

Recap

Demand Prediction \rightarrow NYC taxi dataset.

Rows \rightarrow 1 ride \rightarrow Pickup location
 \rightarrow Drop location
 \rightarrow Pickup time
 \rightarrow No. of passenger

Predict \rightarrow demand \rightarrow At a given time interval \rightarrow 15 min
and
at a particular region
 \downarrow
No. of pickups

Time \rightarrow 8:00 pm
8:01 pm \rightarrow demand \propto



next minute

More rides \rightarrow more share

We can predict / forecast more accurately

25°C \rightarrow 24°C \rightarrow 26°C

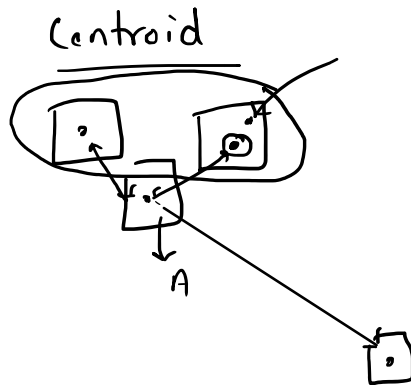
25°C \rightarrow 18°C

25 °C → confidence level ↓

15 min interval → Enough time to move to a nearby region
↓
More accurate

Region → NYC → unsupervised ML → Clustering

↙
k means clustering



Centroids

↓

Identity of our region

distances → Region center

Sort on the basis of distances.

nearly locations

EDA → How we can divide → regions Task 1

How we can divide the time axis Task 2

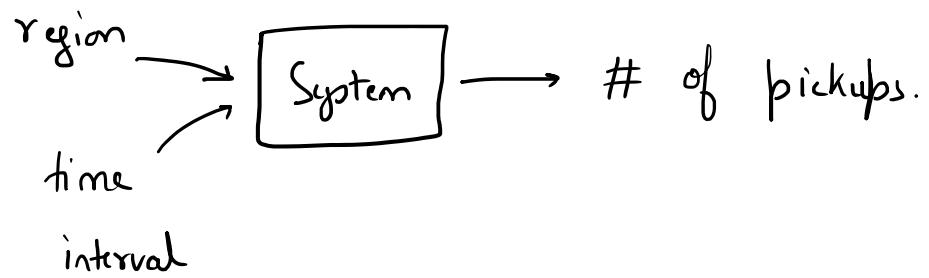
↓

← intervals

15 min intervals

for a given region and for a given time interval

↓
Predict # of pickups



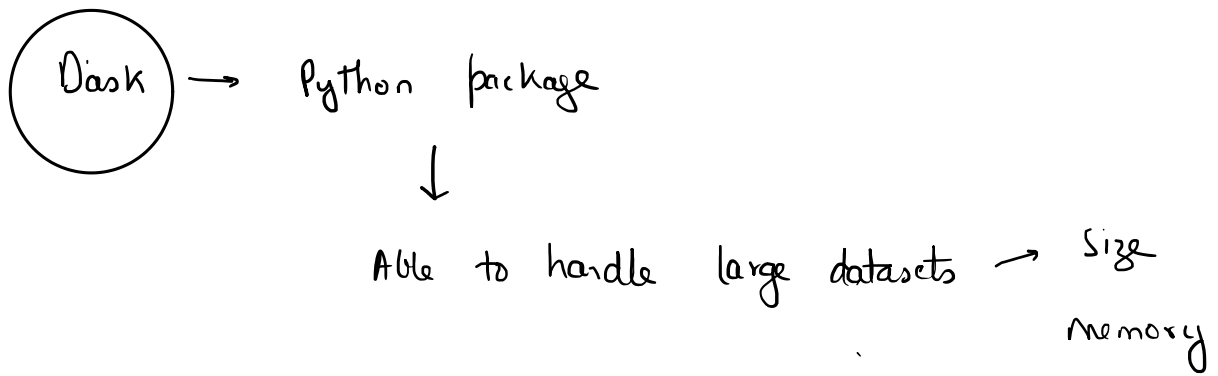
Metric → MAPE

App → Streamlit app

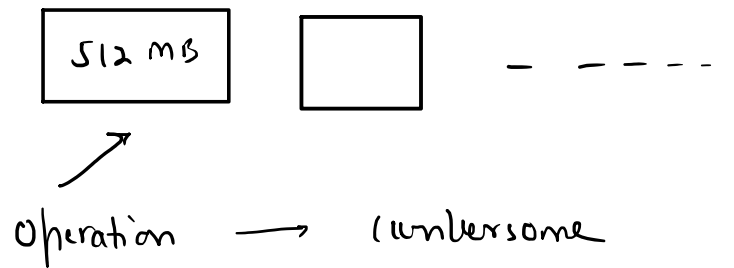
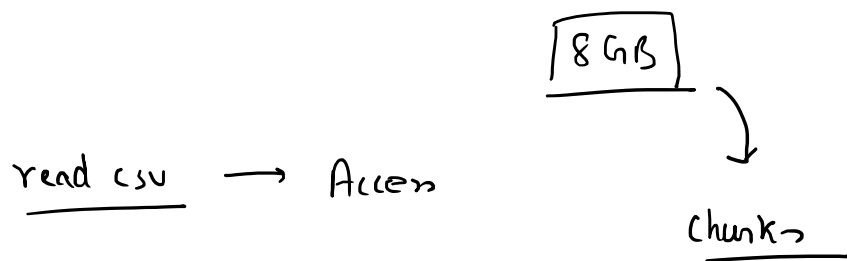
location → Graph → NYC → Regions
time interval

More darker ↑ ← Demand predict
More lighter ↓ ← Color coding

What is Dask?

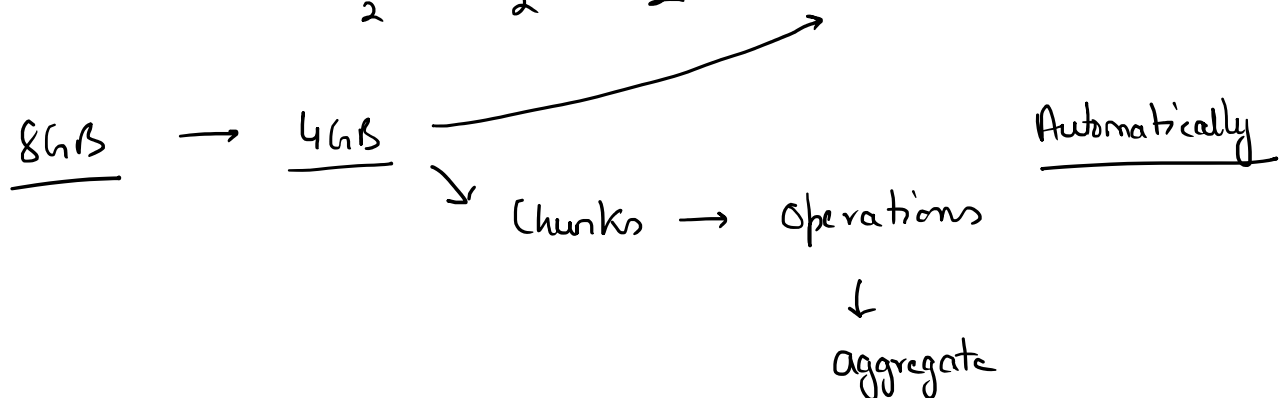


Out of memory → 10 GB



Dataset is huge in size → Dask

2016 → Jan, Feb, Mar → 6 GB
2 2 2



Why use it

Dataset → 2016 → Jan, Feb, Mar

↓
2GB each → 6GB

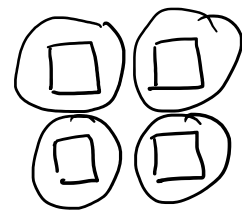
Chunking → Parts

Dask dataframes → Pandas dataframe

Dask arrays → numpy arrays

Handles all the chunking operations automatically

Parallelize → CPU core



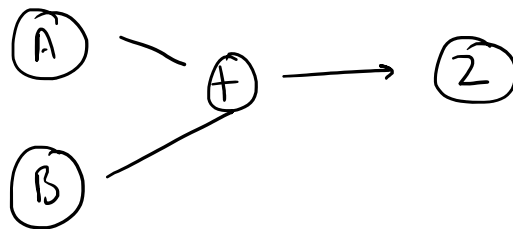
result 1
result 2
result 3
result 4

Result ← aggregate

Task → Pytorch computation graphs

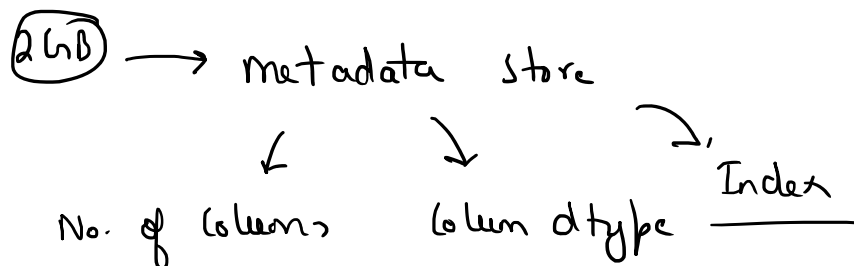
Graph store

$$A + B \rightarrow Z$$



Task → Task graphs

Jan 2016 → load as a task df.



<u>task df</u>	Column A	Column B	Column C
↓	—	—	—
<u>chunks</u>	—	—	—
	—	—	—
	—	—	—
	—	—	—

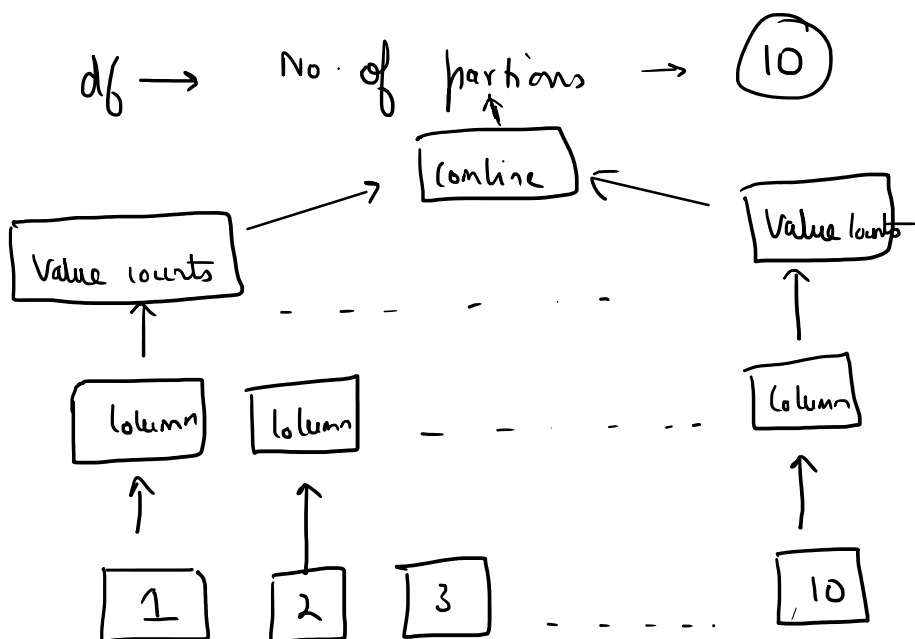
No. of partitions

Categorical → Value counts

`df ['column name'].value_counts().sum()`

Index

category A	—
b	—
c	—

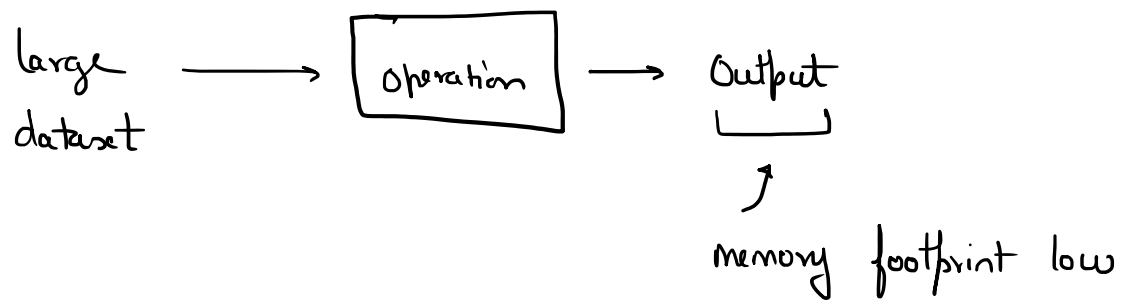


Each chunk
can be parallelized

Return a series

Task → compute() → scalar value

Chain operations



Things to avoid

1) Compute → Uses a lot of compute

2) Try to use pandas → 6GB ← Dask