```
1 #configuración en google colab de spark y pyspark
2 from google.colab import drive
3 drive.mount('/content/gdrive')

→ Mounted at /content/gdrive

1 #instalar java y spark
2 !apt-get install openjdk-11-jdk-headless -qq > /dev/null
3 !wget -q https://dlcdn.apache.org/spark/spark-3.5.3/spark-3.5.3-bin-hadoop3.tgz
4 !tar xf spark-3.5.3-bin-hadoop3.tgz
5 !pip install -q findspark
1 import os
2 os.environ["JAVA_HOME"] = "/usr/lib/jvm/java-11-openjdk-amd64"
3 os.environ["SPARK_HOME"] = "/content/spark-3.5.3-bin-hadoop3"
1 import findspark
2 findspark.init()
3 from pyspark.sql import SparkSession
4 spark = SparkSession.builder.master("local[*]").getOrCreate()
5 sc = spark.sparkContext
1 spark
SparkSession - in-memory
    SparkContext
    Spark UI
    Version
          v3.5.3
    Master
          local[*]
    AppName
          pyspark-shell
1 sc
SparkContext
    Spark UI
    Version
          v3.5.3
    Master
          local[*]
    AppName
          pyspark-shell
    4
1 # Load csv Dataset
2 df=spark.read.csv('gdrive/MyDrive/st0263-242/bigdata/datasets/covid19/Casos_positivos_de_COVID-19_en_Colombia-100K.csv',inferSchema=1
1 df.columns

    ['fecha reporte web',
      'ID de caso',
      'Fecha de notificación',
      'Código DIVIPOLA departamento',
      'Nombre departamento',
      'Código DIVIPOLA municipio',
      'Nombre municipio',
      'Edad',
      'Unidad de medida de edad',
      'Sexo',
      'Tipo de contagio',
      'Ubicación del caso',
      'Estado',
      'Código ISO del país',
      'Nombre del país',
      'Recuperado',
      'Fecha de inicio de síntomas',
      'Fecha de muerte',
      'Fecha de diagnóstico',
      'Fecha de recuperación',
      'Tipo de recuperación',
      'Pertenencia étnica',
      'Nombre del grupo étnico']
```

```
1 df.printSchema()
→ root
     |-- fecha reporte web: string (nullable = true)
      |-- ID de caso: integer (nullable = true)
      |-- Fecha de notificación: string (nullable = true)
      |-- Código DIVIPOLA departamento: integer (nullable = true)
      |-- Nombre departamento: string (nullable = true)
      |-- Código DIVIPOLA municipio: integer (nullable = true)
      |-- Nombre municipio: string (nullable = true)
      |-- Edad: integer (nullable = true)
      -- Unidad de medida de edad: integer (nullable = true)
     |-- Sexo: string (nullable = true)
      |-- Tipo de contagio: string (nullable = true)
      -- Ubicación del caso: string (nullable = true)
      -- Estado: string (nullable = true)
      |-- Código ISO del país: integer (nullable = true)
      |-- Nombre del país: string (nullable = true)
      |-- Recuperado: string (nullable = true)
      |-- Fecha de inicio de síntomas: string (nullable = true)
      |-- Fecha de muerte: string (nullable = true)
      |-- Fecha de diagnóstico: string (nullable = true)
      -- Fecha de recuperación: string (nullable = true)
     |-- Tipo de recuperación: string (nullable = true)
      -- Pertenencia étnica: integer (nullable = true)
      |-- Nombre del grupo étnico: string (nullable = true)
1 df.select('fecha reporte web', 'ID de caso', 'Sexo', 'Tipo de contagio').show(10, False)
    |fecha reporte web|ID de caso|Sexo|Tipo de contagio|
    6/3/2020 0:00:00 1
                              |F |Importado
     9/3/2020 0:00:00 2
                                 IM
                                     Importado
     9/3/2020 0:00:00 3
                                 İΕ
                                      Importado
     11/3/2020 0:00:00 4
                               М
                                     Relacionado
     11/3/2020 0:00:00 5
                                 М
                                      Relacionado
     11/3/2020 0:00:00 6
                                      Relacionado
     11/3/2020 0:00:00 7
                                 F
                                      Importado
    11/3/2020 0:00:00 8
                                      Importado
     11/3/2020 0:00:00 9
                                      Importado
    12/3/2020 0:00:00 10
                                F
                                     Importado
    only showing top 10 rows
1 df=df.withColumnRenamed('fecha reporte web', 'fecha_reporte')
1 df=df.withColumnsRenamed({'ID de caso': 'ID', 'Nombre del país': 'pais', 'Nombre departamento': 'departamento', 'Nombre municipio':
1 df=df.withColumn('Edad_mas_10', df['Edad'] + 10)
2 df.show(5, False)
    |fecha_reporte | ID | Fecha de notificación|Código DIVIPOLA departamento|departamento|Código DIVIPOLA municipio|municipio|Edad|Unic
    6/3/2020 0:00:00 |1 |2/3/2020 0:00:00
                                                                            BOGOTA
                                                                                                                   BOGOTA 19
     9/3/2020 0:00:00 2
                          6/3/2020 0:00:00
                                                76
                                                                             VALLE
                                                                                          76111
                                                                                                                   BUGA
                                                                                                                             34
                                                                                                                   MEDELLIN |50
    9/3/2020 0:00:00 | 3 | 7/3/2020 0:00:00
                                                5
                                                                             IANTIOOUIA
                                                                                         5001
                                                                                                                                  1
     11/3/2020 0:00:00|4
                          9/3/2020 0:00:00
                                                                             ANTIOQUIA
                                                                                         5001
                                                                                                                   MEDELLIN 55
    11/3/2020 0:00:00 5 9/3/2020 0:00:00
                                                                            ANTIOQUIA
                                                                                         5001
                                                                                                                   MEDELLIN 25 1
    only showing top 5 rows
    4
Double-click (or enter) to edit
1 columns_to_remove=['Nombre del grupo étnico', 'Pertenencia étnica']
2 df=df.drop(*columns_to_remove)
1 df.printSchema()
→ root
     |-- fecha_reporte: string (nullable = true)
     -- ID: integer (nullable = true)
      -- Fecha de notificación: string (nullable = true)
      |-- Código DIVIPOLA departamento: integer (nullable = true)
      |-- departamento: string (nullable = true)
      |-- Código DIVIPOLA municipio: integer (nullable = true)
      |-- municipio: string (nullable = true)
      |-- Edad: integer (nullable = true)
```

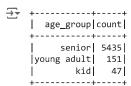
```
|-- Unidad de medida de edad: integer (nullable = true)
     -- Sexo: string (nullable = true)
     |-- Tipo de contagio: string (nullable = true)
     -- Ubicación: string (nullable = true)
     -- Estado: string (nullable = true)
     |-- Código ISO del país: integer (nullable = true)
     |-- pais: string (nullable = true)
     |-- Recuperado: string (nullable = true)
     |-- Fecha de inicio de síntomas: string (nullable = true)
     -- Fecha de muerte: string (nullable = true)
     -- Fecha de diagnóstico: string (nullable = true)
     |-- Fecha de recuperación: string (nullable = true)
     -- Tipo de recuperación: string (nullable = true)
     |-- age_group: string (nullable = true)
     -- Edad_mas_10: integer (nullable = true)
1 filtered_df=df.filter(df.Sexo == 'F')
2 filtered_df.show(5, False)
    |fecha_reporte | ID |Fecha de notificación|Código DIVIPOLA departamento|departamento|Código DIVIPOLA municipio|municipio|Edad|Unic
    6/3/2020 0:00:00 1
                         2/3/2020 0:00:00
                                                11
                                                                             BOGOTA
                                                                                          11001
                                                                                                                    BOGOTA
                                                                                                                              119
    9/3/2020 0:00:00 3
                         17/3/2020 0:00:00
                                                                             IANTIOOUIA
                                                                                                                    MEDELLIN | 50
                                                                                          5001
                                                15
                                                                                                                                   11
                         10/3/2020 0:00:00
    11/3/2020 0:00:00|6
                                                                             IANTIOOUIA
                                                                                                                    ITAGUI
                                                                                                                              27
                                                15
                                                                                          5360
                                                                                                                                   11
    11/3/2020 0:00:00 7
                          8/3/2020 0:00:00
                                                13001
                                                                             CARTAGENA
                                                                                          13001
                                                                                                                    CARTAGENA 85
                                                                                                                                   11
    |11/3/2020 0:00:00|8 |9/3/2020 0:00:00
                                                11
                                                                             IBOGOTA
                                                                                          11001
                                                                                                                    BOGOTA
                                                                                                                             22
                                                                                                                                   1
    only showing top 5 rows
1 filtered2 df=df.filter("Sexo=='F' and Edad>20 and Recuperado like '%Recuperado%' and municipio like 'MEDELLIN'")
2 filtered2 df.show()
        fecha reporte | ID|Fecha de notificación | Código DIVIPOLA departamento | departamento | Código DIVIPOLA municipio | municipio | Edad | Unic
    9/3/2020 0:00:00 3
                               7/3/2020 0:00:00
                                                                                ANTIOQUIA
    14/3/2020 0:00:00 20
                              11/3/2020 0:00:00
                                                                            5 İ
                                                                                ANTIOQUIA
                                                                                                                5001 MEDELLIN
                                                                                                                                 26 İ
    19/3/2020 0:00:00|108|
                              17/3/2020 0:00:00
                                                                            5 l
                                                                                ANTIOOUIA
                                                                                                                5001 MEDELLIN
                                                                                                                                 57
    20/3/2020 0:00:00|131
                                                                            5 l
                                                                                                                5001
                                                                                                                     MEDELLIN
                              15/3/2020 0:00:00
                                                                                ANTIOOUIA
                                                                                                                                 22 l
                                                                            5 l
    20/3/2020 0:00:00 135
                              17/3/2020 0:00:00
                                                                                ANTIOOUTA
                                                                                                                5001 MEDELLIN
                                                                                                                                 44
    20/3/2020 0:00:00 141
                               17/3/2020 0:00:00
                                                                            5
                                                                                ANTIOOUIA
                                                                                                                5001
                                                                                                                     MEDELLIN
                                                                                                                                 62
    20/3/2020 0:00:00 142
                               20/3/2020 0:00:00
                                                                            5
                                                                                ANTIOQUIA
                                                                                                                5001 MEDELLIN
                                                                                                                                 35
    |22/3/2020 0:00:00|238|
                               16/3/2020 0:00:00
                                                                            5
                                                                                ANTIOQUIA
                                                                                                                5001
                                                                                                                     MEDELLIN
                                                                                                                                 61
    23/3/2020 0:00:00 268
                               19/3/2020 0:00:00
                                                                            5
                                                                                ANTIOQUIA
                                                                                                                5001 MEDELLIN
                                                                                                                                 37
    23/3/2020 0:00:00 271
                               18/3/2020 0:00:00
                                                                            5
                                                                                ANTIOQUIA
                                                                                                                5001
                                                                                                                     MEDELLIN
                                                                                                                                 21
    23/3/2020 0:00:00 272
                                                                            5
                                                                                ANTIOQUIA
                                                                                                                5001 MEDELLIN
                               18/3/2020 0:00:00
                                                                                                                                 47
    23/3/2020 0:00:00 275
                               23/3/2020 0:00:00
                                                                            5
                                                                                ANTIOQUIA
                                                                                                                5001 MEDELLIN
                                                                                                                                21
    23/3/2020 0:00:00 292
                               19/3/2020 0:00:00
                                                                            5
                                                                                                                5001
                                                                                                                     MEDELLIN
                                                                                ANTIOQUIA
                                                                                                                                 43
                                                                            5 İ
    23/3/2020 0:00:00 293
                               18/3/2020 0:00:00
                                                                                ANTIOOUIA
                                                                                                                5001 MEDELLIN
                                                                                                                                30 l
                                                                            5
                                                                                                                     MEDELLIN
    123/3/2020 0:00:00 294
                               21/3/2020 0:00:00
                                                                                ANTTOOUTA
                                                                                                                5001
                                                                                                                                 301
    123/3/2020 0:00:0012961
                                                                            5 İ
                                                                                                                5001 MEDELLIN
                               19/3/2020 0:00:00
                                                                                ANTIOOUIA
                                                                                                                                 25 l
    25/3/2020 0:00:00 438
                               20/3/2020 0:00:00
                                                                            5
                                                                                ANTIOQUIA
                                                                                                                5001
                                                                                                                     MEDELLIN
                                                                                                                                 36|
    25/3/2020 0:00:00 450
                               19/3/2020 0:00:00
                                                                            5
                                                                                 ANTIOQUIA
                                                                                                                5001 MEDELLIN
                                                                                                                                 24
    25/3/2020 0:00:00 466
                               19/3/2020 0:00:00
                                                                            5 I
                                                                                 ANTIOOUIA
                                                                                                                5001 MEDELLIN
                                                                                                                                 32 l
    [28/3/2020 0:00:00[595]
                               16/3/2020 0:00:00
                                                                            5 I
                                                                                ANTIOOUIA
                                                                                                                5001 MEDELLIN
                                                                                                                                 27
    only showing top 20 rows
1 from pyspark.sql.functions import udf
2 from pyspark.sql.types import StringType,DoubleType,IntegerType
1 age_udf=udf(lambda age: "kid" if age < 18 else ('young adult' if age<=30 else 'senior'), StringType())
2 df=df.withColumn("age_group", age_udf(df.Edad))
3 df.filter(df['age_group']=='young adult').orderBy('Edad', ascending=False).show(10, False)
    |fecha reporte web|ID de caso|Fecha de notificación|Código DIVIPOLA departamento|Nombre departamento|Código DIVIPOLA municipio|Nombr
    19/6/2020 0:00:00 62406
                                 10/6/2020 0:00:00
                                                                                    IBOGOTA
                                                       111
                                                                                                        11001
                                                                                    VALLE
    21/3/2020 0:00:00 192
                                 20/3/2020 0:00:00
                                                       76
                                                                                                        76001
                                                                                                                                  CALI
                                                                                    IBOGOTA
                                                                                                        111001
                                                                                                                                  I BOGOT
    119/6/2020 0:00:00 62398
                                 17/6/2020 0:00:00
                                                       11
    27/3/2020 0:00:00|517
                                                                                                                                  BOGO1
                                 23/3/2020 0:00:00
                                                       111
                                                                                    BOGOTA
                                                                                                        11001
                                                       11
    19/6/2020 0:00:00 62380
                                 10/6/2020 0:00:00
                                                                                    BOGOTA
                                                                                                        11001
                                                                                                                                  I BOGOT
    22/3/2020 0:00:00 213
                                 19/3/2020 0:00:00
                                                       120
                                                                                    CESAR
                                                                                                        20001
                                                                                                                                  VALLE
    19/6/2020 0:00:00 62617
                                 16/6/2020 0:00:00
                                                       11
                                                                                    BOGOTA
                                                                                                        11001
                                                                                                                                  BOG01
    13/3/2020 0:00:00 15
                                 13/3/2020 0:00:00
                                                                                    META
                                                                                                         50001
                                                                                                                                  VILL/
                                                       150
    19/6/2020 0:00:00 62384
                                 7/6/2020 0:00:00
                                                                                    BOGOTA
                                                                                                        11001
                                                                                                                                  BOGO1
    23/3/2020 0:00:00 248
                                 23/3/2020 0:00:00
                                                                                    BOGOTA
                                                                                                        11001
                                                                                                                                  BOGO1
   only showing top 10 rows
```

```
1 df.groupBy('departamento').count().orderBy('count', ascending=False).show(10, False)
    |departamento|count|
     ------
    BOGOTA
                1300161
    |BARRANQUILLA|13065
    ATLANTICO
                10994
    VALLE
                10404
    CARTAGENA
                 8333
    |ANTIOQUIA
                4554
    NARIÑO
                3520
    CUNDINAMARCA 2827
    AMAZONAS
                2317
    Існосо
                1636
   only showing top 10 rows
1 df.groupBy('municipio').count().orderBy('count', ascending=False).show(10, False)
    |municipio |count|
    BOGOTA
                30016
    BARRANQUILLA 13065
    CARTAGENA
                8333
    CALI
                7747
    SOLEDAD
                6233
    LETICIA
                 2194
    MEDELLIN
                2137
                1501
    TUMACO
    BUENAVENTURA 1453
    QUIBDO
           1367
   only showing top 10 rows
1 df.dropna(subset=['Fecha de inicio de síntomas']).groupBy('Fecha de inicio de síntomas').count().orderBy('count', ascending=False).sh
    |Fecha de inicio de síntomas|count|
    10/6/2020 0:00:00
    16/6/2020 0:00:00
                              2558
    18/6/2020 0:00:00
                              2479
                               2452
    112/6/2020 0:00:00
    1/6/2020 0:00:00
                               2429
    18/6/2020 0:00:00
                               2390
    17/6/2020 0:00:00
                               2344
    5/6/2020 0:00:00
                               2266
    9/6/2020 0:00:00
                               2224
    19/6/2020 0:00:00
                              2162
   only showing top 10 rows
1 df.select('edad', 'age_group', 'ID').groupBy('edad', 'age_group').count().orderBy('count', ascending=False).show()
    |edad| age_group|count|
    | 30|young adult| 2735|
      29 young adult 2611
      31
             senior 2569
      28|young adult| 2540|
      27 young adult
                      2494
      26 young adult | 2436 |
      33
              senior 2371
      32
              senior 2362
      25 young adult | 2335
      34
              senior 2310
      35
              senior 2292
      24|young adult| 2214|
      36
              senior 2175
      37
              senior 2132
              senior 2098
      38
      40
              senior 2050
      23|young adult| 2041|
      39
              senior 1985
      22|young adult| 1879|
              senior 1783
```

only showing top 20 rows

```
1 from pyspark.sql.functions import col
```

 $\label{lem:col} 2 \ \text{df.filter(col('Fecha de muerte').isNotNull()).groupBy('age\_group').count().orderBy('count', ascending=False).show()} \\$ 



1 df.filter(col('Fecha de muerte').isNotNull()).show()

Edac	municipio	Código DIVIPOLA municipio	departamento	Código DIVIPOLA departamento	Fecha de notificación	fecha_reporte II
	+	+	+			+
6!	BOGOTA	11001	BOGOTA	11	18/3/2020 0:00:00	20/3/2020 0:00:00 152
53	BOGOTA	11001	BOGOTA	11	18/3/2020 0:00:00	20/3/2020 0:00:00 153
88	SANTA MARTA	47001	STA MARTA D.E.	47001	20/3/2020 0:00:00	20/3/2020 0:00:00 157
76	YUMBO	76892	VALLE	76	17/3/2020 0:00:00	21/3/2020 0:00:00 188
58	CARTAGENA	13001	CARTAGENA	13001	13/3/2020 0:00:00	16/3/2020 0:00:00 197
59	BOGOTA	11001	BOGOTA	11	20/3/2020 0:00:00	22/3/2020 0:00:00 232
7€	BOGOTA	11001	BOGOTA	11	23/3/2020 0:00:00	23/3/2020 0:00:00 256
61	BOGOTA	11001	BOGOTA	11	23/3/2020 0:00:00	23/3/2020 0:00:00 261
46	CALI	76001	VALLE	76	20/3/2020 0:00:00	23/3/2020 0:00:00 282
83	NEIVA	41001	HUILA	41	21/3/2020 0:00:00	23/3/2020 0:00:00 286
62	IPIALES	52356	NARIÑO	52	20/3/2020 0:00:00	24/3/2020 0:00:00 314
61	PEREIRA	66001	RISARALDA	66	21/3/2020 0:00:00	24/3/2020 0:00:00 326
7:	BOGOTA	11001	BOGOTA	11	23/3/2020 0:00:00	24/3/2020 0:00:00 389
62	BOGOTA	11001	BOGOTA	11	21/3/2020 0:00:00	24/3/2020 0:00:00 415
86	TULUA	76834	VALLE	76	21/3/2020 0:00:00	25/3/2020 0:00:00 425
84	BOGOTA	11001	BOGOTA	11	25/3/2020 0:00:00	26/3/2020 0:00:00 491
72	BOGOTA	11001	BOGOTA	11	23/3/2020 0:00:00	27/3/2020 0:00:00 516
71	BOGOTA	11001	BOGOTA	11	22/3/2020 0:00:00	27/3/2020 0:00:00 514
2!	BARRANQUILLA	8001	BARRANQUILLA	8001	24/3/2020 0:00:00	27/3/2020 0:00:00 533
56	SOLEDAD	8758	ATLANTICO	8	24/3/2020 0:00:00	27/3/2020 0:00:00 535

only showing top 20 rows

## Using SparkSQL

1 df.createOrReplaceTempView("covid19")

1 spark.sql("select departamento, count(\*) as count from covid19 group by departamento order by count desc").show(10)

```
| departamento|count|
| BOGOTA|30016|
|BARRANQUILLA|13065|
| ATLANTICO|10994|
| VALLE|10404|
| CARTAGENA| 8333|
| ANTIOQUIA| 4554|
| NARIÑO| 3520|
|CUNDINAMARCA| 2827|
| AMAZONAS| 2317|
| CHOCO| 1636|
```

only showing top 10 rows

1 spark.sql("select municipio, count(\*) as count from covid19 group by municipio order by count desc").show(10)

only showing top 10 rows

+----+

```
1 spark.sql("""select `Fecha de inicio de síntomas` as fecha, count(*) as count from covid19 where
2 `Fecha de inicio de síntomas` is not NULL group by fecha order by count desc limit 10""").show()
              fecha count
    10/6/2020 0:00:00 2731
    16/6/2020 0:00:00 2558
    18/6/2020 0:00:00 2479
    12/6/2020 0:00:00 2452
    1/6/2020 0:00:00 2429
    8/6/2020 0:00:00 2390
    17/6/2020 0:00:00 2344
    | 5/6/2020 0:00:00| 2266|
    9/6/2020 0:00:00 2224
    19/6/2020 0:00:00 2162
1 spark.sql('select edad, age_group,count(*) as count from covid19 group by edad, age_group order by count desc').show()
    |edad| age_group|count|
    | 30|young adult| 2735|
      29|young adult| 2611|
      31
             senior 2569
      28 young adult | 2540
      27 young adult | 2494 |
      26 young adult | 2436 |
      33 senior 2371
      32
             senior 2362
      25 young adult 2335
      34
             senior 2310
      35 l
             senior 2292
      24|young adult| 2214|
      36 senior 2175
      37
              senior 2132
      38
             senior 2098
      40
             senior 2050
      23 young adult 2041
      39 senior 1985
      22|young adult| 1879|
      41 senior 1783
   only showing top 20 rows
1 spark.sql('''select age_group, count(*) as count from covid19 where
2 `Fecha de muerte` is not null group by age_group order by count desc''').show()
   +----+
    age_group|count|
    +-----+----+
| senior| 5435|
    |young adult| 147|
    | kid| 51|
```