text → "rows" → Array of arrays                    ①
sep=\nd.    ↳ Single array

dataframe ↱                              Map
          ↳ Rows (values, idx).        (works on
                                        these arrays
                                        based iterable)

——————————o——————————×——————————o——————————

FILTER.
                                    → Removes the entire
      ↳ Applied on each row of the parent    array, i|
        array.                    → Array              not
                                                       passed

εdd. Filter (lambda x : True).
      ↳ All pass filter.

Map                                → Here Array of
Filter is applied of    iterable        (Arrays)

                                    Applied on all
                                    (works on object
                                    → inside the parent/
                                      wrapper iterable

——————————o——————————×——————————o——————————

Distinct
      ↳ Returns unique objects inside array
Left to right.

→ .rdd • flatmap ( lambda x : x.split('6 ') ) ) • distinct()

  • collect()

→ Left to Right → Chaining ∴ new instances created + method used.

Group By Key.
  ↳ Requires "input" in a particular format
  INPUT :-
  ↳ Array of tuples.
     ↳ of (key, val)
        ↳ Tuple.

  o ——————— x ——————— o

  (key, value)
     ↳ Tuple.

  [ To attain this we need to analyse + use flatmap and maps ]

  ↳ Group/Aggregate by key.
    + Apply some function.
         ↳ map.
            or
            map values.

Returns → key, statistic of the iterable.
                  Estimate

③

**1. Reduce By Key (fn)**

Requisites ⌐
└→ Input ⟹ Iterable of tuples → (key, val)
└→ Function
  └→ Reducer → lambda x, y : f(x, y) for key → One value.

· Always chk. for the kind of object one is are applying use fn on.
  └→ string
     array
     integer.

Ex: INPUT
text_file. →  split (" ")
  └→ "/n"  └→ object = Array
  sep. "/n"     2 possibilities.
                └→ FLATMAP →

Str: Create appropriate for each.

Array of arrays. → MAP → Array of strings.

Map(lambda x : sum(x)) → Array of tuple [(key, val)]

Group By Key ( ).  (la x : (x, 1)]
  └→ (key, iterate) map(lambda x : sum(x)) → Array of tuples.

[key, final_estimate] ← Reduce By Key ( lambda x, y : x+y) → (Requisite of Group By Key Reduce By Key)
Array of tuples.

# ACTIONS → Will always run the architecture
└ Methods └ Needs/consumes processing from ④
└ Methods that renders the pipeline or.

└ Ex. → o Collect ()
    └ List all the values
       in the wrapper
       iterable.
       └ object

o Count ()

    └ Length of the wrapper list.

o Count By Value ()                [Map (lam x : (x, 1) →
    └ Value - counts().               group By Key () → Map
                                      (lamda , x) & sum (x))]

o save As Text File ()
    └ Save the file to the specified path.

    └ Can also create folder:
                            └ mkdir ().
                            └ Existing
                              folder.
                                      └ saved.

o→ Spark 'idl instance
    └ By default creates partitions /clusters
      of EC2 on which it perform.   └ n (by default)= 2.
                                    └ 2 partition
                                       for one act
                                       ion.
      transformations'.
                        └ will save those
                          partitions separately.

→ We can increase or decrease the # of clusters for working on the data.  **⑤**

↳ rdd.repartition (number -of- clusters = 5)

    ↳ ∴ Now process will be divided across 5 nodes.

rdd.getNumPartitions()5 clusters. ↙
↳ 'Partz' Method
    ↳ Count of Clusters

⊙ Repartition
    ↳ Method →  Fn or → will only affect the
         Non- action based → architec - ture.

         How to distribute memory
         por use.

Each Rdd instance can have different
Works like ⊙ set of functions /
Git hub        transformation
commits. ⊙ Clusters in which the data
         will be distributed.

Average → we need → count of keys + sum of keys
↳ Big. Streaming data

groupBy or Reducer groupBy needs
input in tuple of (key, val).
↳ format.

INPUT (tx → Sep = '\n'
     Row → Name, Rate )
     ↳ Sep = ','

Map ( x : (x.split(',')[0],
     (int(x.split(',')[1]), 1) )

This val will itself be a tuple of sum, count.

Array of tuples
     ↳ name, (tuple of sum, ct)

Reduce groupBy (( x, y : (x[0] + y[0], x[1] + y[1]))
     ↳ will fetch / aggregate — same keys + Apply fn on second part of the tuple

Array of tuples.
     ↳ (name, (sum, ct))

Map ( lambda x : (x[0], x[1][0] / x[1][1]) )
     ↳ (name, avg).

Array of tuples-