

query = df.writeStream \ → checkpointing
• trigger (processingTime="10sec") Interval

① Global aggregation

grouped_data = df.groupby('Country').agg(_{sum(sales)})

Country, Sales

IND, 50
USA, 20
IND, 30

states after 1st sec:

{ "IND": 50 }

state after batch 2

{ "IND": 70,
"USA": 30 }

↓ Data Computed

y

② Window Based

window_grouped = df.groupby(timestamp_col, window_size)

IND, 80 } 10:00:00
USA, 20 } 10:00:15
IND, 60 } 10:00:29
USA, 20] 10:00:35

Window 1 = 10:00:00 to 10:00:30
State after first batch

Window, Country, Value
00:00 - 00:30, IND, 60
00:00 - 00:30, USA, 20

State after Batch 2

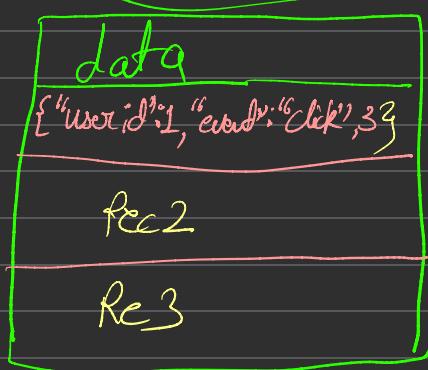
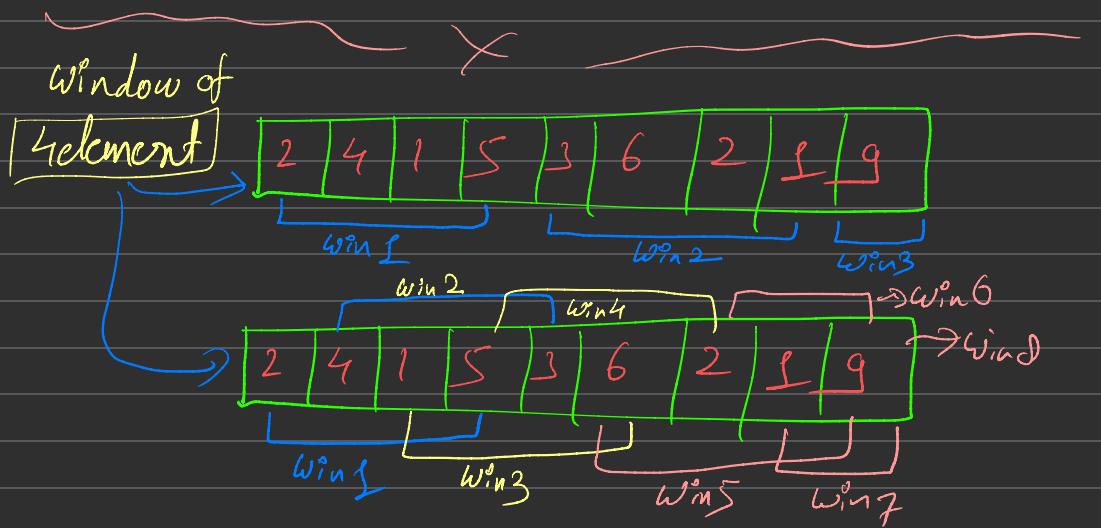
Window 2 \Rightarrow 00:30 to 01:00

Window, Country, Value

00:00 - 00:30, IND, 40

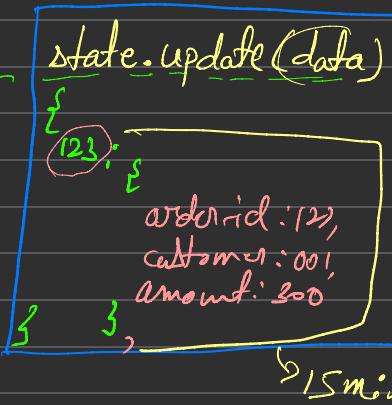
00:00 - 00:30, USA, 20

00:30 - 01:00, USA, 20



State Struct Type

```
def update_state():
```



→ order_id (as key)
→ data (as value)

order_id: 123,
customer_id: 001,
amount: 300

id	Value
A1	20
A3	30
A2	10

→ 2 Sec

→ 1 Sec

update_state([A1], [(A1, 20)], state)

current_state = state.get if state exists else

current_state = { sum: 0,
 count: 0 }

{ sum: 0,
 count: 0 }

total_sum = current_state["sum"]

total_count = current_state["count"]

total_sum = 0

total_count = 0

$\left[\{ "id": "A1", "value": 20 \} \right]$

Q) \Rightarrow for pdf in pdf ifen:

id	value
A1	20
A1	10

total_sum = total_sum + pdf["value"].sum()

total_count = total_count + len(pdf)

id	value
A1	20
A2	50

$$\begin{aligned} \text{total sum} &= 0 + 50 \\ \text{total count} &= 0 + 2 \end{aligned}$$

$\left[\{ "id": "A1", "value": 20 \}, \{ "id": "A1", "value": 10 \} \right]$

$\left\{ \begin{array}{l} \\ \{ "A2": \end{array} \right.$

$\left[\{ "id": A2, "value": 50 \} \right]$

Before Start

State = Empty

(A1)

Current state

= {

Sum: 0,

Count: 0

} 3

(A1)

After Process

= {

Sum: 50

Count: 2

} 3

updated stack

state update (new)

State = (A1)

{

Sum: 50,

Count: 2

} 3