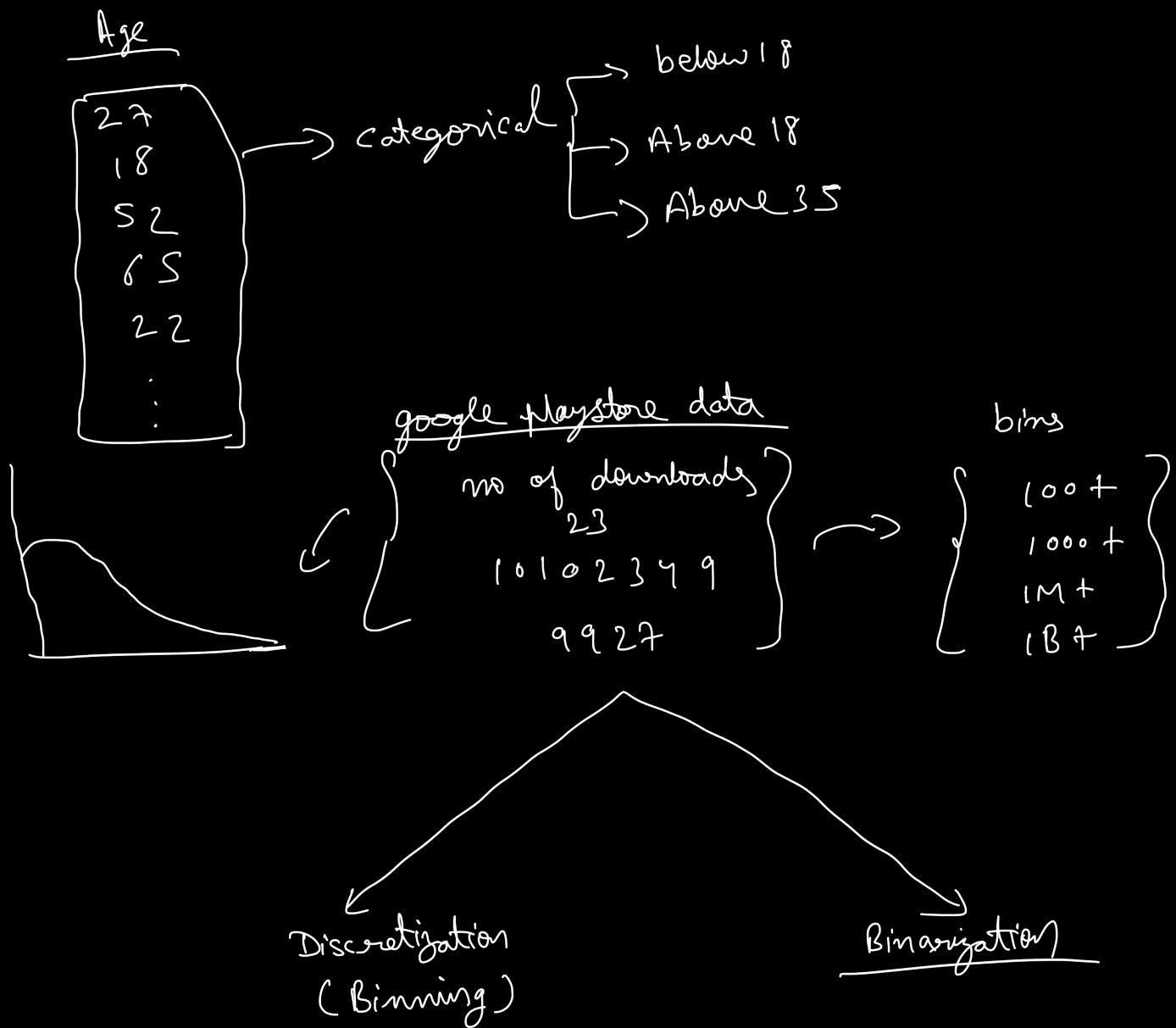


# Binning and Binarization

## Discretization | quantile Binning | KMeans Binning

Encoding Numerical features  $\Rightarrow$  categorical



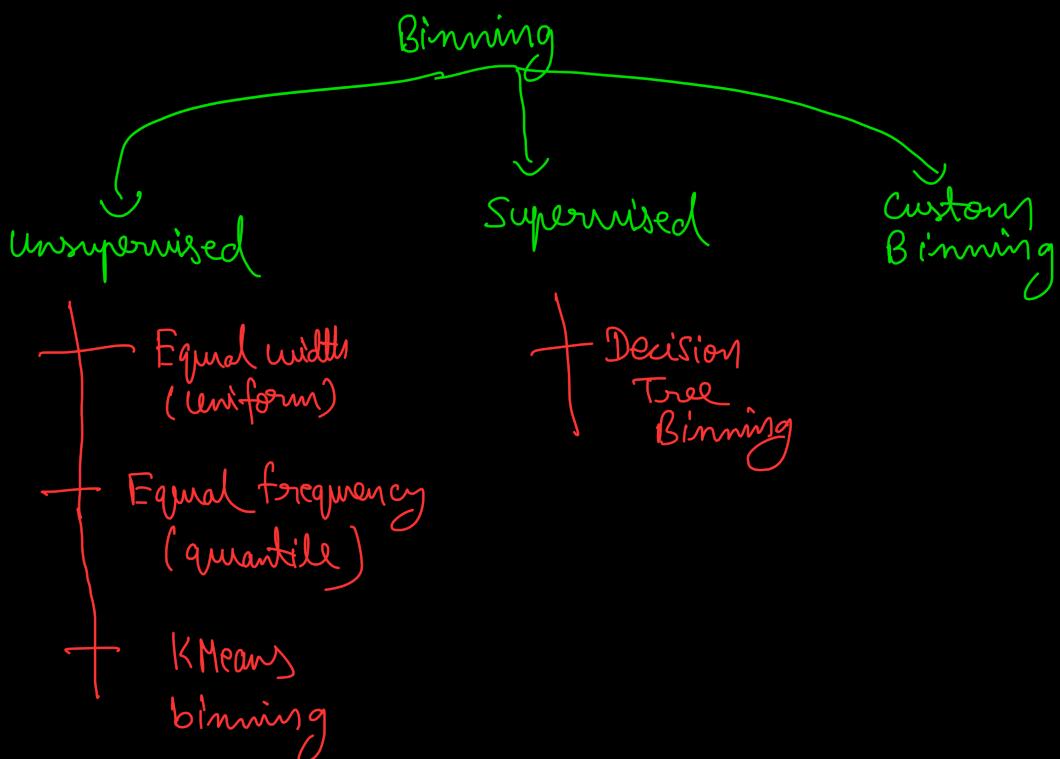
## # Discretization (Binning)

Discretization is the process of transforming continuous variables into discrete variables by creating a set of contiguous intervals that span the range of the variable's values. Discretization is also called, where bin is an alternative name for interval.

→ Why use Discretization

- (I) To handle outliers
- (II) To improve the value spread

# Types of Discretization



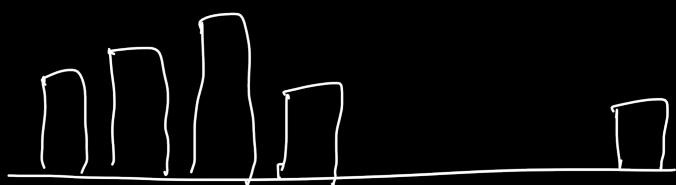
## # Equal width / uniform Binning

Age  
27, 32, 84, 56, . . .  
 $\boxed{\text{Bins} = 10}$

$$\begin{matrix} \max & 100 \\ \min & 0 \end{matrix}$$

$$\frac{\max - \min}{\text{bins}} = \frac{100 - 0}{10} = 10$$

$$(0 - 10), (10 - 20), (20 - 30) . . . (90 - 100)$$



- (i) helps in handling outliers
- (ii) No change in spread of data

## ~~#~~ Equal Frequency / quantile Binning

Intervals = 10

Each interval Contains 10% of total observations

Interval:

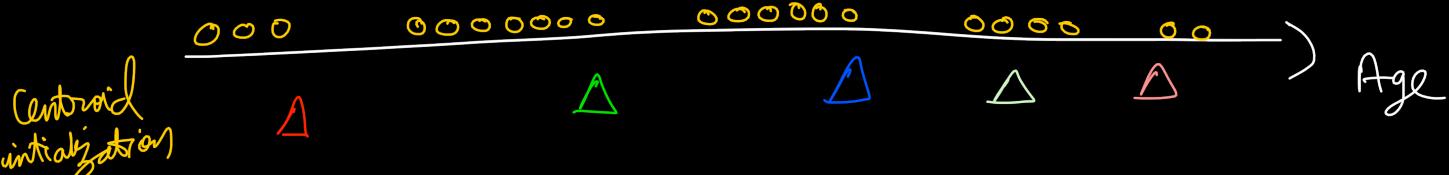
0 - 16; 16 - 20; 20 - 22; 22 - 25; ... 50 - 74

\* This is default value in sklearn  
class.

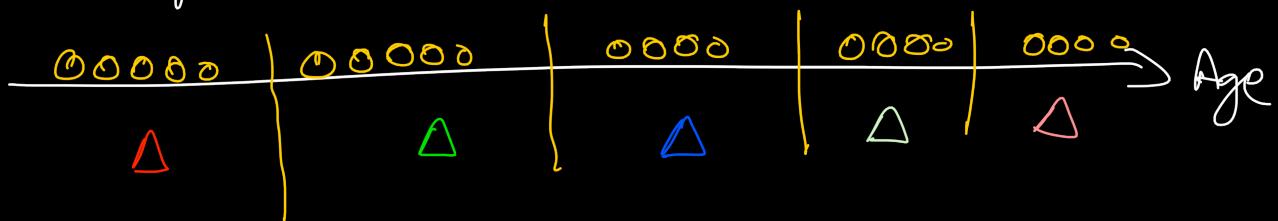
- (i) used mostly (better than previous one)
- (ii) helps in handling outliers
- (iii) Make value spread uniform

## # KMeans Binning

It is used when data is distributed in clusters.



- 1) calculate distance of each point from each centroids and make cluster with nearest centroid.



- 2) Now again centroid of each cluster change according to mean of values inside it and then again repeat step 1 and make cluster.
- 3) we will repeat this until there is no change in previous and current cluster.

4) And finally After getting centroid from final cluster we take centroid value, which is the bin value for our binning.

### Encoding the discretized variable

SK Learn

↳ KBinsDiscretizer()

bins

Strategy  
+ uniform  
+ quantile  
+ Kmean

encoding  
+ ordinal  
+ one hot encoding

### Custom / Domain Based Binning

- \* According to domain knowledge , we decide our bins or category like

[ 0 - 18 ] — Kid

[ 18 - 60 ] — Adult

[ 60 - 80 ] — Old

### # Binarization

Sklearn Binarizer class

threshold

copy = True  
False