

## Spark Dataframes(14/12/2018)

### Carga de datos

```
In []: tabla=[("Marta",25),
              ("Pedro",30),
              ("María",24)]
df=spark.createDataFrame(tabla,["nombre","edad"])
df.show()
```

Creación de un dataframe con encabezados

```
Out[]: +-----+-----+
|nombre|edad|
+-----+-----+
| Marta|  25|
| Pedro|  30|
| María|  24|
+-----+-----+
```

Visualización del dataframe

```
In []: df.printSchema()
```

Muestra la estructura

```
Out[]: root
|-- nombre: string (nullable = true)
|-- edad: long (nullable = true)
```

Estructura del dataframe

```
In []: from pyspark.sql.types import *

estructura=StructType([
    StructField("nombre",StringType()),
    StructField("edad",ByteType())
])

df2=spark.createDataFrame(tabla,estructura)
```

Creación de un dataframe con estructura

### Datos externos

```
In []: dfVentas=spark.read.json("ventas.json")
dfVentas.show()
```

Carga de un archivo json

```
Out[]: +-----+-----+-----+
|Producto| Año| Precio|
+-----+-----+-----+
|    moto| null| 1000.00|
|   coche| null| 2000.00|
|    bici| null|  200.00|
|   coche| 2015| 3000.00|
+-----+-----+-----+
```

```
In []: dfVendedores=spark.read.csv("vendedores.csv",
                                   sep=";",header=True,inferSchema=True)
```

Carga de un archivo csv

### Recuperación de información

```
In []: df.summary().show()
```

```
Out[]: +-----+-----+-----+
|summary|nombre|          edad|
+-----+-----+-----+
|  count|      3|              3|
|   mean|   null|26.333333333333332|
| stddev|   null|3.2145502536643185|
|    min| Marta|             24|
|   25%| null|             24|
|   50%| null|             25|
|   75%| null|             30|
|    max| Pedro|             30|
+-----+-----+-----+
```

```
In []: df.collect()
```

```
Out[]: [Row(nombre='Marta', edad=25),
        Row(nombre='Pedro', edad=30),
        Row(nombre='María', edad=24)]
```

Devuelve el conjunto como lista

```
In []: df.count()
```

Registros en el DF

```
Out[]: 3
```

```
In []: df.select("edad","nombre").show()
```

Selecciona columnas

```
Out[]: +-----+-----+
|edad|nombre|
+-----+-----+
|  25| Marta|
|  30| Pedro|
|  24| María|
+-----+-----+
```

```
In []: df.select("nombre",(df.edad+10)\
                .alias("edad ahora")).show()
```

Selecciona y modifica columnas

```
Out[]: +-----+-----+
|edad|nombre|
+-----+-----+
|  25| Marta|
|  30| Pedro|
|  24| María|
+-----+-----+
```

```
In []: df.where("edad=30").show()
#df.where(df.edad==30).show()
```

Selecciona filas en función de la condición

```
Out[]: +-----+-----+
|nombre|edad|
+-----+-----+
| Pedro|  30|
+-----+-----+
```

```
In []: grupo=dfVentas.groupBy("Producto")
grupo.sum("Precio").show()
```

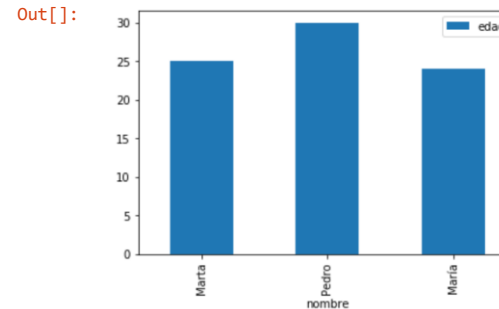
Agrupar por producto y suma el precio

```
Out[]: +-----+-----+-----+
|Producto|sum(Precio)|
+-----+-----+-----+
|   coche|    5000.00|
|    moto|    1000.00|
|    bici|     200.00|
+-----+-----+-----+
```

### Pandas

```
In []: dfp=df.toPandas()
dfp.plot.bar(x="nombre")
```

Cambia a pandas y realiza una gráfica



```
In []: from autovizwidget.widget.utils import
display_dataframe
display_dataframe(df)
```

Exploración de datos con autovizwidget

Out []:

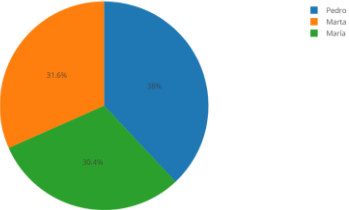
Type: Table Pie Scatter Line Area Bar

Encoding:

X: nombre

Y: edad

force Max



DATAFRAMES [\[doc\]](#)

<code>.cache()</code>	Almacena el Dataframe	<code>.randomSplit(weights, seed=None)</code>	Divide el dataframe aleatoriamente
<code>.collect()</code>	Lista de filas(Row)	<code>.rdd</code>	Devuelve RDD
<code>.columns</code>	Nombres de columnas	<code>.replace(to_replace,value=&lt;novalue&gt;, subset=None)</code>	Realiza reemplazos en el dataframe
<code>.corr(col1, col2)</code>	Correlación de Pearson	<code>.select(*cols)</code>	Selecciona columnas
<code>.count()</code>	Número de registros	<code>.show(n=20, truncate=True, vertical=False)</code>	Muestra el dataframe
<code>.crosstab(col1, col2)</code>	Tabla de contingencia	<code>.stat</code>	Valores estadísticos
<code>.distinct()</code>	DF sin filas repetidas	<code>.summary(*statistics)</code>	Sumario estadist.
<code>.drop(*cols)</code>	Descarta columnas	<code>.take(num)</code>	
<code>.filter(condition)</code> <code>.where(condition)</code>	Selecciona columnas por condición	<code>.toDF(*cols)</code>	Cambia el nombre de las columnas
<code>.groupBy(*cols)</code>	Agrupar por columnas	<code>.toPandas()</code>	Convierte a Pandas
<code>.join(oter, on=None, how=None)</code>	Combina Dataframes	<code>.union(oter)</code>	Une dos dataframes
<code>.na</code>	Métodos para campos NULL	<code>.withColumn(colName, col)</code>	Añade o modifica columna
<code>.orderBy(*cols,**kwargs)</code> <code>.sort(*cols, **kwargs)</code>	Ordena por columnas	<code>.withColumnRenamed(existing, new)</code>	Cambia el nombre de una columna
<code>.printSchema()</code>	Muestra estructura	<code>.write</code>	Graba el DF

GRUPOS [\[doc\]](#) .groupBy(col)

<code>.sum(*cols)</code>	Totales por grupo	<code>.min(*cols)</code>	Mínimo por grupo
<code>.avg(*cols)</code> <code>.mean(*cols)</code>	Medias agrupadas	<code>.agg(*exprs)</code>	Aplica funciones por grupo
<code>.count()</code>	Registros por grupo	<code>.apply(udf)</code>	Aplica una función de usuario a cada grupo
<code>.max(*cols)</code>	Máximo por grupo	<code>.pivot(pivot_col, values=None)</code>	Reestructura el dataframe en función de la columna

COLUMNAS [\[doc\]](#)

Formas de indicar columnas: "nombreCol", df.nombreCol, df["nombreCol"], col("nombreCol")

<code>.alias(*alias, **kwargs)</code>	Renombra columna	<code>.contains(oter)</code>	Comprueba contenido
<code>.asc()</code> <code>.desc()</code>	Orden ascendente Orden descendente	<code>.isNotNull()</code> <code>.isNull()</code>	Comprueba no Null Comprueba si Null
<code>.cast(dataType)</code>	Convierte tipo de dato	<code>.like(oter)</code> <code>.rlike(oter)</code>	Coincidencia patrón SQL Coincidencia exp. regular
<code>.between(lower, upper)</code>	Comprueba intervalo	<code>.substr(start, length)</code>	Recupera subcadena
<code>.endswith(oter)</code> <code>.startswith(oter)</code>	Comprueba inicio cadena	<code>.when(condition, value)</code> <code>.otherwise(value)</code>	Devuelve un valor u otro si se cumple condición

TIPOS DE DATOS [\[doc\]](#)

<code>DataType</code>	<code>BinaryType</code>	<code>TimestampType</code>	<code>FloatType</code>	<code>LongType</code>	<code>MapType</code>
<code>NullType</code>	<code>BooleanType</code>	<code>DecimalType</code>	<code>ByteType</code>	<code>ShortType</code>	<code>StructField</code>
<code>StringType</code>	<code>DateType</code>	<code>DoubleType</code>	<code>IntegerType</code>	<code>ArrayType</code>	<code>StructType</code>

FUNCIONES [\[doc\]](#)

<code>abs</code>	<code>acos</code>	<code>add_months</code>	<code>approxCountDistinctD</code>
<code>ascii</code>	<code>array</code>	<code>array_contains</code>	<code>asc</code>
<code>avg</code>	<code>asin</code>	<code>atan</code>	<code>atan2</code>
<code>broadcast</code>	<code>base64</code>	<code>bin</code>	<code>bitwiseNOT</code>
<code>coalesce</code>	<code>brround</code>	<code>cbirt</code>	<code>ceil</code>
<code>column</code>	<code>col</code>	<code>collect_list</code>	<code>collect_set</code>
<code>corr</code>	<code>concat</code>	<code>concat_ws</code>	<code>conv</code>
<code>countDistinct</code>	<code>cos</code>	<code>cosh</code>	<code>count</code>
<code>create_map</code>	<code>covar_pop</code>	<code>covar_samp</code>	<code>crc32</code>
<code>date_add</code>	<code>cume_dist</code>	<code>current_date</code>	<code>current_timestamp</code>
<code>datediff</code>	<code>date_format</code>	<code>date_sub</code>	<code>date_trunc</code>
<code>decode</code>	<code>dayofmonth</code>	<code>dayofweek</code>	<code>dayofyear</code>
<code>encode</code>	<code>degrees</code>	<code>dense_rank</code>	<code>desc</code>
<code>expml</code>	<code>exp</code>	<code>explode</code>	<code>explode_outer</code>
<code>floor</code>	<code>expr</code>	<code>factorial</code>	<code>first</code>
<code>from_unixtime</code>	<code>format_number</code>	<code>format_string</code>	<code>from_json</code>
<code>grouping</code>	<code>from_utc_timestamp</code>	<code>get_json_object</code>	<code>greatest</code>
<code>hour</code>	<code>grouping_id</code>	<code>hash</code>	<code>hex</code>
<code>instr</code>	<code>hypot</code>	<code>initcap</code>	<code>input_file_name</code>
<code>kurtosis</code>	<code>isnan</code>	<code>isnull</code>	<code>json_tuple</code>
<code>lead</code>	<code>lag</code>	<code>last</code>	<code>last_day</code>
<code>lit</code>	<code>least</code>	<code>length</code>	<code>levenshtein</code>
<code>log1p</code>	<code>locate</code>	<code>log</code>	<code>log10</code>
<code>ltrim</code>	<code>log2</code>	<code>lower</code>	<code>lpad</code>
<code>md5</code>	<code>map_keys</code>	<code>map_values</code>	<code>max</code>
<code>next_day</code>	<code>mean</code>	<code>min</code>	<code>minute</code>
<code>posexplode</code>	<code>month</code>	<code>months_between</code>	<code>nanvl</code>
<code>radians</code>	<code>ntile</code>	<code>pandas_udfE</code>	<code>percent_rank</code>
<code>regexp_extract</code>	<code>posexplode_outer</code>	<code>pow</code>	<code>quarter</code>
<code>rint</code>	<code>rand</code>	<code>randn</code>	<code>rank</code>
<code>rtrim</code>	<code>regexp_replace</code>	<code>repeat</code>	<code>reverse</code>
<code>shiftLeft</code>	<code>round</code>	<code>row_number</code>	<code>rpadd</code>
<code>sin</code>	<code>second</code>	<code>shai</code>	<code>sha2</code>
<code>sort_array</code>	<code>shiftRight</code>	<code>shiftRightUnsigned</code>	<code>signum</code>
<code>sqrt</code>	<code>sinh</code>	<code>size</code>	<code>skewness</code>
<code>struct</code>	<code>soundex</code>	<code>spark_partition_id</code>	<code>split</code>
<code>sumDistinct</code>	<code>stddev</code>	<code>stddev_pop</code>	<code>stddev_samp</code>
<code>toRadiansD</code>	<code>substring</code>	<code>substring_index</code>	<code>sum</code>
<code>to_utc_timestamp</code>	<code>tan</code>	<code>tanh</code>	<code>toDegreesD</code>
<code>udf</code>	<code>to_date</code>	<code>to_json</code>	<code>to_timestamp</code>
<code>upper</code>	<code>translate</code>	<code>trim</code>	<code>trunc</code>
<code>unbase64</code>	<code>unhex</code>	<code>unix_timestamp</code>	<code>weekofyear</code>
<code>var_pop</code>	<code>var_samp</code>	<code>variance</code>	<code>when</code>
<code>window</code>	<code>year</code>		

<https://github.com/bigdatagalicia>