```
permonth and per year query:
---------------------------
SELECT date_format(timestamp_column, 'yyyy-MM') as month_year, count(*) as count
FROM table_name
GROUP BY month_year
--------------------------------------------------------
From pyspark.sql.types import *

The pyspark.sql.types module in PySpark provides a collection of data types that can be
used to define the schema of a DataFrame. Here's a brief explanation of some of the data
types available in this module:

StringType: Represents string values
IntegerType: Represents integer values
LongType: Represents long integer values
DoubleType: Represents double precision floating-point values
FloatType: Represents single precision floating-point values
DecimalType: Represents decimal values with fixed precision and scale
TimestampType: Represents timestamp values
DateType: Represents date values
BooleanType: Represents boolean values
You can use these data types to define the schema of a DataFrame using the StructType
and StructField classes.

For example:

ETL=>Ingestion transfermatuon load reporting:

you can create a DataFrame using the createDataFrame method.

cust_data='dbfs:/filestor/data/cust.csv',
emp_data='dbfs:/filestore/data/emp.csv',
dept_data='dbfs:/filestore/data/dept.csv'

=>cust_schema=StructType([StructField('cust_no',IntegerType(),Nullabale=true)])

schema = StructType([
    StructField("name", StringType(), True),
    StructField("age", IntegerType(), True),
    StructField("salary", DoubleType(), True)
])
----------------------------------------------------------------------------------
--------------------
RDD:
----
from pysark.sql.tyoes import *
from pyspark.sql import sparksession

schema = StructType([
    StructField("name", StringType(), True),
    StructField("age", IntegerType(), True),
    StructField("salary", DoubleType(), True)
])


from pyspark.sql import SparkSession

spark = SparkSession.builder.appName("Example").getOrCreate()

data = [("Alice", 25, 50000.0), ("Bob", 30, 60000.0), ("Charlie", 35, 70000.0)]
rdd = spark.sparkContext.parallelize(data)

df = spark.createDataFrame(rdd, schema)
df.show()

-------------------------------------------------------------------------------------
 from pyspark.sql.functions import *
statement imports all the functions available in the pyspark.sql.functions module in
PySpark.
```

```
65      This allows you to use these functions in your PySpark SQL queries without having to
        prefix them with f..

66
67    Here are some common functions that can be used with this import statement:

68
69    col: returns a Column based on the given column name
70    lit: creates a Column with a literal value
71    when: evaluates a condition and returns one of two values based on the result
72    concat: concatenates two or more columns into a single column
73    substring: extracts a substring from a column based on a starting index and length
74    upper: converts a string column to uppercase
75    lower: converts a string column to lowercase
76    trim: removes leading and trailing whitespace from a string column
77    round: rounds a numeric column to a specified number of decimal places
78    count: counts the number of rows in a DataFrame or the number of non-null values in a
        specific column
79    sum: computes the sum of values in a numeric column
80    avg: computes the average of values in a numeric column
81    min: finds the minimum value in a numeric or string column
82    max: finds the maximum value in a numeric or string column

83
84    By importing all the functions using from pyspark.sql.functions import *, you can use
      these and many other functions in your PySpark SQL queries. However, it's important to
      note that importing all functions can lead to naming conflicts if you have functions
      with the same name as those in the module. In such cases, it may be better to import
      functions individually to avoid any potential issues.

85
86    from pyspark.sql.functions import *
87    each examples
88    Sure, here are some examples of how to use some of the functions available in
      pyspark.sql.functions:

89
90    col: Returns a Column based on the given column name.
91    sql
92    from pyspark.sql.functions import col

93
94    # Create a DataFrame with a column named "name"
95    data = [("Alice", 25), ("Bob", 30), ("Charlie", 35)]
96    df = spark.createDataFrame(data, ["name", "age"])

97
98    # Use "col" to reference the "name" column in a PySpark SQL query
99    df.select(col("name")).show()
100   --------------------------------------------------------------------------------
101   when: Evaluates a condition and returns one of two values based on the result.

102
103   from pyspark.sql.functions import when

104
105   # Create a DataFrame with columns named "age" and "is_adult"

106
107   data = [(25, True), (17, False), (30, True)]
108   df = spark.createDataFrame(data, ["age", "is_adult"])

109
110   # Use "when" to create a new column named "adult_status"
111   # If "is_adult" is True, set "adult_status" to "Adult"
112   # If "is_adult" is False, set "adult_status" to "Minor"

113
114   df.select(
115       "age",
116       "is_adult",
117       when(col("is_adult"), "Adult").otherwise("Minor").alias("adult_status")
118   ).show()
119   ---------------------------------------------------------------------------------------
      ------

120
121   concat: Concatenates two or more columns into a single column.
122   sql

123
124   from pyspark.sql.functions import concat
125
```

```
126    # Create a DataFrame with columns named "first_name" and "last_name"
127
128    data = [("Alice", "Smith"), ("Bob", "Johnson"), ("Charlie", "Brown")]
129    df = spark.createDataFrame(data, ["first_name", "last_name"])
130
131    # Use "concat" to create a new column named "full_name"
132    # Concatenate the "first_name" and "last_name" columns with a space in between
133
134    df.select(concat(col("first_name"), lit(" "),
       col("last_name")).alias("full_name")).show()
135    ------------------------------------------------------------------------------------
       ------
136    substring: Extracts a substring from a column based on a starting index and length.
137    sql
138
139    from pyspark.sql.functions import substring
140
141    # Create a DataFrame with a column named "name"
142
143    data = [("Alice Smith"), ("Bob Johnson"), ("Charlie Brown")]
144    df = spark.createDataFrame(data, ["name"])
145
146    # Use "substring" to create a new column named "last_name"
147    # Extract the last name from the "name" column, assuming it is separated by a space
148
149    df.select(substring(col("name"), instr(col("name"), " ") + 1,
       length(col("name"))).alias("last_name")).show()
150    ------------------------------------------------------------------------------------
       ----------------------------
151    ETL-CONCEPT:
152    ------------
153    #ETL--concepts etc
154
155    from pyspark.sql.types import *
156    from pyspark.sql.functions import col
157
158    data_shopping_cust = 'dbfs:/FileStore/Shopping_CustomerData.csv'
159    data_shopping_index = 'dbfs:/FileStore/Shopping_ShoppingIndexData.csv'
160
161    cust_schema = StructType([StructField('CustomerAge', IntegerType(), nullable=True)])
162    index_schema = StructType([StructField('Kolkata', IntegerType(), nullable=True)])
163
164    df = spark.read.option('header', 'true').csv(path=data_shopping_cust,
       schema=cust_schema).show()
165    df1 = spark.read.option('header', 'true').csv(path=data_shopping_index,
       schema=index_schema).show()
166    ------------------------------------------------------------------------------------
       -----------------------
167    partitions:
168    -----------
169    from pyspark.sql.functions import year
170
171    df = spark.read.format('csv').option('header',
       'true').load('dbfs:/FileStore/Shopping_CustomerData.csv')
172    df = df.withColumn('year', year('AnnualIncome',))
173    df.write.partitionBy('year').mode('overwrite').parquet('/dbfs:/FileStore/output.csv')
174    ------------------------------------------------------------------------------------
       ----------------------
175    Using hash:
176    ----------
177    df_partion = spark.read.format('csv').option('header',
       'true').load('dbfs:/FileStore/Shopping_CustomerData.csv')
178    df_partion  = df.repartition(4)
179    df_partion .write.mode('overwrite').parquet('File//path/to/output')
180    df_partion .show()
181    ----------------------------------------------------------------
182    partition:
183    ----------
184    from pyspark.sql.types import *
```

```
185
186    # Write the DataFrame to disk, partitioned by "Channel_Name" and "Genre"
187    df.write.mode("overwrite").partitionBy("Channel_Name",
       "Genre").csv("dbfs:/FileStore/output")
188
189    # Read the data back from disk
190    df2 = spark.read.csv("dbfs:/FileStore/output", header=True)
191
192    # Show the data
193    df2.show()
194    ----------------------------------------------------------------------
195    # Print the current number of partitions
196    print('Number of current partitions:', str(df.rdd.getNumPartitions()))
197
198    # Reduce the number of partitions to 5 using coalesce
199    df_reduced = df.coalesce(5)
200    print('Number of partitions after reducing using coalesce:',
       str(df_reduced.rdd.getNumPartitions()))
201
202    # Increase the number of partitions to 10 using repartition
203    df_increased = df.repartition(10)
204    print('Number of partitions after increasing using repartition:',
       str(df_increased.rdd.getNumPartitions()))
205    -----------------------------------------------------------------------------------------
       --------------------
206    Broad Cast Join:
207    ----------------
208
209    from pyspark.sql.functions import broadcast
210
211    large_df = spark.read.parquet("path/to/large/dataframe")
212    small_df = spark.read.parquet("path/to/small/dataframe")
213
214    # Mark the small_df for broadcast join
215    small_df_b = broadcast(small_df)
216
217    # Perform the join operation
218    joined_df = large_df.join(small_df_b, "join_column")
219    -----------------------------------------------------------------------------------------
       ----------------
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
```