

SPARK

- CREAR UNA NUEVA SESION DE SPARK: `SPARK = SPARKSESSION.BUILDER.MASTER("LOCAL[*]").GETORCREATE()`
- ABRIR UN ARCHIVO, MOSTRAR EL ESQUEMA Y PROCESAR LA PRIMERA LINEA COMO HEADER:
`CUSTOMERS = SPARK.READ.OPTION("INFERSCHEMA", "TRUE").OPTION("HEADER", "TRUE").CSV("CUSTOMERS.CSV")`
`PRODUCTS = SPARK.READ.OPTION("INFERSCHEMA", "TRUE").OPTION("HEADER", "TRUE").CSV("PRODUCTS.CSV")`
`STOCK = SPARK.READ.OPTION("INFERSCHEMA", "TRUE").OPTION("HEADER", "TRUE").CSV("STOCK.CSV")`
`CUSTOMERS.PRINTSCHEMA()`
`PRODUCTS.PRINTSCHEMA()`
`STOCK.PRINTSCHEMA()`

- MOSTRAR INFORMACIÓN DEL DATAFRAME:

`CUSTOMERS.PRINTSCHEMA()` -> MUESTRA LOS NOMBRES DE LAS COLUMNAS Y EL TIPO DE DATO QUE CONTIENEN

`CUSTOMERS.SHOW(5)` -> MUESTRA LAS CINCO PRIMERAS FILAS

`CUSTOMERS.DTYPES` -> MUESTRA UN VECTOR DE LOS TIPOS DE DATOS

`CUSTOMERS.SUMMARY()` -> MUESTRA UNA TABLA CON LAS OPERACIONES BÁSICAS DE LAS COLUMNAS

EJEMPLO SHOW:

```
+-----+-----+-----+-----+-----+-----+
|      date|      time|customer|product|quantity|price|
|05/10/2018| 2:20 PM|    100|      1|        10|  816|
|06/10/2018| 3:30 PM|    100|      1|        10|    1|
|07/10/2018| 5:20 PM|    100|      1|        10|   10|
|04/08/2018|11:38 PM|    100|      2|         8|   79|
|25/03/2018| 3:52 AM|    100|      3|         1|   91|
+-----+-----+-----+-----+-----+-----+-----+
```

- ¿CUÁLES SON LOS PRODUCTOS EN MI DATASET?
`CUSTOMERS.SELECT(EXPR("PRODUCT AS MY_PRODUCTS")).SHOW`
- ¿CUALES SON LOS PRODUCTOS Y CANTIDADES EN MI DATASET?
`CUSTOMERS.SELECT(EXPR(QUANTITY), EXPR(PRODUCT)).SHOW`
- SELECCIONAR TODA LA TABLA CON `*`
`C.SELECT(EXPR("*")).SHOW()`
- ¿QUÉ COMPRAS SE HAN HECHO A PRODUCTOS CON MÁS DE UNA CANTIDAD?
`CUSTOMERS.SELECTEXPR("*", "QUANTITY>1").SHOW` -> AÑADE UNA NUEVA COLUMNA QUE ES TRUE SI ESA COMPRA TIENE MÁS DE 1 UNIDAD, O FALSE SI NO.
- ¿CUÁNTOS PRODUCTOS SE HAN COMPRADO EN TOTAL?
`CUSTOMERS.SELECTEXPR("SUM(QUANTITY)").SHOW`
- ¿CUÁL ES LA MEDIA DEL PRECIO?
`CUSTOMERS.SELECTEXPR("AVG(PRICE)").SHOW`
- ¿CUÁL ES LA DIFERENCIA ENTRE...?
`C.SELECTEXPR("COUNT(CUSTOMER)").SHOW()` -> CUENTA EL NÚMERO DE FILAS EN LA COLUMNA CUSTOMER
`C.SELECTEXPR("COUNT(DISTINCT(CUSTOMER))").SHOW()` -> CUENTA EL NUMERO DE FILAS DIFERENTES EN LA COLUMNA CUSTOMER.

- ¿CUÁL ES EL PRECIO MEDIO Y EL NÚMERO DE VENTAS?
`CUSTOMERS.SELECTEXPR("AVG(PRICE)", "COUNT(CUSTOMER)").SHOW`
- ¿QUE COMPRAS A PRODUCTOS TIENEN MENOS DE 9 CANTIDADES?
`CUSTOMERS.WHERE("QUANTITY<9").SHOW`
- ¿QUE COMPRAS A PRODUCTOS NO LAS HA HECHO EL CUSTOMER 100?
`CUSTOMERS.WHERE("QUANTITY !=100").SHOW`
- ¿QUE COMPRAS A PRODUCTOS NO SON DEL CUSTOMER 100 Y TIENEN MENOS DE 9 UNIDADES?
`CUSTOMERS.WHERE("CUSTOMER!=100", "QUANTITY<9").SHOW`
- ¿QUÉ COMPRADORES HAN COMPRADO 7 O MÁS ÍTEMS DEL PRODUCTO 8?
`CUSTOMER.WHERE("QUANTITY>=7", "PRODUCT=8").SHOW`
- ¿CUÁL ES EL PRECIO MÁS CARO?
`CUSTOMERS.ORDERBY("PRICE").SHOW(1)`
- ORDENAR DESCENDENTEMENTE POR CUSTOMER Y ASCENDENTEMENTE POR PRECIO DE PRODUCTO:
`CUSTOMER.ORDERBY("DESC(CUSTOMER)", "ASC(PRICE)").SHOW`
- ORDENAR DESCENDENTEMENTE POR CUSTOMER Y ASCENDENTEMENTE POR PRECIO DE PRODUCTO LAS COMPRAS QUE SE HAN HECHO DESPUÉS DEL DÍA 16/09/2018:
`CUSTOMER.WHERE(" DATE > ' 16 / 09 / 2018 ' ".ORDERBY("DESC(CUSTOMER)", "ASC(PRICE)").SHOW`
- ¿CUANTAS COMPRAS TENEMOS POR CADA CLIENTE?
`CUSTOMERS.GROUPBY("CUSTOMER").COUNT() .SHOW`
`.COUNT PORQUE CADA LÍNEA ES UNA COMPRA`
- ¿CUÁL ES EL NUMERO TOTAL DE PRODUCTOS QUE HA COMPRADO CADA CLIENTE?
`CUSTOMERS.GROUPBU("CUSTOMER").AGG(EXPR("SUM(QUANTITY)")).SHOW`
`SE APLICA AGREGACIÓN PORQUE QUEREMOS SABER LA SUMA DE QUANTITY POR CADA CUSTOMER`
- POR CADA CLIENTE CALCULA LA MEDIA DE LA CANTIDAD Y LA DESVIACION TIPICA POBLACIONAL DE CANTIDAD:
`CUSTOMERS.GROUPBY("CUSTOMER").AGG(EXPR("AVG(QUANTITY)"), EXPR("STDDEV_POP(QUANTITY)")).SHOW`
- POR CADA CLIENTE LA MEDIA DE CANTIDAD QUE HA COMPRADO Y EL PRECIO MÁXIMO:
`CUSTOMERS.GROUPBY("CUSTOMER").AGG(EXPR("AVG(QUANTITY)"), EXPR("MAX(PRICE)")).SHOW`
- UNIR DOS DATASETS:
`JOINED = PRODUCTS.JOIN(STOCK, PRODUCTS["ID"] == STOCK["ID"], "INNER")`
`JOINED.SHOW()`

PREGUNTAS DEL LABORATORIO

- ¿CUANTOS ELEMENTOS PODEMOS ENCONTRAR?

```
CUSTOMER.SELECTEXPR( "COUNT( CUSTOMER ) " ).SHOW
```

- ¿CUANTOS CUSTOMERS UNICOS TENEMOS?

```
CUSTOMERS.SELECTEXPR( "COUNT( DISTINCT( CUSTOMER ) ) " ).SHOW
```

- ¿CUANTOS PRODUCTOS HAN SIDO COMPRADOS POR CADA CUSTOMER?

```
CUSTOMERS.GROUPBY( "CUSTOMER" ) .AGG( "SUM( QUANTITY ) AS TOTAL" ) .SHOW
```

- ORDENAR LOS CUSTOMERS POR CANTIDAD

```
CUSTOMERS.GROUPBY( "CUSTOMER" ) .AGG( EXPR( "SUM( QUANTITY ) AS TOTAL" ) ) .ORDERBY( "TOTAL" ) .SHOW( )
```

- ¿CUANTAS VECES EL CUSTOMER CON ID 100 HA COMPRADO MÁS DE 5 ITEMS.

```
CUSTOMERS.WHERE( EXPR( "CUSTOMER = 100" ) ) .WHERE( EXPR( "QUANTITY >5" ) ) .COUNT( )
```

- CUALES HAN SIDO LOS PRODUCTOS COMPRADOS POR EL CUSTOMER CON EL NÚMERO MÁS GRANDE DE TRANSACCIONES? NOS INTERESA EL CUSRTOMER QUE HA COMPRADO MAS VECES.

```
MAX_CUSTOMER = CUSTOMERS.GROUPBY( "CUSTOMER" ) .COUNT( ) .ORDERBY( DESC( "COUNT" ) ) .FIRST( ) .CUSTOMER
```

DATAFRAME NO TIENE FIST PERO GROUPBY SI, ES UN VECTOR INDEXADO POR LAS COLUMNAS A LAS QUE PUEDES ACCEDER HACIENDO .COLUMN

```
CUSTOMERS.WHERE( EXPR( "CUSTOMER = " + STR( MAX_CUSTOMER ) ) ) .SELECT( "PRODUCT" ) .DISTINCT( ) .SHOW( )
```

ES IMPORTANTE HACER STR(MAX_CUSTOMER) PARA QUE LO TOME COMO UNA VARIABLE DEFINIDA. EL DISTINCT NO SE PUEDE APLICAR SIN UNA FUNCIÓN: POR EJEMPLO COUNT(DISTINCT(...)) SI SERIA POSIBLE PERO SI SOLO NECESITAMOS LOS VALORES ÚNICOS DE UNA COLUMNA HAY QUE SELECCIONAR LA COLUMNA SELECT(...) .DISTINCT()

- CUANTAS REGIONES SANITARIAS DIFERENTES HAY Y CUANTOS MEDICAMENTOS DIFERENTES SE HAN DISPENSADO?

```
RECETAS2.SELECT( "RSANITARIA" ) .DISTINCT( ) , COUNT( )
```

```
RECETAS2.SELECT( "MEDICAMENT" ) .DISTINCT( ) .COUNT( )
```

- CUANTAS RECETAS SE HAN DISPENSADO DE CADA MEDICAMENTO?

```
RECETAS2.GROUPBY( "MEDICAMENT" ) .AGG( EXPR( "SUM( NRECEPTES ) AS NUMERO_DE_RECETAS" ) ) .SHOW( )
```

SQL

- MUESTRA LOS DATOS DE MUJERES

```
SPARK.SQL( "SELECT * FROM MY_TABLE WHERE SEXE = 'DONA' " ) .SHOW
```

- ¿CUANTOS MEDICAMENTOS DIFERENTES SE EXPEDIERON EN 2022, CUANTAS RECETAS Y QUÉ COSTE TUVIERON?

```
SPARK.SQL(
```

```
"
```

```
SELECT MEDICAMENT, SUM( CAST( NRECETAS AS INT ) ) AS NUMERO_RECETAS, SUM( CAST( COST AS INT ) ) AS COST
```

```
FROM MYTABLE
```

```
GROUP BY MEDICAMENT
```

```
ORDER BY NUMERO_RECETAS DESC
```

```
"
```

IMPORTANTE HACER EL SUM PORQUE POR CADA MEDICAMENTO TENEMOS COMO MÍNIMO DOS FILAS:

MEDICAMENTO_A MUJER NRECETAS COSTE

MEDICAMENTO_A HOMBRE NRECEAS COSTE

- ¿CUÁL ES EL MEDIAMENTO MAS RECETADO EN HOMBRE Y MUJER Y EN QUE REGIÓN SANITARIA?

```
RECETAS2.CREATEORREPLACETEMPVIEW( "MY_TABLE" )
```

```
CONSULTA = ""
```

```
SELECT MEDICAMENT, RSANITARIA, SEXE, NRECEPTES AS NUMERO_RECETAS
```

```
FROM MY_TABLE
```

```
WHERE SEXE == 'DONA' OR SEXE == 'HOME'
```

```
GROUP BY MEDICAMENT, RSANITARIA, SEXE, NRECEPTES
```

```
ORDER BY SEXE
```

```
""
```

```
RESULT = SPARK.SQL( CONSULTA )
```

```
WINDOW_SPEC = WINDOW.PARTITIONBY( "SEX" ) .ORDERBY( F.DISC( "NUMERO_RECETAS" ) )
```

```
RESULT.WITHCOLUMN( "RANK", F.RANK( ) .OVER( WINDOW_SPEC ) ) .WHERE( EXPR( "RANK=1" ) )
```

SE HACE UNA CONSULTA COGIENDO EL MEDICAMENTO, LA REGIÓN SANITARIA EL SEXO Y EL NUMERO DE RECETAS AGRUPANDO POR TODOS ELLOS FILTRANDO SOLO POR HOMBRE Y MUJER Y ORDENANDO POR SEXO, DE FORMA QUE TENAMOS TODAS LAS MUJERES ARRIBA Y TODOS LOS HOMBRES ABAJO.

SE CREA UNA WINDOW QUE PARTICIONA POR SEXO (POR ESO ORDENAR POR SEXO→MÁS CLARIDAD) Y QUE ORDENA DENTRO DE LAS PARTICIONES(SEXO) POR NUMERO DE RECETAS DE FORMADISCENDENTE. DE ESTA FORMA TENDREMOS EL MEDICAMENTO MÁS PEDIDO PARA MUJERES ARRIBA DE Y EL MEDICAMENTO MAS PEDIDO POR HOMBRES ARRIBA DE LA SECCIÓN DE HOMBRES.

CREAMOS UNA NUEVA COLUMNA LLAMADA RANK QUE SE BASA EN EL ORDEN DEFINIDO EN LA WINDOW. SE APLICA UNA FUNCIÓN DE RANKING QUE ASIGNARÁ DE CADA GRUPO (SEXO) EL VALOR 1 AL QUE TENGA EL VALOR MAS GRANDE DE NUMERO RECETAS. POR ÚLTIMO SOLO NOS QUEDAMOS CON ESAS FILAS QUE TIENEN EN MAYOR NÚMERO DE RECETAS.

- ¿CUÁL ES EL MEDICAMENTO MENOS RECETADO EN HOMBRE Y MUJER Y EN QUE REGIÓN SANITARIA?

```
RECETAS2.CREATEORREPLACETEMPVIEW( "MY_TABLE")
```

```
CONSULTA = """
```

```
SELECT MEDICAMENT, RSANITARIA, SEXE, NRECEPTES AS NUMERO_RECETAS
```

```
FROM MY_TABLE
```

```
WHERE SEXE == 'DONA' OR SEXE == 'HOME'
```

```
GROUP BY MEDICAMENT, RSANITARIA, SEXE, NRECEPTES
```

```
ORDER BY SEXE
```

```
"""
```

```
RESULT = SPARK.SQL( CONSULTA)
```

```
WINDOW_SPEC = WINDOW.PARTITIONBY( "SEXE").ORDERBY( F.DESC( "NUMERO_RECETAS") )
```

```
RESULT.WITHCOLUMN( "RANK", F.RANK( ) .OVER( WINDOW_SPEC ) ).WHERE( EXPR( "RANK=1" ) )
```

- ¿CUÁL ES EL MEDICAMENTO MÁS CARO?

EL MEDICAMENTO MÁS CARO DE MEDIA

```
RECETAS2.CREATEORREPLACETEMPVIEW( "MY_TABLE")
```

```
CONSULTA = """
```

```
SELECT MEDICAMENT, AVG( IMPORT ) AS COSTE_MEDIO_MEDICAMENTO
```

```
FROM MY_TABLE
```

```
GROUP BY MEDICAMENT
```

```
ORDER BY COSTE_MEDIO_MEDICAMENTO DESC
```

```
LIMIT 1
```

```
"""
```

```
SPARK.SQL( CONSULTA)
```

EL MEDICAMENTO MÁS CARO DE TODOS Y EN QUÉ REGIÓN ES TAN CARO

```
RECETAS2.CREATEORREPLACETEMPVIEW( "MY_TABLE")
```

```
CONSULTA = """
```

```
SELECT MEDICAMENT, RSANITARIA
```

```
FROM MY_TABLE
```

```
ORDER BY IMPORT DESC
```

```
LIMIT 1
```

```
"""
```

```
SPARK.SQL( CONSULTA ).SHOW( )
```

- CONSIDERANDO LOS 10 MEDICAMENTOS MAS RECETADOS DURANTE 2022 INDICAR EL COSTE DE ESTE POR REGIÓN SANITARIA.

```
RECETAS2.CREATEORREPLACETEMPVIEW( "MY_TABLE")
```

```
CONSULTA_SUB = """ (
```

```
SELECT MEDICAMENT
```

```
FROM MY_TABLE
```

```
WHERE MEDICAMENT != 'SENSE ESPECIFICAR'
```

```
GROUP BY MEDICAMENT
```

```
ORDER BY SUM( CAST( NRECEPTES AS INT ) ) DESC
```

```
LIMIT 10) """
```

```
CONSULTA = """
```

```
SELECT RSANITARIA, MEDICAMENT, AVG( CAST( IMPORT AS INT ) ) AS COSTE_MEDIO_MEDICAMENTO
```

```
FROM MY_TABLE
```

```
WHERE MEDICAMENT IN""" + CONSULTA_SUB + """
```

```
GROUP BY RSANITARIA, MEDICAMENT
```

```
ORDER BY RSANITARIA, MEDICAMENT
```

```
"""
```

```
SPARK.SQL( CONSULTA ).SHOW( )
```