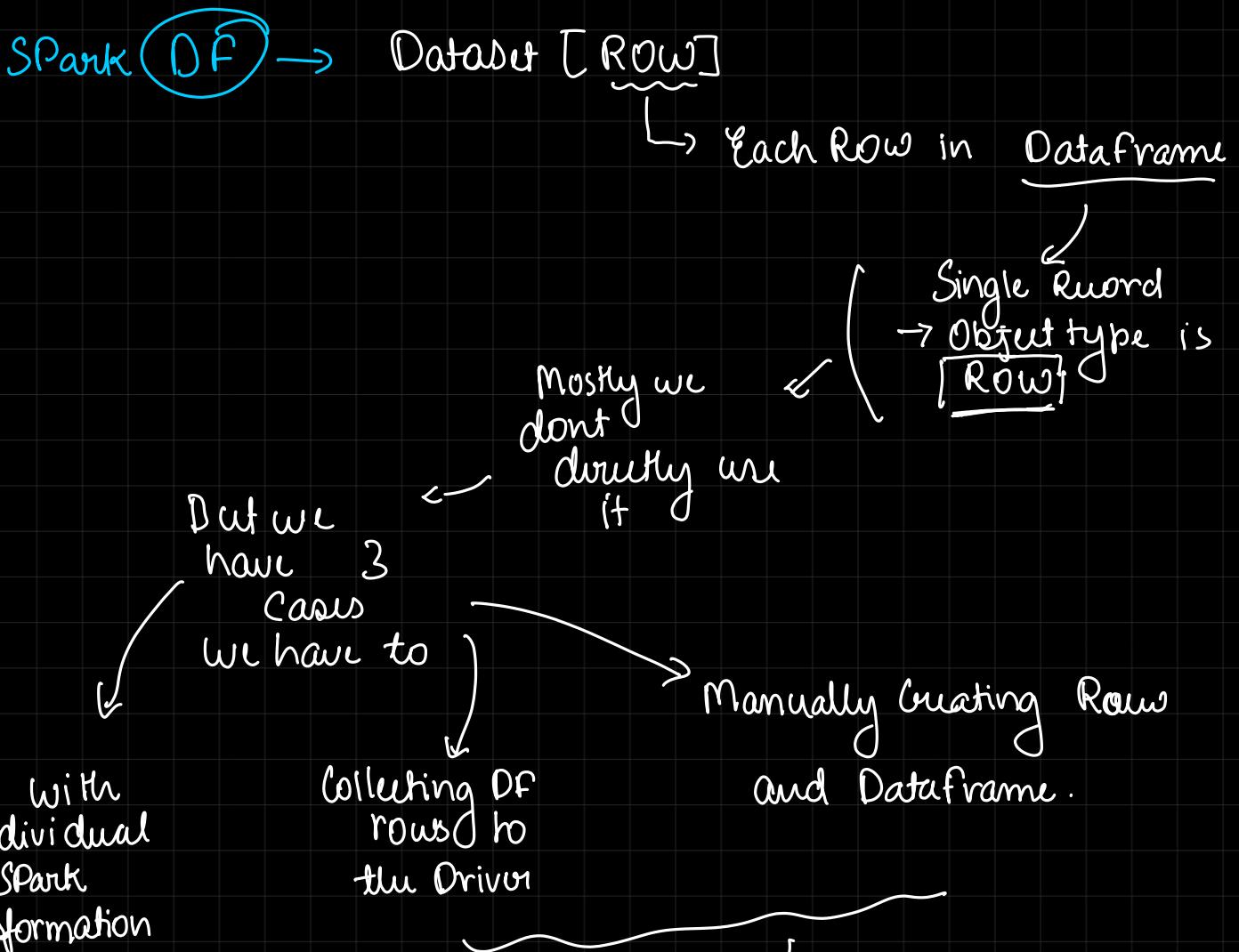


SPark DataSet / Data Transform →

- Combining Data Frame
- Aggregation and Summarizing
- Applying function and Built in transform
- Built-in and Column-level Transformation
- Creating using ↳ UDFs



Mostly used in
Developing
or
Unit testing.

```

my_schema = StructType([
    StructField("ID", StringType()),
    StructField("EventDate", StringType())
])

my_rows = [Row("123", "04/05/2020"), Row("124", "4/5/2020"), Row("125", "04/5/2020"), Row("126", "4/05/2020")]
my_rdd = spark.sparkContext.parallelize(my_rows, 2)
my_df = spark.createDataFrame(my_rdd, my_schema)

my_df.printSchema()
my_df.show()
new_df = to_date_df(my_df, "M/d/y", "EventDate")
new_df.printSchema()
new_df.show()

```

Spark → Comes with Bunch of Transformation →

agg (*exprs)

coV (col1, col2)

crossTab (cols, col2)

cube (*cols)

filter (condition)

groupBy (*cols)

join (other, on=None, how=None)

orderBy (*cols, *kwargs)

replace (to_replace, Value, subset)

rollup (*cols)

→ when DF is
not columnar
Struct

first we
need to
Create

columnar
format

Then perform
Transformation
function

for this we
need Row

```

root
|-- ID: string (nullable = true)
|-- EventDate: string (nullable = true)

+---+-----+
| ID| EventDate|
+---+-----+
|123|04/05/2020|
|124| 4/5/2020|
|125| 04/5/2020|
|126| 4/05/2020|
+---+-----+

root
|-- ID: string (nullable = true)
|-- EventDate: date (nullable = true)

+---+-----+
| ID| EventDate|
+---+-----+
|123|2020-04-05|

```

For Example we have
Unstructured Data



We Cannot have Column as
its unstructured
DF

↳ this will String.

```
log_regex = r'^(\S+) (\S+) (\S+) \[(\w:/.+\s[+-]\d{4})\] "( \S+ ) (\S+) (\S+)" (\d{3}) (\S+) "( \S+ )" ("[^"]*" )'

log_df = file_df.select(regexp_extract(str: 'value', log_regex, idx: 1).alias('ip'),
                        regexp_extract(str: 'value', log_regex, idx: 4).alias('date'),
                        regexp_extract(str: 'value', log_regex, idx: 6).alias('request'),
                        regexp_extract(str: 'value', log_regex, idx: 10).alias('referrer'))
```

```
log_df.printSchema()
log_df \
    .withColumn(colName: "referrer", substring_index(str: "referrer", delim: "/", count: 3)) \
    .groupBy("referrer") \
    .count() \
    .show(n: 100, truncate=False)
```

```
root
|-- ip: string (nullable = true)
|-- date: string (nullable = true)
|-- request: string (nullable = true)
|-- referrer: string (nullable = true)

+-----+-----+
|referrer|count|
+-----+-----+
|http://ijavascript.cn|1   |
|http://www.google.co.tz|1   |
|http://www.google.ca|6   |
```

Working with DF Columns ↗

```
%fs head /databricks-datasets/airlines/part-00000  
[Truncated to first 65536 bytes]  
Year,Month,DayofMonth,DayofWeek,DepTime,CRSDepTime,ArrTime,CRSArrTime,UniqueCarrier,FlightNum,TailNum,ActualElapsedTime,Dest,Distance,TaxiIn,TaxiOut,Cancelled,CancellationCode,Diverted,CarrierDelay,WeatherDelay,NASDelay,SecurityDelay,  
1987,10,14,3,741,738,912,849,PS,1451,NA,91,79,NA,23,11,SAN,SFO,447,NA,NA,0,NA,0,NA,NA,NA,NA,NA,YES,YES  
1987,10,15,4,729,730,903,849,PS,1451,NA,94,79,NA,14,-1,SAN,SFO,447,NA,NA,0,NA,0,NA,NA,NA,NA,NA,YES,NO
```

↓
After this Creating Dataframe.

```
airlinedf = spark.read \  
.format("csv") \  
.option("header","true") \  
.option('inferSchema','true') \  
.option('samplingRatio','0.0001') \  
.load('/databricks-datasets/airlines/part-00000')
```

Infering Schema with
small Sampling ratio

helps to run quickly.

we now have
DF to show
Learn Column
expression

Spark DataFrame Columns

↳ object of types Column

— we Cannot Manipulate
Independently

— get get by Spark Transformation
only.

Thru are

2 ways to
Do transform

Column
String

Column
Object

Simplest way
to access
Column.

It is Complex

```

Just now (1s)

airlinedf.select('Origin', 'Dest', 'Distance').show(10)

(1) Spark Jobs
+---+---+---+
|Origin|Dest|Distance|
+---+---+---+
| SAN| SFO| 447|
+---+---+---+

```

As there is many ways
to create Column object.

Most common Method
to use

col()

Column()

```
from pyspark.sql.functions import *
```

```
airlinedf.select(column("Origin"), col("Dest"), airlinesDF.Distance).show(10)
```

↓
3 types 3 way.

How to Create Column Expression :→

SparkDF → 2 way to Create Column Expression

String or
SQL

Column Object
Expression

For Example
[year, Month, Date]

3 column n ↴

I want to combine
this
make as
Date Format

```

Last execution failed
1 airlinedf.select('Origin', 'Dest', 'Distance', 'to_date(concat(year,month,dayofmonth), 'yyyyMMd') as FlightDate').show(10)
0 > AnalysisException: [UNRESOLVED_COLUMN_WITH_SUGGESTION] A column or function parameter with name 'to_date(concat(concat(Year,Month,DayofMonth),'yyyyMMd')) as FlightDate' cannot be resolved. Did you mean one of the following? [DayofMonth, LateArrivalCraftDelay, CancellationCode, WeatherDelay, ActualElapsedTime]; Project [Origin#104, Dest#105, Distance#106, 'to_date(concat(concat(Year,Month,DayofMonth), 'yyyyMMd')) as FlightDate']...

```

This will through Error

To_date (concat(year,month,DayofMonth))

As Select () → only use

Column String
Column Object.

for this we need `Expr()`

Expr(To_date(concat(Year,Month,Day) , yyyyMMdd))

This will Return
Column Object

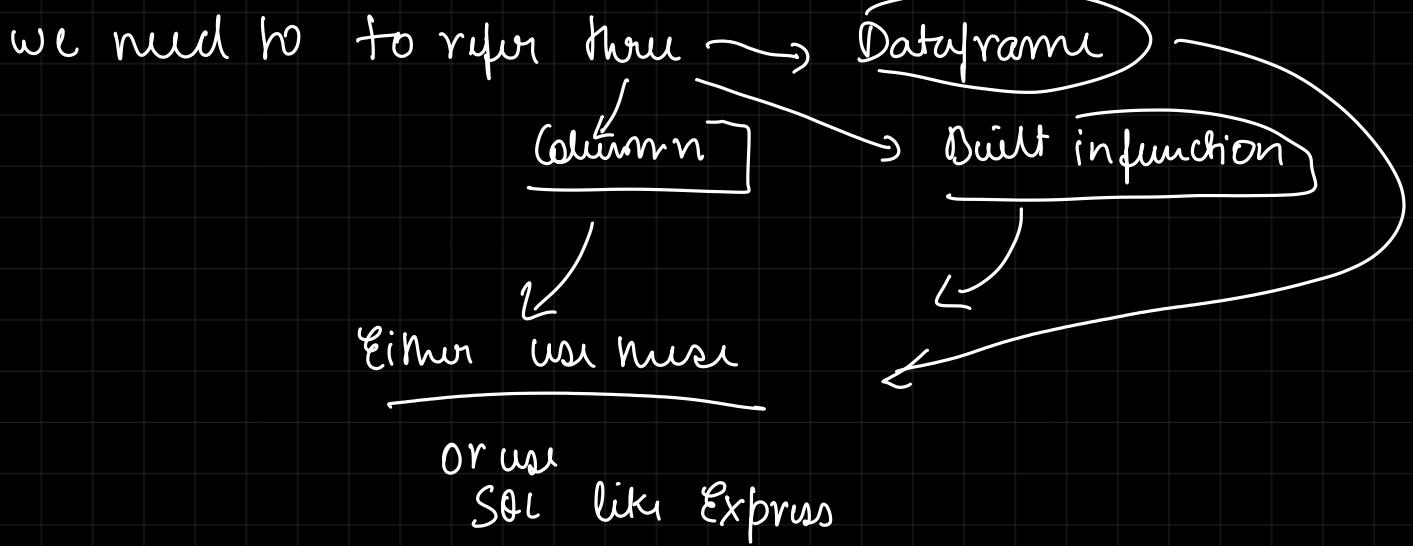
This is all
String or SQL Expression

```
Just now (2s)
airlinedf.select("Origin","Dest","Distance",
|           |   expr("to_date(concat(Year,Month,DayofMonth), 'yyyyMMdd') as FlightDate" ).show(10)
▶ (1) Spark Jobs
+---+---+-----+
|Origin|Dest|Distance|FlightDate|
+---+---+-----+
| SAN| SFO| 447|1987-10-14|
| SAN| SFO| 447|1987-10-15|
| SAN| SFO| 447|1987-10-17|
| SAN| SFO| 447|1987-10-18|
| SAN| SFO| 447|1987-10-19|
| SAN| SFO| 447|1987-10-21|
| SAN| SFO| 447|1987-10-22|
| SAN| SFO| 447|1987-10-23|
| SAN| SFO| 447|1987-10-24|
```

```
Just now (2s)
airlinedf.select("Origin","Dest","Distance",
|           |   expr("to_date(concat(Year,Month,DayofMonth), 'yyyyMMdd') as FlightDate" ).show(10)
▶ (1) Spark Jobs
+---+---+-----+
|Origin|Dest|Distance|FlightDate|
+---+---+-----+
| SAN| SFO| 447|1987-10-14|
| SAN| SFO| 447|1987-10-15|
| SAN| SFO| 447|1987-10-17|
| SAN| SFO| 447|1987-10-18|
| SAN| SFO| 447|1987-10-19|
| SAN| SFO| 447|1987-10-21|
| SAN| SFO| 447|1987-10-22|
| SAN| SFO| 447|1987-10-23|
| SAN| SFO| 447|1987-10-24|
```

Column Object]

```
airlinedf.select("Origin","Dest","Distance",
|           |   to_date(concat("Year", "Month", "DayofMonth"), 'yyyyMMdd').alias("FlightDate")
|           | ).show(10)
▶ (1) Spark Jobs
+---+---+-----+
|Origin|Dest|Distance|FlightDate|
+---+---+-----+
| SAN| SFO| 447|1987-10-14|
| SAN| SFO| 447|1987-10-15|
| SAN| SFO| 447|1987-10-17|
| SAN| SFO| 447|1987-10-18|
| SAN| SFO| 447|1987-10-19|
| SAN| SFO| 447|1987-10-21|
| SAN| SFO| 447|1987-10-22|
| SAN| SFO| 447|1987-10-23|
| SAN| SFO| 447|1987-10-24|
| SAN| SFO| 447|1987-10-25|
```



Create User Define function ↗

```

if __name__ == "__main__":
    spark = SparkSession.builder \
        .appName('UDFFunction') \
        .master('local[2]') \
        .getOrCreate()

    logger = Log4J(spark)

    survey_df = spark.read \
        .option(key: "header", value: "true") \
        .option(key: "inferSchema", value: "true") \
        .csv('D:\\SPARK\\pythonProject\\SparkSQL\\data\\survey.csv')

    survey_df.show(10)
  
```

Gender
Female
Male
Male
Male
Male
Female

↓ How to make Standard

```

def parse_gender(gender):
    female_pattern = r"^\f$|f.m|w.m"
    male_pattern = r"^\m$|ma|m.l"

    if re.search(female_pattern, gender.lower()):
        return "Female"
    elif re.search(male_pattern, gender.lower()):
        return "Male"
    else:
        return "Unknown"
  
```

for this
create Custom
function.

- } 1 - Male
2 - Female
3 - Unknown

↳ How to use them ⇒ we have to create
expression and use it

Column object
Expression

} We have 2 approach → to
develop the Expression.

String Expression.

Column Object Transformation

With Column() : → Allows to transform with single column.

Without impacting other DF

```
survey_df2 = survey_df.withColumn(colName: "Gender", parse_gender("Gender"))
```

I need to register my custom function

new
UDF

UDF()

to use if

UDF (Parse gender , StringType())

Column Object Expression

```
parse_gender_udf = udf(parse_gender, StringType())
survey_df2 = survey_df.withColumn(colName: "Gender", parse_gender_udf("Gender"))
```

default is String type
and it optional

This register as Dataframe UDF()

UDF

SQl Expression

This will register in Catalog.

Act as SQl function.

It will not register in Catalog

It will create → UDF

Serialized in Executor.

```

data_list = [
    ("Ravi", "28", "1", "2020"),
    ("avi", "26", "1", "81"),
    ("Sony", "27", "12", "6"),
    ("lg", "8", "12", "220"),
]

```

```

raw_df = spark.createDataFrame(data_list).toDF("name", "day", "month", "year")
raw_df.printSchema()

```

This issue we can fix it by "digit problem"

To add Unique Identifier.

It will Unique across all Identifier

If the Column don't Exist It will add the new column.

```

df1 = raw_df.withColumn("id", monotonically_increasing_id())
df1.show()
✓ 2.9s

+---+---+---+---+
|name|day|month|year|
+---+---+---+---+
|Ravi| 28| 1| 2020| |
|Sony| 27| 12| 6|
| avi| 26| 1| 81| 8589934592|
| lg| 8| 12| 220| 17179869184|
+---+---+---+---+

```

```

df2 = df1.withColumn("year" , expr(
    ...
    case when year < 21 then year + 2000
    when year < 100 then year + 1900
    else year
    end
    ...
))

✓ 0.1s

```

```

df2.show()
✓ 1.3s

+---+---+---+---+
|name|day|month| year| id|
+---+---+---+---+
|Ravi| 28| 1| 2020| 0|
|Sony| 27| 12| 2006.0| 1|
| avi| 26| 1| 1981.0| 8589934592|
| lg| 8| 12| 2020.0| 17179869184|
+---+---+---+---+

```

Putting .0 decimal → Now to fix this

This can be done by

By Casting?

Inline Cast

Change the schema

```

df3 = df1.withColumn("year" , expr(
    ...
    case when year < 21 then cast(year as int) + 2000
    when year < 100 then cast(year as int) + 1900
    else cast(year as int)
    end
    ...
))

df3.show()

```

✓ 1.3s

name	day	month	year	id
Ravi	28	1	2020	0
Sony	27	12	2006	1
avi	26	1	1981	8589934592
lg	8	12	2020	17179869184

```

df4 = df4.withColumn("year", col("year").cast(IntegerType()))
df4.show()
df4.printSchema()

✓ 1.2s

```

name	day	month	year	id
Ravi	28	1	2020	0
Sony	27	12	2006	1
avi	26	1	1981	8589934592
lg	8	12	2020	17179869184

```

root
|-- name: string (nullable = true)
|-- day: string (nullable = true)
|-- month: string (nullable = true)
|-- year: integer (nullable = true)
|-- id: long (nullable = false)

```

Adding and Removing Columns and duplicate

```

df8 = df7.withColumn('dob', expr("to_date(concat(day,'/',month,'/',year),'d/M/y')"))
df8.show()

```

✓ 1.2s

name	day	month	year	id	dob
Ravi	28	5	2020	0	2020-05-28
milton	2	4	2020	1	2020-04-02
avi	26	7	1981	8589934592	1981-07-26
lg	8	12	2020	8589934593	2020-12-08
DELL	17	12	2020	8589934594	2020-12-17
Sony	27	1	2006	17179869184	2006-01-27
Sam	18	1	2020	17179869185	2020-01-18

```

df9= df7.withColumn('dob',to_date(expr("concat(day,'/',month,'/', year)"),'d/M/y'))
df9.show()

```

✓ 1.3s

name	day	month	year	id	dob
Ravi	28	5	2020	0	2020-05-28
milton	2	4	2020	1	2020-04-02
avi	26	7	1981	8589934592	1981-07-26
lg	8	12	2020	8589934593	2020-12-08
DELL	17	12	2020	8589934594	2020-12-17
Sony	27	1	2006	17179869184	2006-01-27
Sam	18	1	2020	17179869185	2020-01-18

Drop Column

```

df8 = df8.drop('day','month','year')
df8.show()

```

✓ 1.2s

name	id	dob
Ravi	0	2020-05-28
DELL	1	2020-12-17
Sony	2	2006-01-27
avi	8589934592	1981-07-26
lg	8589934593	2020-12-08
Sony	8589934594	2006-01-27
Sony	17179869184	2006-01-27
Sam	17179869185	2020-01-18
milton	17179869186	2020-04-02

Drop Duplicate

```

df8 = df8.dropDuplicates(['name','dob'])
df8.show()

```

✓ 1.5s

name	id	dob
DELL	1	2020-12-17
Sony	2	2006-01-27
Ravi	0	2020-05-28
avi	8589934592	1981-07-26
lg	8589934593	2020-12-08
Sam	17179869185	2020-01-18
milton	17179869186	2020-04-02

Short Desc

```

df8 = df8.sort(expr('dob desc'))
df8.show()

```

✓ 1.4s

name	id	dob
avi	8589934592	1981-07-26
Sony	2	2006-01-27
Sam	17179869185	2020-01-18
milton	17179869186	2020-04-02
Ravi	0	2020-05-28
lg	8589934593	2020-12-08
DELL	1	2020-12-17

Aggregation



All aggregation → in Spark
are implemented
by built-in
function

It's classified into (3)

Simple Aggregation

Windowing
Aggregation

Grouping
Aggregation

- head()
- lag()
- rank()
- dense_rank()
- cumu_dist()

- avg()
- count()
- max()
- min()
- sum()

```
invoice_df.select(count("*").alias("total_invoice"),  
                  sum("Quantity").alias("TotalQuantity"),  
                  avg("UnitPrice").alias("Avg_UnitPrice"),  
                  count_distinct("InvoiceNo").alias("CountDistinct"))  
                .show()
```

We can use them as
Column Object
Expression.

Result

total_invoice	TotalQuantity	Avg_UnitPrice	CountDistinct
541909	5176450	4.611113626088481	25900

We can also use a SQL Expression

instead of Select → Select Expr ()

```
invoice_df.selectExpr(  
    *expr: "count(1) as count 1",  
    "count(StockCode) as `count field`",  
    "sum(Quantity) as TotalQuantity",  
    "avg(UnitPrice) as AvgPrice"  
)  
  .show()
```

Count (*) = Count(1)

this will not count
null Value

>>> SQL Expression <<<

count 1	count field	TotalQuantity	AvgPrice
541909	541908	5176450	4.61111362608849

Groupping Aggregation ➔

By using SQL Query ↴

```
invoice_df.createOrReplaceTempView("sales")
summary_Sql = spark.sql(
    """
    SELECT Country, InvoiceNo,
           SUM(Quantity),
           ROUND(SUM(Quantity*UnitPrice),2) as InvoiceValue
      FROM sales
     GROUP BY 1,2
    """
)
```

Country	InvoiceNo	sum(Quantity)	InvoiceValue
United Kingdom	536446	329	440.89
United Kingdom	536508	216	155.52
United Kingdom	537018	-3	0.0
United Kingdom	537401	-24	0.0
United Kingdom	537811	74	268.86
United Kingdom	C537824	-2	-14.91
United Kingdom	538895	370	247.38
United Kingdom	E10/EZ1	741	700.52

We can do the same for grouping → Mechanism.
in 

DataFrame Grouping

```
summary_df = invoice_df.groupBy("Country", "InvoiceNo") \
    .agg(func.sum("Quantity").alias("Total_quantity"),
         func.round(func.sum(func.expr("Quantity*UnitPrice")),2).alias("InvoiceValue"))
summary_df.show()
```

0.8s

Country	InvoiceNo	Total_quantity	InvoiceValue
United Kingdom	536446	329	440.89
United Kingdom	536508	216	155.52
United Kingdom	537018	-3	0.0
United Kingdom	537401	-24	0.0
United Kingdom	537811	74	268.86


```
invoi_up = invoice_df.groupBy('Country',
    func.weekofyear(func.timestamp(col='InvoiceDate',format='dd-MM-yyyy H.mm')).alias("WeekNumber"),
    .agg(func.count_distinct('InvoiceNo').alias('NumInvoices'),
        func.sum('Quantity').alias('TotalQuantity'),
        func.round(func.sum(func.expr("Quantity*UnitPrice")),scale=2).alias('Invoice_Value'))
invoi_up.show()
```

1.5s

Country	WeekNumber	NumInvoices	TotalQuantity	Invoice_Value
EIRE	17	1	163	362.9
Norway	47	1	594	2016.78
Channel Islands	46	1	78	211.63
France	10	4	471	998.49

Example 2 ↗

Country	WeekNumber	NumInvoices	TotalQuantity	InvoiceValue
EIRE	48	7	2822	3147.23
EIRE	49	5	1280	3284.11
EIRE	50	5	1184	2321.78
EIRE	51	5	95	276.84
France	48	41	1299	2808.16
France	49	9	2303	4527.01
France	50	6	529	537.32
France	51	5	847	1702.87

1st → Partitioning by Country first

2nd → Order by week number.

Dataframe Joins ↗

Is all about → bringing together 2DF

Combine 2DF together by 2 Mings

Join Condition / Expression

Join Type:

Join type = "inner"

Join_Expr = Order-df.Prod.id ==

Product-df.prod-id

Order df .Join (Product_renamed_df , Join Expr , "inner")

```
order_df.join(product_renamed_df,join_expr,"left") \
    .drop(product_renamed_df.prod_id) \
    .select("order_id","prod_id","prod_name","unit_price","list_price","qty") \
    .sort("order_id") \
    .show()
```

✓ 1.7s

order_id	prod_id	prod_name	unit_price	list_price	qty
01	07 32 GB Flash Storage	320	320	2	
01	02 Optical Mouse	350	350	1	
01	04 Wireless Keyboard	580	580	1	
02	03 Wireless Mouse	450	450	1	
02	06 16 GB Flash Storage	220	240	1	
03	01 Scroll Mouse	195	250	1	
04	09 NULL	270	NULL	3	
04	08 64 GB Flash Storage	410	430	2	
05	02 Optical Mouse	350	350	1	

```
order_df.join(product_renamed_df,join_expr,"left") \
    .drop(product_renamed_df.prod_id) \
    .select("order_id","prod_id","prod_name","unit_price","list_price","qty") \
    .withColumn("prod_name", func.expr("coalesce(prod_name,prod_id)")) \
    .withColumn("list_price", func.expr("coalesce(list_price,unit_price)")) \
    .show()
```

✓ 1.7s

order_id	prod_id	prod_name	unit_price	list_price	qty
01	07 32 GB Flash Storage	320	320	2	
02	03 Wireless Mouse	450	450	1	
01	02 Optical Mouse	350	350	1	
01	04 Wireless Keyboard	580	580	1	
03	01 Scroll Mouse	195	250	1	
04	09 NULL	270	270	3	
04	08 64 GB Flash Storage	410	430	2	
05	02 Optical Mouse	350	350	1	
02	06 16 GB Flash Storage	220	240	1	

Spark offers two approaches to Join

Shuffle Join

most commonly used in
Implementation Done by
MapReduce kind of
Hadoop.

ID	FL_DATE	OP_CARRIER	OP_CARRIER_FL_NUM	ORIGIN_CITY_NAME	DEST_CITY_NAME
00	1/1/2000	DL	1451	BOS	Boston, MA
01	1/1/2000	DL	1479	BOS	Boston, MA
02	1/1/2000	DL	1857	BOS	Boston, MA
03	1/1/2000	DL	1997	BOS	Boston, MA
04	1/1/2000	DL	2060	BOS	Boston, MA
05	1/1/2000	US	2419	BOS	Boston, MA
06	1/1/2000	US	2862	BOS	Boston, MA
07	1/1/2000	DL	346	BTR	Baton Rouge, LA
08	1/1/2000	DL	412	BTR	Baton Rouge, LA
09	1/1/2000	DL	299	BUF	Buffalo, NY
10	1/1/2000	DL	495	BUF	Buffalo, NY

ID	CRS_DEP_TIME	CRS_DEP_TIME	WHEELS_ON_TAXI_IN	CRS_ARR_TIME	ARR_TIME	CANCELLED	DISTANCE
05	1115	1113	1343	5	1400	1340	0
01	1315	1311	1536	7	1559	1543	0
02	1415	1414	1642	9	1721	1651	0
03	1715	1720	1955	10	2013	2005	0
04	2015	2010	2230	10	2300	2240	0
05	0000	0000	0056	7	955	930	0
06	1440	1446	1713	8	1739	1717	0
07	1740	1744	1957	9	2008	2006	0
08	1345	1345	1552	9	1622	1601	0
09	1245	1245	1443	5	1455	1448	0
10	2035	2035	2226	9	2241	2235	0

ID	FL_DATE	OP_CARRIER	OP_CARRIER_FL_NUM	ORIGIN_CITY_NAME	DEST_CITY_NAME
11	1/1/2000	DL	677	BUF	Buffalo, NY
12	1/1/2000	DL	253	BWI	Baltimore, MD
13	1/1/2000	DL	1003	BWI	Baltimore, MD
14	1/1/2000	DL	1501	BWI	Baltimore, MD
15	1/1/2000	DL	1907	BWI	Baltimore, MD
16	1/1/2000	DL	2063	BWI	Baltimore, MD
17	1/1/2000	DL	1943	BWI	Baltimore, MD
18	1/1/2000	US	2952	BWI	Baltimore, MD
19	1/1/2000	US	2967	BWI	Baltimore, MD
20	1/1/2000	DL	540	CAE	Columbia, SC

ID	CRS_DEP_TIME	DEP_TIME	WHEELS_ON_TAXI_IN	CRS_ARR_TIME	ARR_TIME	CANCELLED	DISTANCE
11	710	740	925	947	947	0	576
12	2040	2235	7	2233	2232	0	576
13	1635	1838	2020	12	1832	2032	0
14	1430	1435	1673	12	1634	1635	0
15	530	530	716	4	723	720	0
16	1250				1449	1	576
17	1845	1855	2041	9	2046	2050	0
18	710	710	905			0	576
19	1700	1700	1845	6	1853	1851	0
20	955	952	1052	5	1104	1057	0

ID	FL_DATE	OP_CARRIER	OP_CARRIER_FL_NUM	ORIGIN_CITY_NAME	DEST_CITY_NAME
21	1/1/2000	DL	544	CAE	Columbia, SC
22	1/1/2000	DL	1289	CAE	Columbia, SC
23	1/1/2000	DL	5332	CAE	Columbia, SC
24	1/1/2000	DL	2315	CAE	Columbia, SC
25	1/1/2000	DL	1564	CAE	Columbia, SC
26	1/1/2000	DL	2030	CAE	Columbia, SC
27	1/1/2000	DL	393	CHS	Charleston, SC
28	1/1/2000	DL	423	CHS	Charleston, SC
29	1/1/2000	DL	717	CHS	Charleston, SC
30	1/1/2000	DL	773	CHS	Charleston, SC

ID	CRS_DEP_TIME	DEP_TIME	WHEELS_ON_TAXI_IN	CRS_ARR_TIME	ARR_TIME	CANCELLED	DISTANCE
21	2130	2125	2208	10	2232	2218	0
22	1515	1514	1619	9	1633	1628	0
23	1735	1737	1844	7	1852	1851	0
24	720	720	811	10	821	0	191
25	1330	1331	1931	2019	9	2047	2028
26	1340					1450	1
27	1005	1002	1100	7	1118	1107	0
28	600	556	654	17	707	711	0
29	1340	1340	1432	7	1454	1439	0
30	2135	2125	2219	5	2242	2224	0

3 Partition → distribute across different Executor

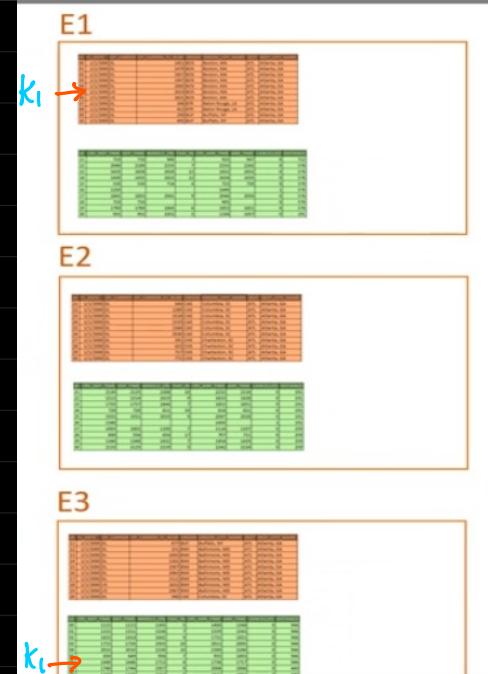
If $k_1 = \epsilon_L$ $k_1 = \epsilon_3$
 ↳ Now to Perform Join its not possible

Because keys are in different Executor

↳ It need to be in same Executor

So Join Can be Performed ↴

In 2 Stages



1st → Each Executor Map those records using Join key → Send it to the Exchange

(Map Phase)

Map Exchange

This Map Exchange be available to be picked up by [Spark] Framework

Spark Framework will pick this record

Send it to the Reduce Exchange

going to collect all the record for the same key.

For Example

→ All the record with same key → 0-1

will go to [Same Exchange.]

This Exchange Known as Shuffle Partition

we need to decide how many Shuffle Partition needed.

Shuffle Partition

=

Number of Nodes

So Each Partition handle
by one Executor
so all Partition Process Parallelly.

Now Data got transformed

from

Map Exchange

)

To Reduce Exchange

do Shuffle operation

↳ It can block network

↳ Slow down

Join operation

This Record
combines

And Create new

Data frame Partition

For this

we need to use

Sort Merge
Join

-algo.

```
if __name__ == "__main__":
```

```
spark = SparkSession.builder \  
.appName("ShuffleJoin") \  
.master("local[3]") \  
.getOrCreate()
```

3 nodes

```
flight_time_df1 = spark.read.json('data/d1')  
flight_time_df2 = spark.read.json('data/d2')
```

3 Partition

```
spark.conf.set("spark.sql.shuffle.partitions", 3)
```

```

join_expr = flight_time_df1.id == flight_time_df2.id
join_df = flight_time_df1.join(flight_time_df2, join_expr, "inner")
join_df.foreach(lambda f : None)

```

this is transformation

↳ It will Only activate when we take an action

Joining Dataframes

are of 2 type

→ Large to Large
 → Large to Small

Tiny → or fixed with single Executor

Large - Large

↳ filter Join → Perform Almost all aggregation

↳ filtering out necessary data.

global_online_sales			
order_id	city	product_code	store_code
S001	New York	ABX-512	abc.com
S002	London	GCY-213	abc.com

us_warehouse_events				
order_id	city	warehouse_code	event_type	event_date
S001	New York	NY516	order placed	15-Jan-2020
S001	New York	NY516	dispatched	16-Jan-2020

This not needed

so filter As we are Only looking for us warehousing Only

look for all Possible opportunities
to Reduce dataframe size

European Market detail will not help.

)
Then perform JOINS

→ look for → Number Shuffle Partitions
Number of Executors { }

helps with Parallelism of my Joins.

What maximum Possible Parallelism

Executor = 500
Shuffle Partition = 400
Unique Key = 200

Even you have
(500 Nodes)
(400 Shuffle)

If I have 200 Unique Key
then I will have Only 200

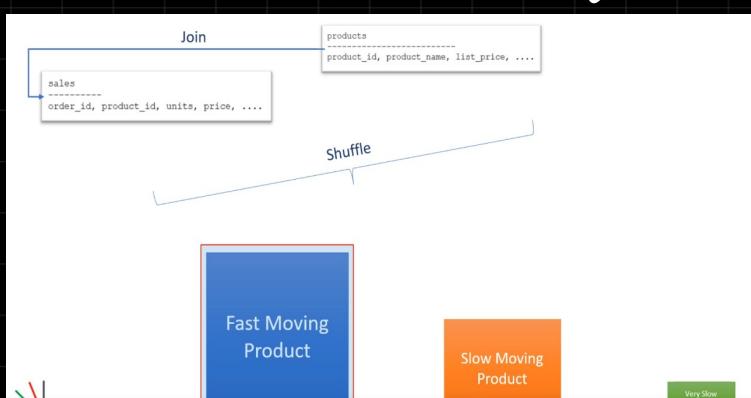
Shuffle Partition
other 200 will remain Blank.

It will make
Only 200.

→ So → In Unique Key → Limiting Scalability

need to increase

Join Keys

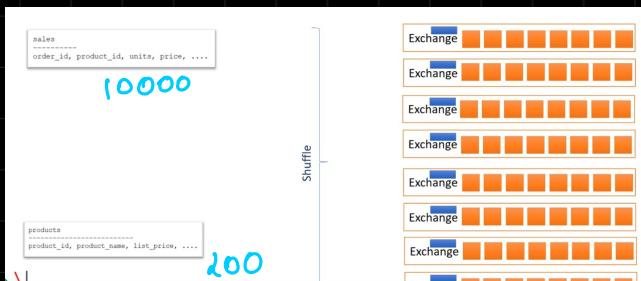


So better partition

Some are slow moving
Some are fast moving

Large - Large Join → Shuffle Joins

Large - Small Join → Broadcast Joins



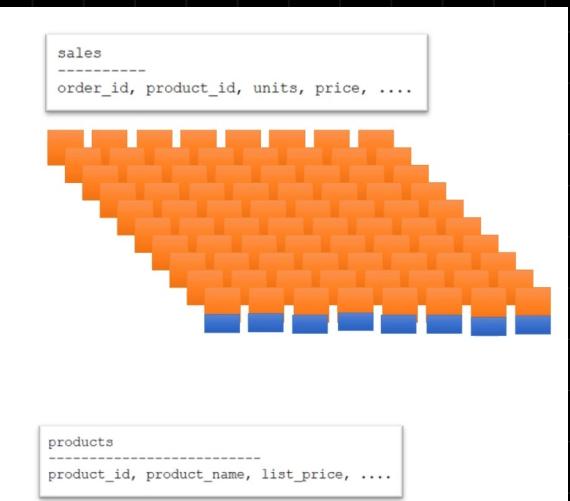
Shuffle Join

So for this it will
create

200 K Reduce
Exchange

Shuffle the 10000

data to this its
huge take more
bandwidth



here 9 shuffle 200

to
10000

as its small
and it perform
faster than the
other

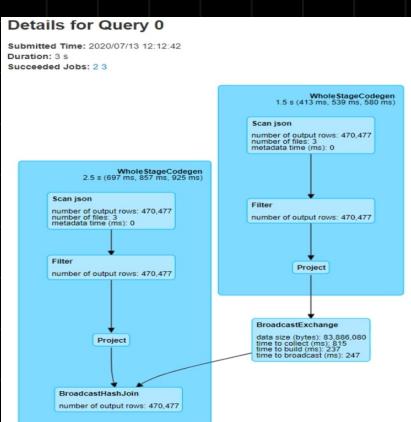
Broadcast Join

↓

It Eliminate the need
of Shuffle

Perform a Join with
out
Shuffle

↙



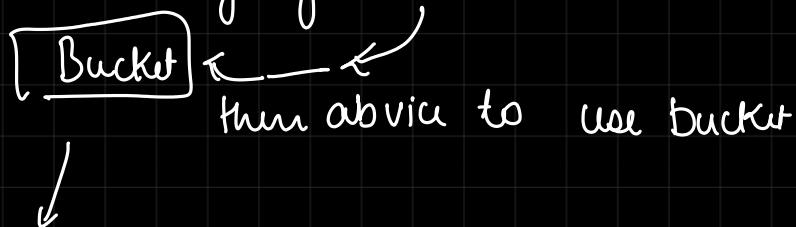
```
join_expr = flight_time_df1.tu == flight_time_df2.tu
join_df = flight_time_df1.join(func.broadcast(flight_time_df2), join_expr, "inner")
```

Implementing Bucketing

Even after doing Broadcast Join
Shuffle Join

Java Memory
not enough out of memory.

If we know we are going to Join 2 Dataset



Bucketing Dataset also need Shuffle however
that needed
Only once

It can be done → into 2 steps

— bucket Dataset

What is a best Performance ?

- ↳ It's not possible to getting ideal Performance
- ↳ Predictable Performance

