

2.1 Columnas:

Google Colab	<pre>1 df.columns</pre> <pre>['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']</pre>
JupyterHub	<pre>In [5]: df.columns</pre> <pre>['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']</pre>

2.2 Tipos de datos:

Google Colab	JupyterHub
<pre>1 df.printSchema()</pre> <pre>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)</pre>	<pre>In [6]: df.printSchema()</pre> <pre>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)</pre>

2.3 Seleccionando columnas

GoogleColab

```
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	M	Importado
9/3/2020 0:00:00	3	F	Importado
11/3/2020 0:00:00	4	M	Relacionado
11/3/2020 0:00:00	5	M	Relacionado
11/3/2020 0:00:00	6	F	Relacionado
11/3/2020 0:00:00	7	F	Importado
11/3/2020 0:00:00	8	F	Importado
11/3/2020 0:00:00	9	F	Importado
12/3/2020 0:00:00	10	F	Importado

only showing top 10 rows

JupyterHub

```
In [7]: 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	M	Importado
9/3/2020 0:00:00	3	F	Importado
11/3/2020 0:00:00	4	M	Relacionado
11/3/2020 0:00:00	5	M	Relacionado
11/3/2020 0:00:00	6	F	Relacionado
11/3/2020 0:00:00	7	F	Importado
11/3/2020 0:00:00	8	F	Importado
11/3/2020 0:00:00	9	F	Importado
12/3/2020 0:00:00	10	F	Importado

only showing top 10 rows

GoogleColab

JupyterHub

```
1 df=df.withColumnRenamed('fecha_reporte web', 'fecha_reporte')

1 df=df.withColumnRenamed({'ID de caso': 'ID', 'Nombre del país': 'pais', 'Nombre departamento': 'departamento', 'Nombre municipio': 'municipio', 'Ubicación del caso': 'ubicación'})
```

```
In [8]: df=df.withColumnRenamed('fecha_reporte web', 'fecha_reporte')
```

```
In [9]: df=df.withColumnsRenamed({'ID de caso': 'ID', 'Nombre del país': 'pais', 'Nombre departamento': 'departamento',
<

```

```
In [10]: 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 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)
```

```
df = df.withColumn('edad_mas_18', df['edad'] > 18)
df.show(5, false)
```

ubicación	estado	código ISO del país	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	edad_mas_18
Casa	Leve	1280	ITALIA	recuperado	27/2/2020 0:00:00	NULL	14/3/2020 0:00:00	13/3/2020 0:00:00	PCR	6	NULL	128
Casa	Leve	1724	ESPAÑA	recuperado	4/3/2020 0:00:00	NULL	19/3/2020 0:00:00	19/3/2020 0:00:00	PCR	15	NULL	129
Casa	Leve	1724	ESPAÑA	recuperado	26/2/2020 0:00:00	NULL	15/3/2020 0:00:00	15/3/2020 0:00:00	PCR	16	NULL	130
Casa	Leve	NULL	NULL	recuperado	4/3/2020 0:00:00	NULL	11/3/2020 0:00:00	26/3/2020 0:00:00	PCR	16	NULL	165
Casa	Leve	NULL	NULL	recuperado	1/3/2020 0:00:00	NULL	23/3/2020 0:00:00	23/3/2020 0:00:00	PCR	16	NULL	35

JupyterHub	<pre>In [11]: df=df.withColumn('Edad_mas_10', df['Edad'] + 10) df.show(5, False)</pre> <pre> +-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+ fecha_reporte ID fecha de notificación Código DIVIPOLA departamento departamento Código DIVIPOLA municipio 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 Edad_mas_10 +-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+ 6/3/2020 0:00:00 1 2/3/2020 0:00:00 11 Casa Leve 300 ITALIA Recuperado 27/2/2020 0:00:00 PCR 6 NULL 29 </pre>
------------	---

2.6 Eliminando columnas

GoogleColab	<pre>1 columns_to_remove=['Nombre del grupo étnico', 'Pertenencia étnica'] 2 df=df.drop(*columns_to_remove)</pre>
JupyterHub	<pre>j): columns_to_remove=['Tipo de contagio', 'Unidad de medida de edad'] df=df.drop(*columns_to_remove) df.columns</pre> <pre> ['fecha_reporte', 'ID', 'fecha de notificación', 'Código DIVIPOLA departamento', 'departamento', 'Código DIVIPOLA municipio', 'municipio', 'Edad', 'Sexo', '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'] </pre>

2.7 Filtrando datos

Google Colab	<pre>1 filtered_df=df.filter(df.Sexo == 'F') 2 filtered_df.show(5, False)</pre> <pre> +-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+ fecha_reporte ID fecha de notificación Código DIVIPOLA departamento departamento Código DIVIPOLA municipio 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 Edad_mas_10 +-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+ 6/3/2020 0:00:00 1 2/3/2020 0:00:00 11 Casa Leve 300 ITALIA Recuperado 27/2/2020 0:00:00 PCR 6 NULL 29 9/3/2020 0:00:00 3 7/3/2020 0:00:00 5 ANTIOQUIA 5001 MEDELLIN 50 1 F Importado 6/3/2020 0:00:00 13/3/2020 0:00:00 PCR 6 NULL 29 11/3/2020 0:00:00 6 10/3/2020 0:00:00 5 ANTIOQUIA 5300 ITAGUI 27 1 F Importado 6/3/2020 0:00:00 13/3/2020 0:00:00 PCR 6 NULL 29 11/3/2020 0:00:00 7 10/3/2020 0:00:00 13001 CARTAGENA 13001 CARTAGENA 15 1 F Importado 6/3/2020 0:00:00 13/3/2020 0:00:00 PCR 6 NULL 29 11/3/2020 0:00:00 8 9/3/2020 0:00:00 11 BOGOTA 11001 BOGOTA 22 1 F Importado 6/3/2020 0:00:00 13/3/2020 0:00:00 PCR 6 NULL 29 </pre> <p>only showing top 5 rows</p>
JupyterHub	<pre>In [14]: filtered_df=df.filter(df.Sexo == 'F') filtered_df.show(5, False)</pre> <pre> +-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+ fecha_reporte ID fecha de notificación Código DIVIPOLA departamento departamento Código DIVIPOLA municipio 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 Edad_mas_10 +-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+ 6/3/2020 0:00:00 1 2/3/2020 0:00:00 11 Casa Leve 300 ITALIA Recuperado 27/2/2020 0:00:00 PCR 6 NULL 29 9/3/2020 0:00:00 3 7/3/2020 0:00:00 5 ANTIOQUIA 5001 MEDELLIN 50 1 F Importado 6/3/2020 0:00:00 13/3/2020 0:00:00 PCR 6 NULL 29 11/3/2020 0:00:00 6 10/3/2020 0:00:00 5 ANTIOQUIA 5300 ITAGUI 27 1 F Importado 6/3/2020 0:00:00 13/3/2020 0:00:00 PCR 6 NULL 29 11/3/2020 0:00:00 7 10/3/2020 0:00:00 13001 CARTAGENA 13001 CARTAGENA 15 1 F Importado 6/3/2020 0:00:00 13/3/2020 0:00:00 PCR 6 NULL 29 11/3/2020 0:00:00 8 9/3/2020 0:00:00 11 BOGOTA 11001 BOGOTA 22 1 F Importado 6/3/2020 0:00:00 13/3/2020 0:00:00 PCR 6 NULL 29 </pre>

Google Colab	<pre>1 filtered2_df=df.filter("Sexo=='F' and Edad>20 and Recuperado like '%Recuperado%' and municipio like 'MEDELLIN'") 2 filtered2_df.show()</pre> <pre> +-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+ fecha_reporte ID fecha de notificación Código DIVIPOLA departamento departamento Código DIVIPOLA municipio 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 Edad_mas_10 +-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+ 9/3/2020 0:00:00 3 7/3/2020 0:00:00 5 ANTIOQUIA 5001 MEDELLIN 50 1 F Importado 6/3/2020 0:00:00 13/3/2020 0:00:00 PCR 6 NULL 29 14/3/2020 0:00:00 20 11/3/2020 0:00:00 5 ANTIOQUIA 5001 MEDELLIN 50 1 F Importado 6/3/2020 0:00:00 13/3/2020 0:00:00 PCR 6 NULL 29 19/3/2020 0:00:00 108 17/3/2020 0:00:00 5 ANTIOQUIA 5001 MEDELLIN 50 1 F Importado 6/3/2020 0:00:00 13/3/2020 0:00:00 PCR 6 NULL 29 20/3/2020 0:00:00 131 15/3/2020 0:00:00 5 ANTIOQUIA 5001 MEDELLIN 50 1 F Importado 6/3/2020 0:00:00 13/3/2020 0:00:00 PCR 6 NULL 29 20/3/2020 0:00:00 135 17/3/2020 0:00:00 5 ANTIOQUIA 5001 MEDELLIN 50 1 F Importado 6/3/2020 0:00:00 13/3/2020 0:00:00 PCR 6 NULL 29 20/3/2020 0:00:00 141 17/3/2020 0:00:00 5 ANTIOQUIA 5001 MEDELLIN 50 1 F Importado 6/3/2020 0:00:00 13/3/2020 0:00:00 PCR 6 NULL 29 20/3/2020 0:00:00 142 20/3/2020 0:00:00 5 ANTIOQUIA 5001 MEDELLIN 50 1 F Importado 6/3/2020 0:00:00 13/3/2020 0:00:00 PCR 6 NULL 29 22/3/2020 0:00:00 238 16/3/2020 0:00:00 5 ANTIOQUIA 5001 MEDELLIN 50 1 F Importado 6/3/2020 0:00:00 13/3/2020 0:00:00 PCR 6 NULL 29 23/3/2020 0:00:00 268 19/3/2020 0:00:00 5 ANTIOQUIA 5001 MEDELLIN 50 1 F Importado 6/3/2020 0:00:00 13/3/2020 0:00:00 PCR 6 NULL 29 23/3/2020 0:00:00 271 18/3/2020 0:00:00 5 ANTIOQUIA 5001 MEDELLIN 50 1 F Importado 6/3/2020 0:00:00 13/3/2020 0:00:00 PCR 6 NULL 29 23/3/2020 0:00:00 272 18/3/2020 0:00:00 5 ANTIOQUIA 5001 MEDELLIN 50 1 F Importado 6/3/2020 0:00:00 13/3/2020 0:00:00 PCR 6 NULL 29 23/3/2020 0:00:00 275 23/3/2020 0:00:00 5 ANTIOQUIA 5001 MEDELLIN 50 1 F Importado 6/3/2020 0:00:00 13/3/2020 0:00:00 PCR 6 NULL 29 23/3/2020 0:00:00 292 19/3/2020 0:00:00 5 ANTIOQUIA 5001 MEDELLIN 50 1 F Importado 6/3/2020 0:00:00 13/3/2020 0:00:00 PCR 6 NULL 29 23/3/2020 0:00:00 293 18/3/2020 0:00:00 5 ANTIOQUIA 5001 MEDELLIN 50 1 F Importado 6/3/2020 0:00:00 13/3/2020 0:00:00 PCR 6 NULL 29 23/3/2020 0:00:00 294 21/3/2020 0:00:00 5 ANTIOQUIA 5001 MEDELLIN 50 1 F Importado 6/3/2020 0:00:00 13/3/2020 0:00:00 PCR 6 NULL 29 23/3/2020 0:00:00 296 19/3/2020 0:00:00 5 ANTIOQUIA 5001 MEDELLIN 50 1 F Importado 6/3/2020 0:00:00 13/3/2020 0:00:00 PCR 6 NULL 29 25/3/2020 0:00:00 438 20/3/2020 0:00:00 5 ANTIOQUIA 5001 MEDELLIN 50 1 F Importado 6/3/2020 0:00:00 13/3/2020 0:00:00 PCR 6 NULL 29 25/3/2020 0:00:00 450 19/3/2020 0:00:00 5 ANTIOQUIA 5001 MEDELLIN 50 1 F Importado 6/3/2020 0:00:00 13/3/2020 0:00:00 PCR 6 NULL 29 25/3/2020 0:00:00 466 19/3/2020 0:00:00 5 ANTIOQUIA 5001 MEDELLIN 50 1 F Importado 6/3/2020 0:00:00 13/3/2020 0:00:00 PCR 6 NULL 29 28/3/2020 0:00:00 595 16/3/2020 0:00:00 5 ANTIOQUIA 5001 MEDELLIN 50 1 F Importado 6/3/2020 0:00:00 13/3/2020 0:00:00 PCR 6 NULL 29 </pre> <p>only showing top 20 rows</p>
--------------	---

JupyterHub	<pre> filtered2_df=df.filter("Sexo=='F' and Edad>20 and Recuperado like '%Recuperado%' and municipio like 'MEDELLIN'") filtered2_df.show() </pre> <pre> +-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+ +-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+ +-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+ +-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+ +-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+ fecha_reporte ID Fecha de notificación Código DIVIPOLA departamento departamento Código DIVIPOLA municipio municipio Edad -----+-----+-----+-----+-----+-----+-----+-----+-----+-----+ 9/3/2020 0:00:00 3 7/3/2020 0:00:00 5 ANTIOQUIA 724 5001 MEDELLIN 50 29/2/2020 0:00:00 1 F Importado Casa Leve NULL PCR 6 14/3/2020 0:00:00 20 11/3/2020 0:00:00 5 ANTIOQUIA 724 5001 MEDELLIN 26 9/3/2020 0:00:00 1 F Relacionado Casa Leve NULL PCR 6 19/3/2020 0:00:00 108 17/3/2020 0:00:00 5 ANTIOQUIA 724 5001 MEDELLIN 57 14/3/2020 0:00:00 67 19/3/2020 0:00:00 28/3/2020 0:00:00 724 PCR 6 20/3/2020 0:00:00 131 15/3/2020 0:00:00 5 ANTIOQUIA 724 5001 MEDELLIN 22 13/3/2020 0:00:00 32 28/3/2020 0:00:00 28/3/2020 0:00:00 724 PCR 6 20/3/2020 0:00:00 135 17/3/2020 0:00:00 5 ANTIOQUIA 276 5001 MEDELLIN 44 16/3/2020 0:00:00 54 20/3/2020 0:00:00 1/4/2020 0:00:00 724 PCR 6 20/3/2020 0:00:00 141 17/3/2020 0:00:00 5 ANTIOQUIA 191 5001 MEDELLIN 62 1 F Importado Casa Leve 191 PCR 6 </pre>
------------	---

2.8 Usando lambda

Google Colab	<pre> 1 age_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) </pre> <pre> +-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+ +-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+ +-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+ +-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+ +-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+ fecha_reporte ID Fecha de notificación Código DIVIPOLA departamento departamento Código DIVIPOLA municipio municipio Edad -----+-----+-----+-----+-----+-----+-----+-----+-----+-----+ 19/6/2020 0:00:00 62406 10/6/2020 0:00:00 11 BOGOTA 11001 NULL BOGOTA 30 1 F Relacionado Casa Leve NULL PCR 6 19/6/2020 0:00:00 40 young adult 19/6/2020 0:00:00 12/7/2020 0:00:00 PCR 6 21/3/2020 0:00:00 192 20/3/2020 0:00:00 76 VALLE 76001 NULL CALI 30 1 F Relacionado Casa Leve NULL PCR 6 17/3/2020 0:00:00 40 young adult 21/3/2020 0:00:00 2/4/2020 0:00:00 PCR 6 19/6/2020 0:00:00 62398 17/6/2020 0:00:00 11 BOGOTA 11001 NULL BOGOTA 30 1 F Relacionado Casa Leve NULL PCR 6 14/6/2020 0:00:00 40 young adult 19/6/2020 0:00:00 23/7/2020 0:00:00 Tiempo 5 27/3/2020 0:00:00 517 23/3/2020 0:00:00 11 BOGOTA 11001 NULL BOGOTA 30 1 F Relacionado Casa Leve NULL PCR 6 17/3/2020 0:00:00 40 young adult 27/3/2020 0:00:00 15/4/2020 0:00:00 PCR 6 19/6/2020 0:00:00 62380 10/6/2020 0:00:00 11 BOGOTA 11001 NULL BOGOTA 30 1 F Comunitaria Casa Leve NULL PCR 6 19/6/2020 0:00:00 11001 19/6/2020 0:00:00 11/7/2020 0:00:00 PCR 6 </pre>
JupyterHub	<pre> age_udf=udf(lambda age: "kid" if age < 18 else 'young adult' if age<=30 else 'senior'), StringType() df=df.withColumn("age_group", age_udf(df.Edad)) df.filter(df['age_group']=='young adult').orderBy('Edad', ascending=False).show(10, False) </pre> <pre> +-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+ +-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+ +-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+ +-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+ +-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+ fecha_reporte ID Fecha de notificación Código DIVIPOLA departamento departamento Código DIVIPOLA municipio municipio Edad -----+-----+-----+-----+-----+-----+-----+-----+-----+-----+ 19/6/2020 0:00:00 62406 10/6/2020 0:00:00 11 BOGOTA 11001 NULL BOGOTA 30 1 F Relacionado Casa Leve NULL PCR 6 19/6/2020 0:00:00 40 young adult 19/6/2020 0:00:00 12/7/2020 0:00:00 PCR 6 21/3/2020 0:00:00 192 20/3/2020 0:00:00 76 VALLE 76001 NULL CALI 30 1 F Relacionado Casa Leve NULL PCR 6 17/3/2020 0:00:00 40 young adult 21/3/2020 0:00:00 2/4/2020 0:00:00 PCR 6 19/6/2020 0:00:00 62398 17/6/2020 0:00:00 11 BOGOTA 11001 NULL BOGOTA 30 1 F Relacionado Casa Leve NULL PCR 6 14/6/2020 0:00:00 40 young adult 19/6/2020 0:00:00 23/7/2020 0:00:00 Tiempo 5 27/3/2020 0:00:00 517 23/3/2020 0:00:00 11 BOGOTA 11001 NULL BOGOTA 30 1 F Relacionado Casa Leve NULL PCR 6 17/3/2020 0:00:00 40 young adult 27/3/2020 0:00:00 15/4/2020 0:00:00 PCR 6 19/6/2020 0:00:00 62380 10/6/2020 0:00:00 11 BOGOTA 11001 NULL BOGOTA 30 1 F Comunitaria Casa Leve NULL PCR 6 19/6/2020 0:00:00 11001 19/6/2020 0:00:00 11/7/2020 0:00:00 PCR 6 </pre>

Punto 3 en Google Colab

3.1

Spark Dataframes	<pre> 1 df.groupby('departamento').count().orderBy('count', ascending=False).show(10, False) </pre> <pre> +-----+-----+ departamento count +-----+-----+ BOGOTA 30016 BARRANQUILLA 13865 ATLANTICO 18994 VALLE 18404 CARTAGENA 8333 ANTIOQUIA 4554 INDIERO 3528 CUNDINAMARCA 2857 AMAZONAS 2317 CHOCO 1636 +-----+-----+ only showing top 10 rows </pre>
------------------	--

SparkSQL	<pre>[] 1 df.createOrReplaceTempView("covid19")</pre> <pre>1 spark.sql("select departamento, count(*) as count from covid19 group by departamento order by count desc").show(10)</pre> <pre>+-----+-----+ 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</pre>
----------	---

3.2

Spark Dataframes	<pre>1 df.groupBy('municipio').count().orderBy('count', ascending=False).show(10, False)</pre> <pre>+-----+-----+ municipio count +-----+-----+ BOGOTA 30016 BARRANQUILLA 13065 CARTAGENA 8333 CALI 7747 SOLEDAD 6233 LETICIA 2194 MEDELLIN 2137 TUMACO 1501 BUENAVENTURA 1453 QUIBDO 1367 +-----+-----+ only showing top 10 rows</pre>
Spark SQL	<pre>1 spark.sql("select municipio, count(*) as count from covid19 group by municipio order by count desc").show(10)</pre> <pre>+-----+-----+ municipio count +-----+-----+ BOGOTA 30016 BARRANQUILLA 13065 CARTAGENA 8333 CALI 7747 SOLEDAD 6233 LETICIA 2194 MEDELLIN 2137 TUMACO 1501 BUENAVENTURA 1453 QUIBDO 1367 +-----+-----+ only showing top 10 rows</pre>

3.3

Dataframes	<pre>1 df.dropna(subset=['Fecha de inicio de síntomas']).groupBy('Fecha de inicio de síntomas').count().orderBy('count', ascending=False).show(10, False)</pre> <pre>+-----+-----+ Fecha de inicio de síntomas 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 +-----+-----+ only showing top 10 rows</pre>
SparkSQL	<pre>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()</pre> <pre>+-----+-----+ 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 +-----+-----+</pre>

3.4

Dataframes	<pre>1 df.select('edad', 'age_group', 'ID').groupBy('edad', 'age_group').count().orderBy('count', ascending=False).show()</pre> <pre>+-----+-----+-----+ 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 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</pre>
SparkSQL	<pre>1 spark.sql('select edad, age_group,count(*) as count from covid19 group by edad, age_group order by count desc').show()</pre> <pre>+-----+-----+-----+ 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 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</pre>

3.5

Muertes por grupo de edad

Dataframes	<pre>1 from pyspark.sql.functions import col 2 df.filter(col('Fecha de muerte').isNotNull()).groupBy('age_group').count().orderBy('count', ascending=False).show()</pre> <pre>+-----+-----+ age_group count +-----+-----+ senior 5435 young adult 151 kid 47 +-----+-----+</pre>
SparkSQL	<pre>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()</pre> <pre>+-----+-----+ age_group count +-----+-----+ senior 5435 young adult 147 kid 51 +-----+-----+</pre>

Punto 3 JupyterHub

3.1

Spark Dataframes	<pre>In [22]: df.groupBy('departamento').count().orderBy('count', ascending=False).show(10, False)</pre> <pre>+-----+-----+ 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</pre>
------------------	--

SparkSQL	<pre>In [79]: spark.sql("select departamento, count(*) as count from covid19 group by departamento order by count desc").show(10)</pre> <pre> +-----+ departamento count +-----+ BOGOTA 30016 BARRANQUILLA 13065 ATLANTICO 10994 VALLE 10404 CARTAGENA 8333 ANTIOQUIA 4554 NARIÑO 3520 CUNDINAMARCA 3027 AMAZONAS 2317 CHOCO 1636 +-----+ only showing top 10 rows </pre>
----------	---

3.2

Spark Dataframes	<pre>In [24]: df.groupBy('municipio').count().orderBy('count', ascending=False).show(10, False)</pre> <pre> +-----+ municipio count +-----+ BOGOTA 30016 BARRANQUILLA 13065 CARTAGENA 8333 CALI 7747 SOLEDAD 6233 LETICIA 2194 MEDELLIN 2137 TUMACO 1501 BUENAVENTURA 1453 QUIBDO 1367 +-----+ only showing top 10 rows </pre>
Spark SQL	<pre>In [80]: spark.sql("select municipio, count(*) as count from covid19 group by municipio order by count desc").show(10)</pre> <pre> +-----+ municipio count +-----+ BOGOTA 30016 BARRANQUILLA 13065 CARTAGENA 8333 CALI 7747 SOLEDAD 6233 LETICIA 2194 MEDELLIN 2137 TUMACO 1501 BUENAVENTURA 1453 QUIBDO 1367 +-----+ only showing top 10 rows </pre>

3.3

Dataframes	<pre>In [71]: df.dropna(subset=['fecha_reporte']).groupBy('fecha_reporte').count().orderBy('count', ascending=False).show(10, False)</pre> <pre> +-----+ fecha_reporte count +-----+ 27/6/2020 0:00:00 4149 26/6/2020 0:00:00 3543 24/6/2020 0:00:00 3541 25/6/2020 0:00:00 3486 29/6/2020 0:00:00 3274 28/6/2020 0:00:00 3178 18/6/2020 0:00:00 3171 19/6/2020 0:00:00 3059 21/6/2020 0:00:00 3019 30/6/2020 0:00:00 2803 +-----+ only showing top 10 rows </pre>
SparkSQL	<pre>In [53]: spark.sql("""select fecha_reporte as fecha, count(*) as count from covid19 where fecha_reporte is not NULL group by fecha order by count desc limit 10""").show()</pre> <pre> +-----+ fecha count +-----+ 27/6/2020 0:00:00 4149 26/6/2020 0:00:00 3543 24/6/2020 0:00:00 3541 25/6/2020 0:00:00 3486 29/6/2020 0:00:00 3274 28/6/2020 0:00:00 3178 18/6/2020 0:00:00 3171 19/6/2020 0:00:00 3059 21/6/2020 0:00:00 3019 30/6/2020 0:00:00 2803 +-----+ </pre>

3.4

Dataframes	<pre>In [20]: df.select('edad', 'age_group', 'ID').groupBy('edad', 'age_group').count().orderBy('count', ascending=False).show()</pre> <pre> +-----+ edad age_group count +-----+ 30 young adult 2735 28 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 38 senior 2090 40 senior 2050 23 young adult 2041 39 senior 1985 22 young adult 1879 41 senior 1793 +-----+ only showing top 20 rows </pre>
------------	--

SparkSQL	<pre>In [38]: spark.sql('select edad, age_group, count(*) as count from covid19 group by edad, age_group order by count desc').show()</pre> <pre> +-----+-----+ 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 2336 34 senior 2310 35 senior 2292 24 young adult 2214 36 senior 2175 37 senior 2132 38 senior 2090 40 senior 2050 23 young adult 2041 39 senior 1985 22 young adult 1879 41 senior 1703 +-----+-----+ only showing top 20 rows </pre>
----------	---

3.5

Muertes por grupo de edad

Dataframes	<pre>In [77]: from pyspark.sql.functions import col df.filter(col('Fecha de muerte').isNotNull()).groupBy('age_group').count().orderBy('count', ascending=False).show()</pre> <pre> +-----+-----+ age_group count +-----+-----+ senior 5435 young adult 151 kid 47 +-----+-----+ </pre>
SparkSQL	<pre>In [36]: spark.sql('select age_group, count(*) as count from covid19 where Fecha de muerte is not null group by age_group order by count desc').show()</pre> <pre> +-----+-----+ age_group count +-----+-----+ senior 5435 young adult 151 kid 47 +-----+-----+ </pre>

URI S3: s3://davidnotebook/jupyter/