

Histogram

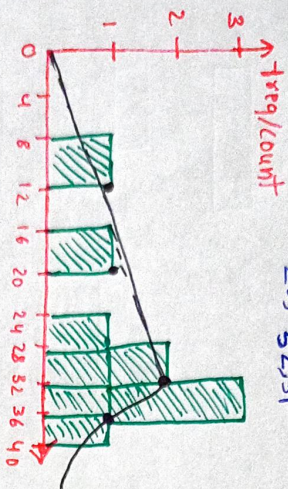
To create histogram do:

- Sort the values (asc)
- Create bins (groups)

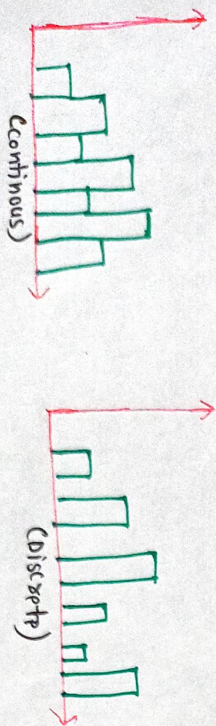
$$\text{No. of Bins} = \text{My choice}$$

$$\text{No. of bins} = \frac{\text{Last value (Max)} - \text{First value (Min)}}{\text{or No. of Bins}}$$

Ex: Values = [10, 20, 25, 30, 35, 40]
2, 6, 3, 2, 3, 1



Plotting Continuous values as well as discrete.



PDF: Probability density function

PMF: Probability mass function

PDF use to smooth discrete histogram and PMF used to smoothen continuous histogram.

Measure of Central Tendency (MCT)

MCT is single value that attempts to describe a set of data identifying the central tendency.

① Mean/Average:

$$\text{Population mean} = \sum_{i=1}^N \frac{[x_i]}{[N]}$$

$$N \gg n$$

$$\text{Sample mean} = \sum_{i=1}^n \frac{[x_i]}{[n]}$$

But, $\bar{x} > \mu$ or $\mu > \bar{x}$ depending on sampling technique.

➔ In practical, we use mean to replace NaN values [although not efficient] due to outliers]

② Median:

To calculate it do

➔ Sort the values

➔ Find central numbers

➔ if no. of values even, find average of central numbers
➔ if no. of values odd, find the central number

➔ central number is value in mid of the all values.

➔ If no outliers use mean, if outliers use median to replace the NaN values.

③ Mode:

most frequent occurring element
➔ Used to replace the NaN values in categorical data within the most frequent value.

④ Measure of Dispersion:

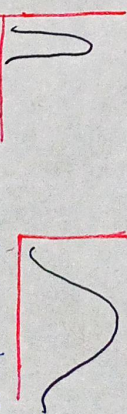
$$\text{Variance } (\sigma^2)$$

• spread of data

$$\text{Population variance} = \sum_{i=1}^N \frac{(x_i - \mu)^2}{N}$$

$$\text{Sample variance} = \sum_{i=1}^n \frac{(x_i - \bar{x})^2}{n-1}$$

Variance \propto spread of data



less variance,
less spread

$$\text{Standard Deviation } (\sqrt{\sigma^2})$$

• Tells how many standard deviation a value is away from mean.

Ex: $\rightarrow [1, 2, 3, 4, 5, 3, 4, 3, 2]$, $\mu = 3$, $\sigma^2 = 2$, $\sigma = 1.41$

