
30        sum_hash = item_ids.get(pair[0], 0) + item_ids.get(pair[1], 0)
31        return sum_hash % self.hash_table_size
32
33    def create_hash_df(self, pairs, item_ids):
34        # Xử lý bảng băm
35        hash_values = [self.hash_function(pair, item_ids) for pair in pairs]
36        hash_values_df = self.spark.createDataFrame(zip(pairs, hash_values), ["items", "id_bucket"])
37        hash_values_df = hash_values_df.selectExpr("items_1 as item1", "items_2 as item2", "id_bucket")
38        hash_table_df = hash_values_df.groupBy("id_bucket").count() \
39        .filter(F.col("count") > (self.s * self.basket.count())) \
40        .orderBy(F.col("count").asc())
41
42    # Cập nhật bảng băm
43    for row in hash_table_df.collect():
44        id_bucket = row["id_bucket"]
45        self.hash_table[id_bucket] = 1
46
47    # Lọc dựa các bucket thoả mãn
48    hash_values_df = hash_values_df.filter(F.col("id_bucket") \
49    .isin([id_bucket for id_bucket, value in enumerate(self.hash_table) if value == 1]))
50    return hash_values_df
51
52    def find_frequent_pairs(self):
53        frequency_df = self.frequency()
54        exploded_df = self.basket.select("STT", F.explode("Basket").alias("item"))
55
56        # Loại bỏ các item không nằm trong frequency_df
57        exploded_df_alias = exploded_df.alias("edf")
58        modified_df = exploded_df_alias.join(
59            frequency_df,
60            exploded_df_alias["item"] == frequency_df["item"],
61            "inner"
62        ).groupBy("STT").agg(F.expr("collect_list(edf.item) as Basket"))
63
64        # Tạo các cặp pair dựa trên baskets đã cập nhật
65        pairs_with_freq = modified_df.rdd.flatMap(lambda row:
66            [(item1, item2)
67             for item1, item2 in combinations(sorted(row["Basket"]), 2)].collect())
68
69        # Hash table (Xử lý bảng băm)
70        item_ids = self.frequency_df.select("item", "id").rdd.collectAsMap()
71        df_hash_table = self.create_hash_df(pairs_with_freq, item_ids)
72
73        # Count all pairs (Tìm frequent tất cả các pair thoả mãn điều kiện)
74        df_hash_table = df_hash_table.groupBy("item1", "item2", "id_bucket").count()
75        df_hash_table = df_hash_table.withColumnRenamed("count", "frequency")
76        filtered_pairs_df = df_hash_table.filter(F.col("frequency") > (self.s * self.basket.count()))
77
78        joined_df = filtered_pairs_df.join(frequency_df, filtered_pairs_df["item1"] == frequency_df["item"], "inner") \
79        .withColumnRenamed("count", "fr1") \
80        .withColumnRenamed("id", "id1") \
81        .drop("item")
82
83        joined_df = joined_df.join(frequency_df, joined_df["item2"] == frequency_df["item"], "inner") \
84        .withColumnRenamed("count", "fr2") \
```

```
84         .withColumnRenamed('item1', 'fr1') \
85         .withColumnRenamed("id", "id2") \
86         .select("id1", "item1", "id2", "item2","frequency", "fr1", "fr2", "id_bucket" ) \
87         .orderBy(F.lower(F.col("item1")))
88     return joined_df
89
90 def find_association_rules(self, pairs_df):
91     # Xứ lý và tìm association rules
92     total_baskets = self.basket.count()
93
94     association_rules1_df = pairs_df.withColumn("confidence", F.col("frequency") / F.col("fr1")) \
95                                     .withColumn("support", F.col("frequency") / total_baskets) \
96                                     .withColumn("lift", F.col("confidence") / F.col("fr2") * total_baskets) \
97                                     .withColumnRenamed('item1', 'antecedent').withColumnRenamed('item2', 'consequent') \
98                                     .select("antecedent", "consequent", "confidence", "lift", "support") \
99
100     association_rules2_df = pairs_df.withColumn("confidence", F.col("frequency") / F.col("fr2")) \
101                                     .withColumn("support", F.col("frequency") / total_baskets) \
102                                     .withColumn("lift", F.col("confidence") / F.col("fr1") * total_baskets) \
103                                     .withColumnRenamed('item2', 'antecedent').withColumnRenamed('item1', 'consequent') \
104                                     .select("antecedent", "consequent", "confidence", "lift", "support")
105
106     association_rules_df = association_rules1_df.union(association_rules2_df) \
107                                     .filter( (F.col("confidence")>= self.c) \
108                                             & (F.col("support") >= self.s)) \
109                                     .orderBy(F.lower(F.col("antecedent")))
110     return association_rules_df
111
112 def run(self):
113     # Find frequent pairs
114     frequent_pairs = self.find_frequent_pairs()
115     frequent_pairs.show(truncate=False)
116
117     # Find association rules
118     association_rules_df = self.find_association_rules(frequent_pairs)
119     association_rules_df.show(truncate=False)
120
121     frequent_pairs.toPandas().to_csv('pcy_frequent_pairs.csv')
122     association_rules_df.toPandas().to_csv('pcy_association_rules.csv')
123
124 def close(self):
125     self.spark.stop()
126
127 spark = SparkSession.builder.getOrCreate()
128 pcy_instance = PCY('/content/baskets/*.csv',0.012,0.01)
129 pcy_instance.run()
130 pcy_instance.close()
```

id1	item1	id2	item2	frequency	fr1	fr2	id_bucket
104	other vegetables	167	whole milk	222	1827	2363	1
125	rolls/buns	167	whole milk	209	1646	2363	2

antecedent	consequent	confidence	lift	support
other vegetables	whole milk	0.12151067323481117	0.7694304712706219	0.014836596939116486
rolls/buns	whole milk	0.12697448359659783	0.8040284376030018	0.013967787208447505
whole milk	other vegetables	0.09394837071519255	0.7694304712706219	0.014836596939116486
whole milk	rolls/buns	0.08844688954718578	0.8040284376030018	0.013967787208447505