

result

result.coalesce(1).write()

↓ ①

05x160
50

num = 100

df1 > (id, name, dept)

↓

df1 > (id, name, dept, salted_key)

id, name, dept, salted_key

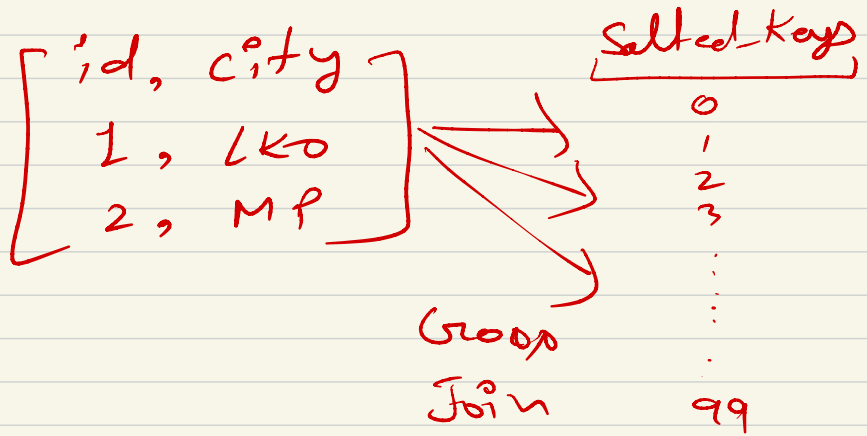
1, a, n, 50

2, b, y, 10

1, a, n, 12

$$\left[\begin{array}{l} 1, a, x, (18) \\ 2, b, y, (21) \end{array} \right]$$

$df2 \Rightarrow (id, city)$



$df2_exploded =$

id	$city$	$salted_key$
1	LKO	0
1	LKO	1
1	LKO	...
1	LKO	18
1	LKO	...
1	LKO	50
1	LKO	...
1	LKO	99

2, MP, 0
2, MP, 1
2, MP, 2
2, MP, 10
2, MP, 99

df1 ⇒

id	name	dept	saltedKey
1	X	a	50
2	X	b	2
1	X	a	12
1	X	a	18
2	X	b	10

df2-exploded

id	city	saltedKey
1	Lko	0
1	Lko	12
1	Lko	18
2	Uko	50
2	MP	2
2	MP	10
2	MP	99

from df1

inner join df2-exploded

On df1.id = df2.id and df1.salt
= df2.salt

--master local[*]

\hookrightarrow All vcores
 \hookrightarrow n vcores

spark-submit

--master yarn

--num-executors 10

--executor-memory 4G

user memory

$$= \text{usable} * (1 - 0.6)$$

$$= 3700 * (0.4)$$

$$= 370 \times 4$$

$$= 1484 \text{ MB}$$

$$0 + 1 + 2 = 3$$

