

Recap

25 February 2025 17:07

Demand Prediction \rightarrow # of pickups

Dataset \rightarrow NYC Taxi \rightarrow 2016 \rightarrow Jan
Feb
March } 2GB

Row \rightarrow single ride \rightarrow fare
distance
Pickup coord
Drop coord

Region
Time \rightarrow \rightarrow System \rightarrow # of pickups

Outliers \rightarrow Remove

Task \rightarrow chunks \rightarrow operations

\downarrow

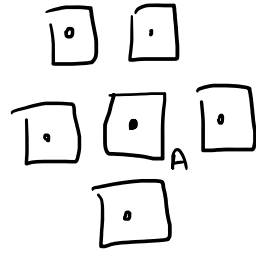
Removed the outliers.

✓ Pickup datetime lat pickup long pickup
time, date

3.3 crore

Task 1 \rightarrow Break our NYC into regions
unsupervised ML \rightarrow clustering

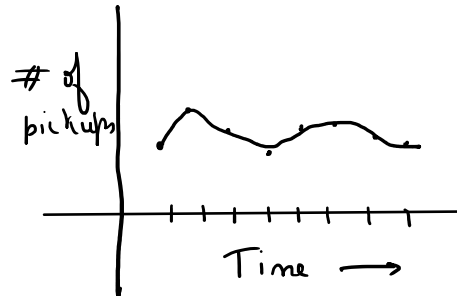
Task 2 → Historical data for pickups for each region.



Time interval

what are the
of pickups

Time axis → 15 min intervals



3.3 crore lat long

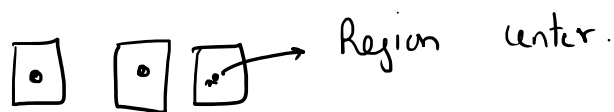
Task 1 → Mini Batch K Means → Batch size

Batch → Centroid distance

Pandas chunking → Partial fit

chunk → Batch → centroids

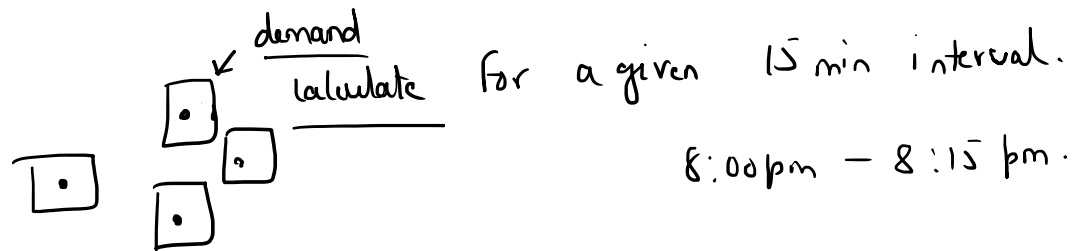
Regions → Centroid coord (Region centers)



Task 2 → Region Datetime

15 min interval → Resample

15 min interval \rightarrow Resample



extra feature \rightarrow avg demand avg - pickups
for each region and for each
15 min interval.

EWMA \rightarrow Historical data + Current observation O_t
 \downarrow \downarrow
weight weight

\rightarrow exponentially decaying in nature.

Datetime	Region	<u>total pickups</u>	<u>avg - pickups</u>	derived feature.
—	—	—	—	
—	—	—	—	
—	—	—	—	
—	—	—	—	

	datetime	Region	avg pickups	total pickups (target)
✓ <u>Region</u>	—	—	—	—
✓ <u>Time interval</u>	—	—	—	—
	—	—	—	—
	—	—	—	—

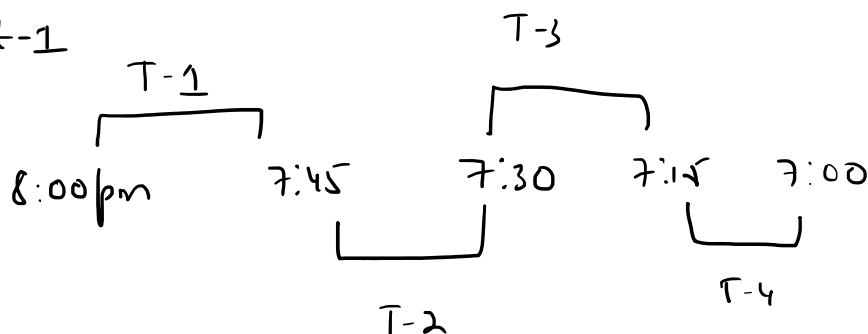
Features → lagged features.

T₀ → current time interval.

T-1 T-2 T-3 T-4

Total
pickups

t-1



input data T-1 T-2 T-3 T-4 Region Avg Pickups

train test split

time axis

Month → Jan, Feb (Train set)

March (Test set)

for every region and every 15 minute interval.

$$\text{EWMA} \quad \alpha = \underline{0.9} \rightarrow \underline{0.4}$$

$$\text{last } N \text{ timesteps} \quad \text{EWMA} \quad \alpha = \frac{2}{N+1}$$

last 4 timesteps.

$$t-1, t-2, t-3, t-4 \quad \alpha = \frac{2}{4+1} = \frac{2}{5} = 0.4$$

$$\text{EWMA} = \underline{0.4}$$

We are actually looking into historical data

EWMA Avg \rightarrow Some part of history in the average

$$\frac{\text{Data leakage}}{\text{features}} \quad \text{⑩} \rightarrow \text{EWMA}$$

$$\alpha = \underline{0.4} \rightarrow \text{Historical avg of prev 4 timesteps.}$$