```
from pyspark.sql.functions import col
from pyspark.sql.types import IntegerType, DoubleType, BooleanType, DateType

Accountkey = "9nsX5Nm3u+bDVap4WVhERJ3nJHN9sQGOaTb9eodL8AKDMzzDyw2siAOqB+lJlSpbD/2yVsMy8CLA+AStZXJDXQ==" # Rep

dbutils.fs.mount(
    source= "wasbs://tokyo-olympics-dataset@tokyoolympicsdataset01.blob.core.windows.net", # Remove trailing
    mount_point = "/mnt/tokyo-olympics-dataset",
    extra_configs = {"fs.azure.account.key.tokyoolympicsdataset01.blob.core.windows.net": Accountkey}
)

Dut[4]: True

display(dbutils.fs.mounts())
```



mountPoint	source	encryptionType		
/databricks-datasets	databricks-datasets			
/Volumes	Volumes UnityCatalogVolumes			
/databricks/mlflow-tracking	databricks/mlflow-tracking			
/databricks-results	databricks-results databricks-results			
/databricks/mlflow-registry	databricks/mlflow-registry			
/Volume	DbfsReserved			
/mnt/tokyo-olympics-dataset	wasbs://tokyo-olympics-dataset@tokyoolympicsdataset01.blob.core.windows.net			
/volumes	DbfsReserved			
1	DatabricksRoot			
Wolume	NhfsReserved			

%fs

ls "/mnt/tokyo-olympics-dataset"



path	name	size	modificationTime
dbfs:/mnt/tokyo-olympics-dataset/Raw_Data/	Raw_Data/	0	0
dbfs:/mnt/tokyo-olympics-dataset/Transformed_Data/	Transformed_Data/	0	0

athletes = spark.read.format("csv").option("header","true").option("inferSchema","true").load("dbfs:/mnt/tokyo-olympics-dataset/Raw_coaches = spark.read.format("csv").option("header","true").option("inferSchema","true").load("dbfs:/mnt/tokyo-olympics-dataset/Raw_centriesgender = spark.read.format("csv").option("header","true").option("inferSchema","true").load("dbfs:/mnt/tokyo-olympics-dataset/Raw_Datems = spark.read.format("csv").option("header","true").option("inferSchema","true").option("dbfs:/mnt/tokyo-olympics-dataset/Raw_Datems = spark.read.format("csv").option("header","true").option("header","true").option("header","true").option("header","true").option("header","true").option("header","true").option("header","true").option("hea

athletes.show()



+		
PersonName	Country	Discipline
AALERUD Katrine	Norway	Cycling Road
ABAD Nestor	Spain	Artistic Gymnastics
ABAGNALE Giovanni	Italy	Rowing
ABALDE Alberto	Spain	Basketball
ABALDE Tamara	Spain	Basketball
ABALO Luc	France	Handball
ABAROA Cesar	Chile	Rowing
ABASS Abobakr	Sudan	Swimming
ABBASALI Hamideh	Islamic Republic	Karate
ABBASOV Islam	Azerbaijan	Wrestling
ABBINGH Lois	Netherlands	Handball
ABBOT Emily	Australia	Rhythmic Gymnastics
ABBOTT Monica	United States of	Baseball/Softball
ABDALLA Abubaker	Qatar	Athletics
ABDALLA Maryam	Egypt	Artistic Swimming
ABDALLAH Shahd	Egypt	Artistic Swimming
ABDALRASOOL Mohamed	Sudan	Judo

ABDEL LATIF Radwa	Egypt	Shooting
ABDEL RAZEK Samy	Egypt	Shooting
ABDELAZIZ Abdalla	Egypt	Karate
+		+
only showing top 20 rows		

athletes.printSchema()

\rightarrow root

|-- PersonName: string (nullable = true)
|-- Country: string (nullable = true)
|-- Discipline: string (nullable = true)

coaches.show()



L	L	L	L	_	
	Event	Discipline	Country	Name	
Ī	null	Football	Egypt	ABDELMAGID Wael	•
	null	Volleyball	Japan	ABE Junya	
	null	Basketball	Japan	ABE Katsuhiko	
-	[nul]	Footbal	C�te d'Ivoire	ADAMA Cherif	
ĺ	null	Volleyball	Japan	AGEBA Yuya	
	Men	Hockey	Japan	AIKMAN Siegfried	
	Men	Hockey	Germany	AL SAADI Kais	
	Softball	Baseball/Softball	Canada	ALAMEDA Lonni	
	Men	Volleyball	Islamic Republic	ALEKNO Vladimir	
	Women	Handball	ROC	ALEKSEEV Alexey	
	null	Basketball	Spain	ALLER CARBALLO Ma	
	Men	Football	Saudi Arabia	ALSHEHRI Saad	
	null	Football	Egypt	ALY Kamal	
	null	Basketball	Puerto Rico	AMAYA GAITAN Fabian	
	null	Football	Spain	AMO AGUADO Pablo	
	Women	Football	United States of	ANDONOVSKI Vlatko	
	Women	Hockey	Netherlands	ANNAN Alyson	
	Women	Hockey	Japan	ARNAU CREUS Xavier	

```
ARNOLD Graham | Australia | Football | Men | AXNER Tomas | Sweden | Handball | Women | Ha
```

coaches.printSchema()

→ root

|-- Name: string (nullable = true)

|-- Country: string (nullable = true)

|-- Discipline: string (nullable = true)

|-- Event: string (nullable = true)

entriesgender.show()

\rightarrow	+		-	+					
	Discipline Female Male Total								
	+		+	+					
	3x3 Basketball	32	32	64					
	Archery	64	64	128					
	Artistic Gymnastics	98	98	196					
	Artistic Swimming	105	0	105					
	Athletics	969	1072	2041					
	Badminton	86	87	173					
	Baseball/Softball	90	144	234					
	Basketball	144	144	288					
	Beach Volleyball	48	48	96					
	Boxing	102	187	289					
	Canoe Slalom	41	41	82					
	Canoe Sprint	123	126	249					
	Cycling BMX Frees	10	9	19					
	Cycling BMX Racing	24	24	48					
	Cycling Mountain	38	38	76					
	Cycling Road	70	131	201					
	Cycling Track	90	99	189					
	Diving	72	71	143					
	Equestrian	73	125	198					

```
Fencing|
                            107 | 108 | 215 |
     only showing top 20 rows
entriesgender.printSchema()
     root
      |-- Discipline: string (nullable = true)
      |-- Female: integer (nullable = true)
      |-- Male: integer (nullable = true)
      |-- Total: integer (nullable = true)
entriesgender = entriesgender.withColumn("Female",col("Female").cast(IntegerType()))\
    .withColumn("Male",col("Male").cast(IntegerType()))\
    .withColumn("Total",col("Total").cast(IntegerType()))
entriesgender.printSchema()
\rightarrow
    root
      |-- Discipline: string (nullable = true)
      |-- Female: integer (nullable = true)
      |-- Male: integer (nullable = true)
      |-- Total: integer (nullable = true)
medals.show()
     +---+
                  TeamCountry|Gold|Silver|Bronze|Total|Rank by Total|
     |Rank|
        1|United States of ...| 39|
                                             33 | 113 |
        2|People's Republic...| 38| 32|
                                                                 2
                                             18 88
                                      14
                                                                 5 |
                        Japan | 27 |
                                             17|
                                                   58
```

22

21

22

65

Great Britain

5	ROC	20	28	23	71	3
6	Australia	17	7	22	46	6
7	Netherlands	10	12	14	36	9
8	France	10	12	11	33	10
9	Germany	10	11	16	37	8
10	Italy	10	10	20	40	7
11	Canada	7	6	11	24	11
12	Brazil	7	6	8	21	12
13	New Zealand	7	6	7	20	13
14	Cuba	7	3	5	15	18
15	Hungary	6	7	7	20	13
16	Republic of Korea	6	4	10	20	13
17	Poland	4	5	5	14	19
18	Czech Republic	4	4	3	11	23
19	Kenya	4	4	2	10	25
20	Norway	4	2	2	8	29

only showing top 20 rows

medals.printSchema()

```
→ root
```

|-- Rank: integer (nullable = true)

|-- TeamCountry: string (nullable = true)

|-- Gold: integer (nullable = true)

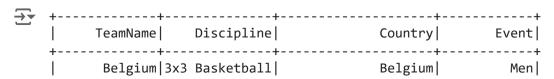
|-- Silver: integer (nullable = true)

|-- Bronze: integer (nullable = true)

|-- Total: integer (nullable = true)

|-- Rank by Total: integer (nullable = true)

teams.show()



```
China|3x3 Basketball|People's Republic...|
                                                          Men
        China|3x3 Basketball|People's Republic...|
                                                        Women
       France | 3x3 Basketball |
                                          Francel
                                                        Women
        Italv|3x3 Basketball|
                                           Italv
                                                        Women
        Japan | 3x3 Basketball |
                                                          Men
                                           Japan
        Japan 3x3 Basketball
                                           Japan
                                                        Womenl
       Latvia | 3x3 Basketball |
                                          Latvia
                                                          Men
                                        Mongolia
     Mongolia | 3x3 Basketball |
                                                        Women
  Netherlands | 3x3 Basketball |
                                     Netherlands
                                                          Men
       Poland | 3x3 Basketball |
                                          Poland|
                                                          Men
          ROC|3x3 Basketball|
                                             ROC |
                                                          Men
          ROC|3x3 Basketball|
                                             ROC I
                                                        Women
      Romania | 3x3 Basketball |
                                         Romania
                                                        Women
       Serbia 3x3 Basketball
                                          Serbial
                                                          Menl
|United States|3x3 Basketball|United States of ...|
                                                        Women
    Australia|
                     Archery
                                       Australia
                                                   Men's Team
    Australial
                    Archery
                                     Australia
                                                   Mixed Team
   Bangladesh|
                                      Bangladesh|
                    Archery
                                                   Mixed Team
      Belarus|
                    Archery
                                         Belarus | Women's Team |
+-----
```

only showing top 20 rows

```
teams.printSchema()
```

```
\rightarrow
    root
      |-- TeamName: string (nullable = true)
      |-- Discipline: string (nullable = true)
      |-- Country: string (nullable = true)
```

|-- Event: string (nullable = true)

Find the top countries with the highest number of gold medals top gold medal countries = medals.orderBy("Gold", ascending=False).select("TeamCountry", "Gold").show()

```
TeamCountry|Gold|
|United States of ...| 39|
```

```
|People's Republic...|
                       38
                       27
               Japan
       Great Britain
                       22
                       20
                 ROC |
           Australia|
                       17|
                       10
         Netherlands
              France
                       10
                       10
             Germany
               Italv
                       10
                        7
              Canada
              Brazil
                        7|
                        7
         New Zealand
                        7
                Cubal
                        6
             Hungary|
   Republic of Korea
                        6
              Poland|
                        4
      Czech Republic
                        4
               Kenya
                        4
              Norway
                        4
only showing top 20 rows
```

```
# Calculate the average number of entries by gender for each discipline
average_entries_by_gender = entriesgender.withColumn(
    'Avg_Female', entriesgender['Female'] / entriesgender['Total']
).withColumn(
    'Avg_Male', entriesgender['Male'] / entriesgender['Total']
)
average_entries_by_gender.show()
```

\rightarrow	+				+		
<u></u>	<u>į</u>	 Discipline	Female	Male T	otal	Avg_Female	Avg_Male
	3x3	Basketball	32	32	64	0.5	0.5
		Archery	64	64	128	0.5	0.5
	Artistic	Gymnastics	98	98	196	0.5	0.5
	Artist:	ic Swimming	105	0	105	1.0	0.0

Athletics	969	1072	2041	0.4747672709456149	0.5252327290543851
Badminton	86	87	173	0.49710982658959535	0.5028901734104047
Baseball/Softball	90	144	234	0.38461538461538464	0.6153846153846154
Basketball	144	144	288	0.5	0.5
Beach Volleyball	48	48	96	0.5	0.5
Boxing	102	187	289	0.35294117647058826	0.6470588235294118
Canoe Slalom	41	41	82	0.5	0.5
Canoe Sprint	123	126	249	0.4939759036144578	0.5060240963855421
Cycling BMX Frees	10	9	19	0.5263157894736842	0.47368421052631576
Cycling BMX Racing	24	24	48	0.5	0.5
Cycling Mountain	38	38	76	0.5	0.5
Cycling Road	70	131	201	0.3482587064676617	0.6517412935323383
Cycling Track	90	99	189	0.47619047619047616	0.5238095238095238
Diving	72	71	143	0.5034965034965035	0.4965034965034965
Equestrian	73	125	198	0.3686868686868687	0.6313131313131313
Fencing	107	108	215	0.49767441860465117	0.5023255813953489
+		+		+	

only showing top 20 rows

athletes.repartition(1).write.mode("overwrite").option("header", 'true').csv("dbfs:/mnt/tokyo-olympics-dataset/Transformed_Data/Athle coaches.repartition(1).write.mode("overwrite").option("header", "true").csv("dbfs:/mnt/tokyo-olympics-dataset/Transformed_Data/Coache entriesgender.repartition(1).write.mode("overwrite").option("header", "true").csv("dbfs:/mnt/tokyo-olympics-dataset/Transformed_Data/Medals.teams.repartition(1).write.mode("overwrite").option("header", "true").csv("dbfs:/mnt/tokyo-olympics-dataset/Transformed_Data/Medals.teams.repartition(1).write.mode("overwrite").option("header", "true").csv("dbfs:/mnt/tokyo-olympics-dataset/Transformed_Data/Teams.cs