# Week 5: Batch Processing

1. intro
2. instulation [spark]
3. spark SQL and Data frames
4. spark internals.

---

## Batch vs streaming   [based on time]

| Batch | Streaming |
|---|---|
| → process chunck of data at regural interval. | – process data on fly |
| – process data each [ day "most common", month hour, year ] | – ong fly immediately. |
| – process taxi-trip data each month<br>– 80% of work is in batch. | process taxi data as soon as generated. |

– we will focous with **Batch** which can be acheived using diffirent methods like:

1. python scripts like data pipeline in lesson1 that was done using corn jobs , (Airflow)

2. (SQL) like (dbt) which we used GUI to control time and scheduling of work. & also using SQL for schedule

3. spark in this lesson.

Pros of batch Jobs:

- Easy to use and scale.
- Re-executable.

.80% of work is in batch.

Disadvantage.
- Jobs takes a lot of time to process as we work with a lot of data.

---

what's spark .........?    "Clusteres, Big data"
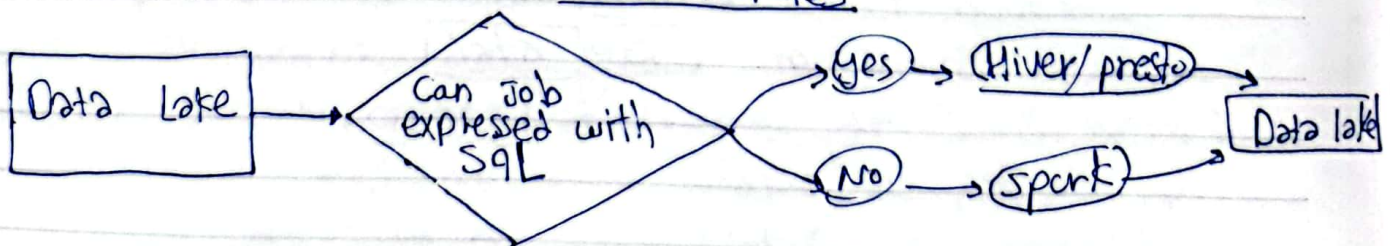                                [Multi Thread]

- Appeach spark is open source multi-language engine For large scale data processing Engine as it process data.
→ = Spark can run Clusters with multi nodes each pulling and transforming data

- Multi language as we can run Java, Scala natieviely Python, R    [py spark]
- Spark can deals with
   • Batches.
   • Streaming " seeing streaming of data as a small batches "

```
Data Lake → [Can Job expressed with SQL] →yes→ (Hiver/presto) → Data lake
                                          →No→ (Spark)
```

- If Job can be done with (SQL) we use Hiver, Athena  '[expressed]
- If not ex[ ML work need to be done then we can use spark]

# Installing Spark

- we can install spark from lunix, windows
- Spark needs Java.
- we can use pyspark, that can enable us to use Jupiter note book.

ex   df = spark . read \ . option ("header","true" ) \
              . csv ('taxi_data.csv')

- spark has GUI also.
- distriputed systems &works

· It's very common to use (Hadoop) as data storage
                          (spark) as data processing
· spark is most open source Data project.
· Spark can do every thing from storage to ML &graphs .
· spark naitvely uses scola .

# Pyspark

"A libirary used to work with spark with python"

```
from pyspark.sql import sparksession

spark= sparksession.builder \ .master("local[*]") \
        .appName('test') \ .getorCreate()
```

- master(["local"]) ; spark will run in (all) aviliable CPU.
- spark has **UI** that can enable us to view jobs.
- spark is very good solution when we work with large dataset (unlike) pandas.                  Large CSV

> df = spark.read \ .option("header", "true") \
        . csv('file-name.csv')

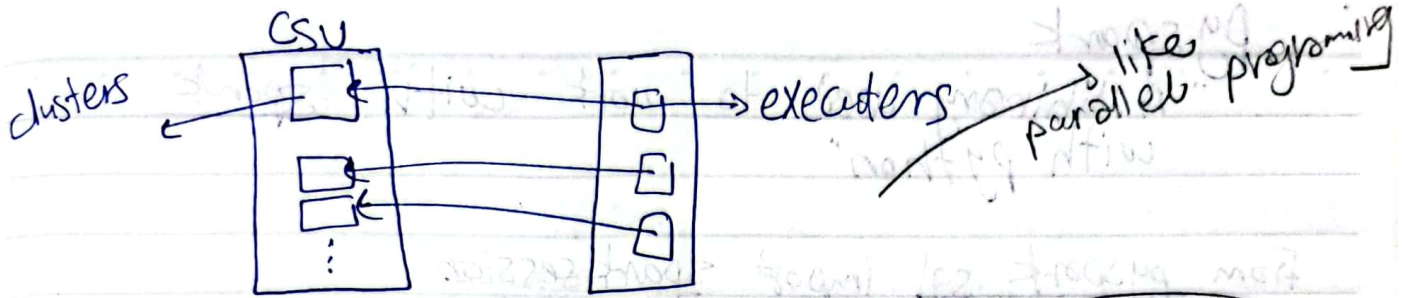- problem is that spark can't handel or infer data type to solue that
> spark.create Dataframe((df)) . schema
        get the some shape of df
        save subset of data as pandass
        then import it's data type.

get datatype from pandas

clusters    CSU      → executers   like parallel programing

Each executer should access a certin (partition)
this is the mesning of partitioning .
parallel and Speed up .

> df = df . repartation (24)
> df . write . parquet ( 'path\')

- By default Spark partation files to
  Same number of (cpu) cores .

_____

work with data frames :
> df = spark . read . parquet ('path')

- parquet files save it's schema unlike
  csv files
- parquet files are smaller than csv
  because they store data according to
  data type so
  - integer values take less space
    than long or string .

# Actions vs Transformations

- some of spark methods are **lazy** or (Transformations)

   e-x   df. repartation ( . 12 )
       df . select ( 'pickup time', 'drop_ ')

- oThers are <u>actions</u> or <u>eager</u>
   e-x   df.show()

---

## Functions & user defined functions ( UDFS )

- Spark enable us to create functions..
   → From pyspark.sql import functions as F

> df \
       → creat new column   → convert to date builtin function
.with column ( 'pickup-date' , (F) todate ( df. pickup-datatine))
  .select ( 'pickup-date ' , 'pulacrtion ID )
   .show (s)
     ↳ eager action      select this columns From data.

. spark has alot of builtin functions, But also you can add your own (UDfs)

```
→ def  crazy-stuff (base-num):
        num = int ( base num [1:] )
        if num % 7 == 0 :
            return F(S)/ {num:03x}
        elif num % 3 == 0
            return F'@/ {num:03x}'
```

crazy-stuff-udf = $\boxed{F.udf}$( crazy-stuff, returnType=Types.stringType())

> df
.with column ('pickup-date', F.to-date ( df.pickup-datetime))
.with column ( 'base-id', ⟨crazy-stuff-udf⟩(df.dispatching-base-m)
  .show()
                          ↳ send value of base to
                            function, which check for it
                            and return string charachter
                            which will be saved in new
                            column called base-id

---

5.3.3 optional "preparing yellow & green taxi data"

- using __bash script__ to download the dataset.
- bash script is similar to __lynix__
- we also can use Air flow DAG
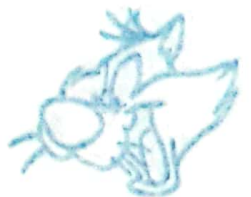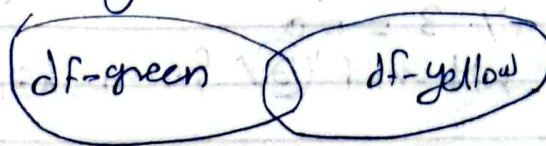  - Also can use python.

5.3.4 spark sql "Combining 2 datasets"
- spark can also run __sql quers__

> df-green = df-green
-withColumn renamed ('1pep-pickup-datatime', 'pickup-datatine')
                          ↳ rename

> set ( df-green.columns) & set( df-yellow.columns)
  get common values between columns.

Another solution using for loop.

( df-green ⨉ df-yellow )

- we will add new column with name of the data set
  to make it easy for insights.

```
> df green
  .select ( common - columns)
  .with column ( 'service_type', F. lit ('green'))
  ↳add string to this column.
```

combine data set
```
> df_a = df_green .union All ( df_yellow)
> dF-a .group By ( 'service_type') . count() . show()
  show values of green & yellow.
```

_Query sql with spark_
- spark expect <u>table view</u>, but we will pass df
  ```
  > df . regiter Temptable ( 'trips_data')
    . make table to use, spark sql.
  ```

```
> spark.sql ( " " "
SELECT
      service_type,
         count(1)
    FROM
       trips_data
    GROUP BY
       service_type
         " " " ).show ()
```
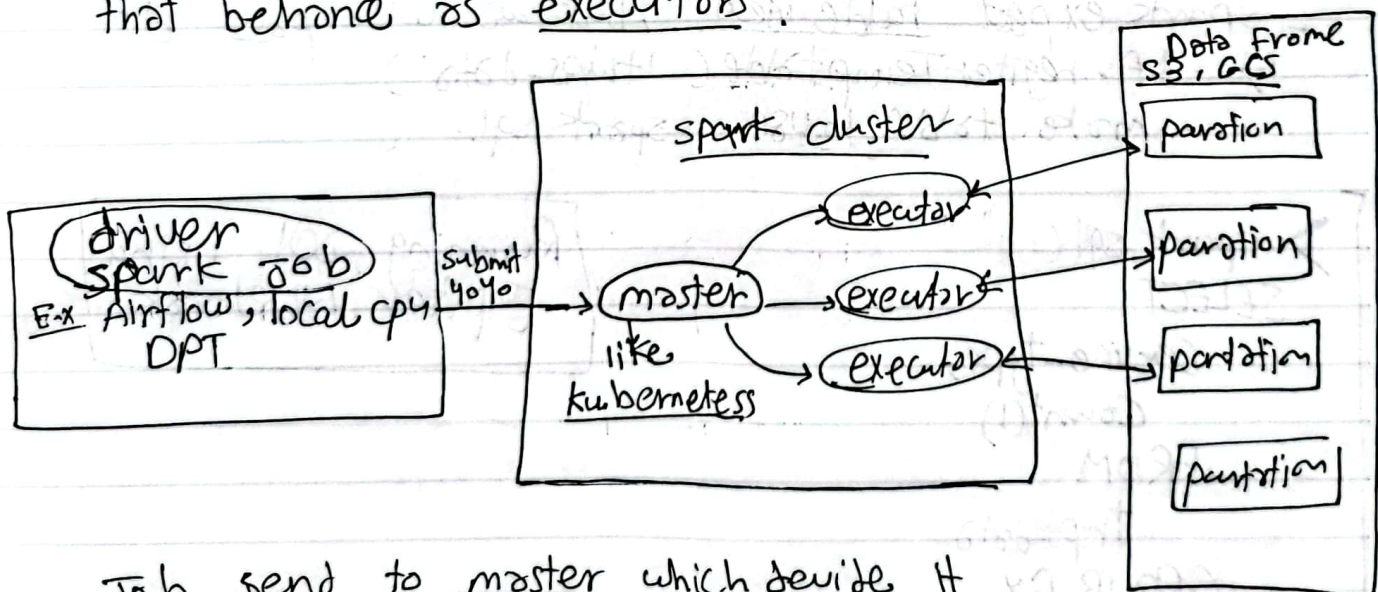
Running SQL using
Jupiter Notebook

saving result as parquet file :
> df.write.parquet (' data/report')

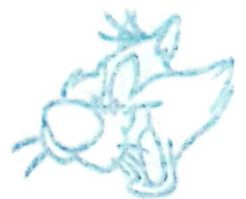This will Save and create more than 200 parquet File due to partitioning

So we will combine or coalesce opposite of partioing :
> df.coalesce(1).write.parquet('data/report', mode='overwrite')

---

## Spark Internals

- Spark clusters often contains multiple computers that behonce as executors.

Job send to master which devide it to exeutors and acess partations on S3 or GCS...

- _Hadoop_ is a good tool, but as _data volumes_ and _analytics_ requirments increase in complexity it simply _dasen't_ _work_ any more.

- _Hadoop_ stores data locally for _data locality_, _portations_ are _duplicated_ across several executors for _redundency_.

- Hadoop has _fallen out of fashion_ as all go now for cloud.

## Graub by - spark

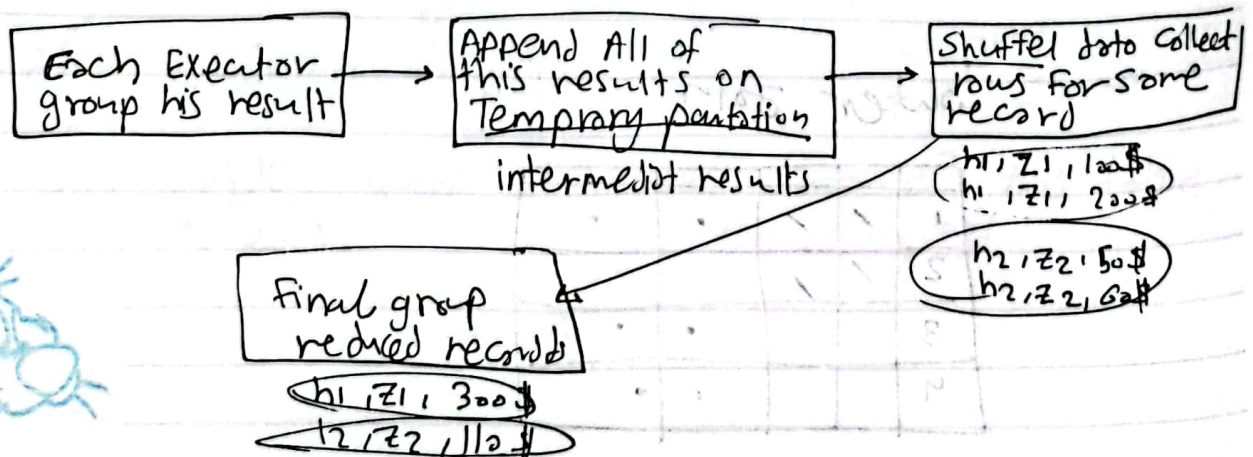- I want to select _total revenu for each hour_ for _each zone_ also.
  - So I must select _each_ area and hour.
  - After that I should sum this values.

- So we need to _group data_ which is in _separate portations_ But executors deals only with _individual portations_. How to solve ----- ?
  - Spark solved that by speroting grouping into 2 stages.

| Each Executor group his result | → | Append All of this results on Temprory pantation | → | Shuffel data Colleet rous for some record |
|---|---|---|---|---|

intermediot results

$h_1, Z_1, 100\$$
$h_1, Z_1, 200\$$

$h_2, Z_2, 50\$$
$h_2, Z_2, 60\$$

Final grop redued records

$h_1, Z_1, 300\$$
$h_2, Z_2, 110\$$

- we can see that from spark UI
  Exchange = shuffel :

- If we added order by that will add another step
  to check order of final sum.

---

## Joins in spark

Joins is smillor to group-by, but instead joins
have 2 distinct Cases:
  1. Joining 2 large Tables.
  2. Joining large & small Table.

---

1. Joining 2 large Tables:

> df_join = df_green-reveny.Join (df_yellow-reveny
          on = [ 'hour', Zone' ] ; how = 'outer)

| Table 1 | | |
|---|---|---|
| 1 | / | / |
| 2 | / | / |

| Table 2 | | |
|---|---|---|
| 1 | • | • |
| 3 | • | • |
| 4 | • | • |

outer Join

| | | | | |
|---|---|---|---|---|
| 1 | / | / | • | • |
| 2 | / | / | | |
| 3 | | | • | • |
| 4 | | | • | • |

For this large values would be simillar to groub by with shuffeling and multi steps.

$K_1$: composite key becouse we join based on 2 values [Zone, Hour]



Yellow-taxi
$K_1, Y_1$
$K_2, Y_2$

Green
$K_4, G_3$

shuffled
$K_1, Y_1$
$K_4, G_3$

$K_1, Y_1, \emptyset$
$K_4, \emptyset, G_3$
outer result

2- large & small table:
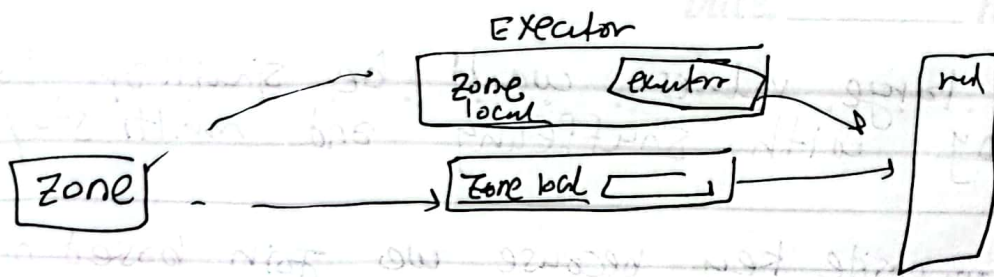
Smal table here is [ Zones-lookup ]

df-yellow has code of Zone which we should merge with small table Zone-lookup

(df-Join. join ( df-Zones, df-Joine. Zone = df-Zones. locationID)

df → df
default inner join
tell Joiner that both columns are same.

small-table spark will send copy of this table to each executor No Shuffel (broadcasting)

Executor



- we will pass "RDD" part as it optional
  and go to "Google Cloud storage"

## 5-6-1    Connecting to GCP

- Some of instruction to connect to GCS
  from local cluster in spark.

- to make it easy for spark to featch and
  process data from cloud.

```
> gs util (-m) (cp) (-r) <Folder>    url
```

Google Storage utilly    copy    reacursivly

→ Multithey
→ download url to folder

Now you can read remote file:

```
>dF-green = spark.read.parquet('gs:// Folder/pq/green/*/*')
```

## 5.6.2 Creating local spark cluster

- run code in cloud
- setting data proc cluster

- spark is very important and you
  need to practice
  - pyspark
  - scala

  } ⟶ Etislate

Recommended    Data comp