

## Handling Numerical Data

### → Encoding Numerical Data

When Encode Numerical to Categorical?

Say for ex you have column like  
No of downloads in playstore.

Few Apps has 100 + downloads, few have  
1000 +, few have 10000 + . 1000000 + etc

So now if we plot this we will have  
graph



Now this is sometimes troublesome

so we can create bins where based  
upon bin condition we can add our  
data accordingly

1st Bin → 100 +                          2nd → 1000 +  
3rd → 10000 + etc etc.

Two famous Techniques to convert  
Numerical data to Categorical

1) Discretization which is called  
Binning

2) Binarization

## Discretization or Binning

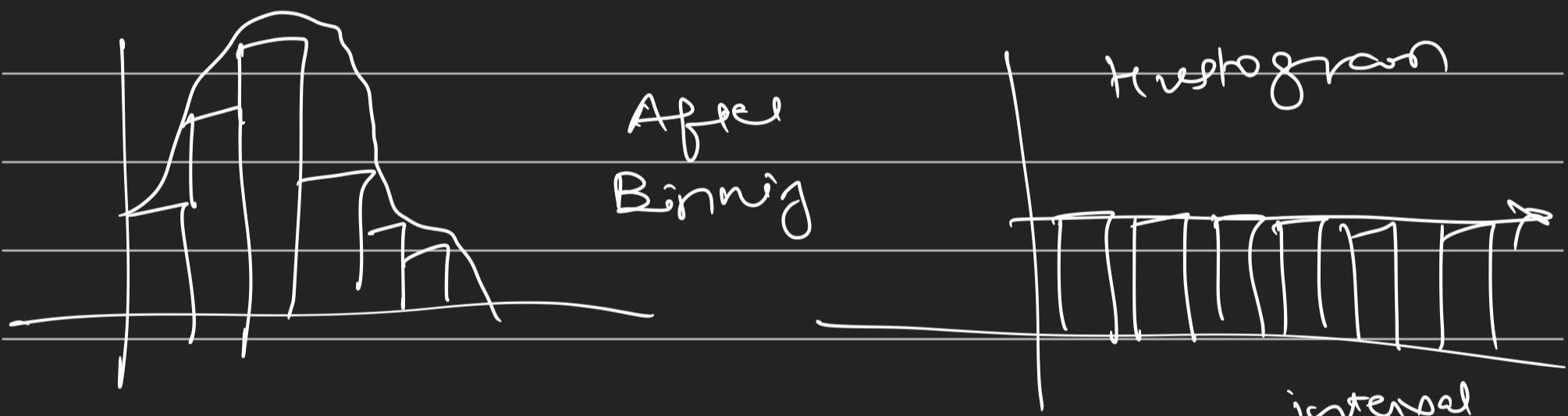
process of transforming continuous value

into discrete variables by creating set of continuous intervals that span on range of variable value.

Why Discretization is called Binning?  
→ When intervals can be called as bins for storing continuous variable we call it binning.

### Advantage of Binning

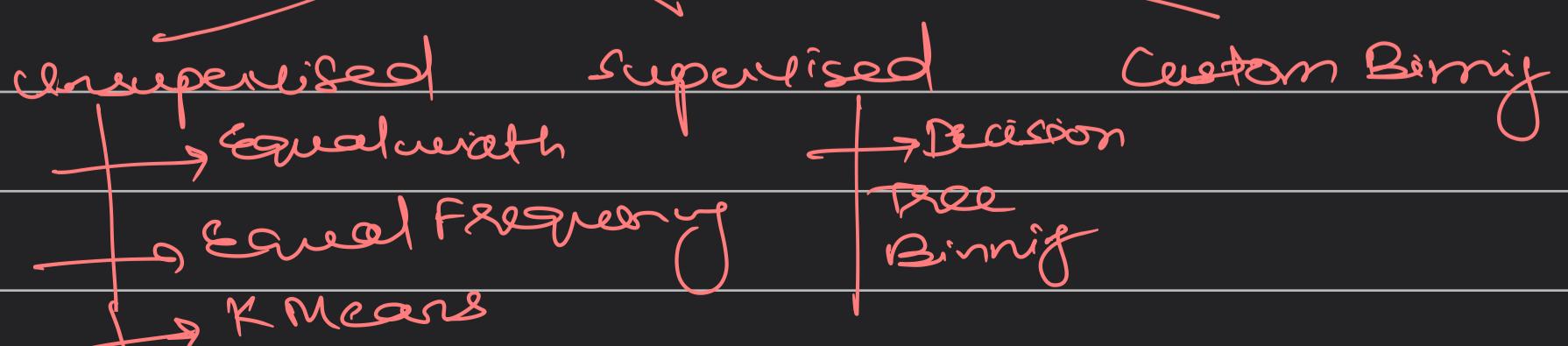
say we have a age distribution 0 - 100 years & for that our histogram is like this



- Adv → 1) We handle outliers  
2) Improve value spread

will come under 1st or last bin accordingly

### Types of Binning



## ▲ Unsupervised Binning

### • Equal width or Uniform Binning

we give No of Bins upfront

$$\text{Interval Range} = \frac{\text{Max} - \text{Min}}{\text{No of Bins}}$$

$$\text{Ex} = \frac{100 - 0}{10} = 10$$

Intervals would be

$$0^1_{-10}, 10^2_{-20}, 20^3_{-30}, 30^4_{-40}, 40^5_{-50}, \\ 50^6_{-60}, 60^7_{-70}, 70^8_{-80}, 80^9_{-90}, 90^10_{-100}$$

total 10 bins

So for a particular variable / column value we will put in related bin accordingly.

somewhat

Now outliers gets handled <sup>in</sup> easily,  
histogram remains same as no change  
in spread.

### • Equal Frequency / Quartile Binning

we give no of Bins upfront.

Say 10

Now each interval or bin will contain 10% of total observation

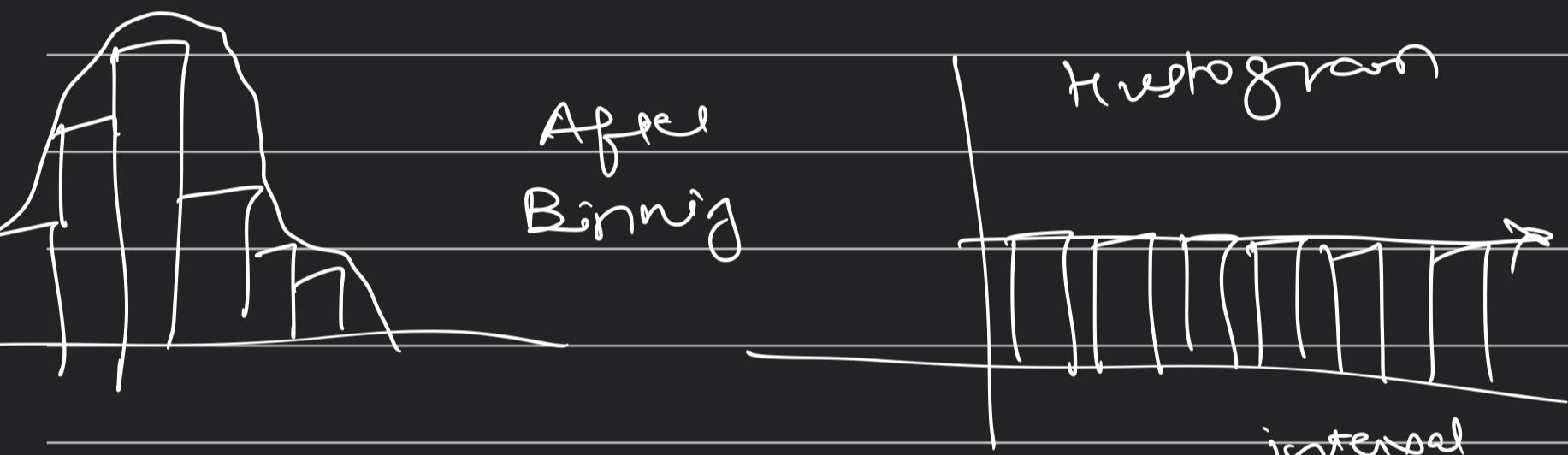
value

Interval would be 0 to 10<sup>th</sup> percentile  
Outlier data would be null  
is not present

0 - 16, 16 - 20, 20 - 22 etc  
<sub>(0<sup>th</sup> 20<sup>th</sup> percentile value)</sub>

Here the spread gets charged, perform well on outliers, here the spread becomes uniform

Most used Binning Method.



## • K Means Binning

used on Categorical data & works on the fundamentals of K-means Clustering algorithm.

Read more when you study K Means Clustering Algorithm.

What after this?

After the process we do Encoding using ordinal or nominal based upon selection in variable values

## • Custom Based Binning

we create bins on the basis of our domain knowledge.

Ex → Age we created bins like  
[0 - 18] → Young  
[18 - 60] → Adult  
[60+] → Senior

## • Binarization

is a special case of binning, here we convert continuous value to binary [0 or 1]

Ex Annual Income

If Income  $\leftarrow$  100,000 K  $\rightarrow$  0 ]  
Income  $>$  100,000 K  $\rightarrow$  1 ]  
At this value no tax  
Taxing ↙

Mostly used in Image Processing.

