# Random variable

$$X = 10 \smile$$

$$Y = Pass \smile$$

| age $X_1$ | weight $X_2$ | hieght $X_3$ | $Y$ |
|---|---|---|---|
| $\not=$ | $-$ | | |
| $\not=$ | $-$ | | |
| $\not=$ | $-$ | | |
| $-$ | $-$ | | |

$$X \approx Y$$

Random variable

$$X = \text{coin (Toss)}$$

$$X = T/H$$

$$X = Dice$$

$$X = 1/2/3/4/5/6$$

# Type of Random variable

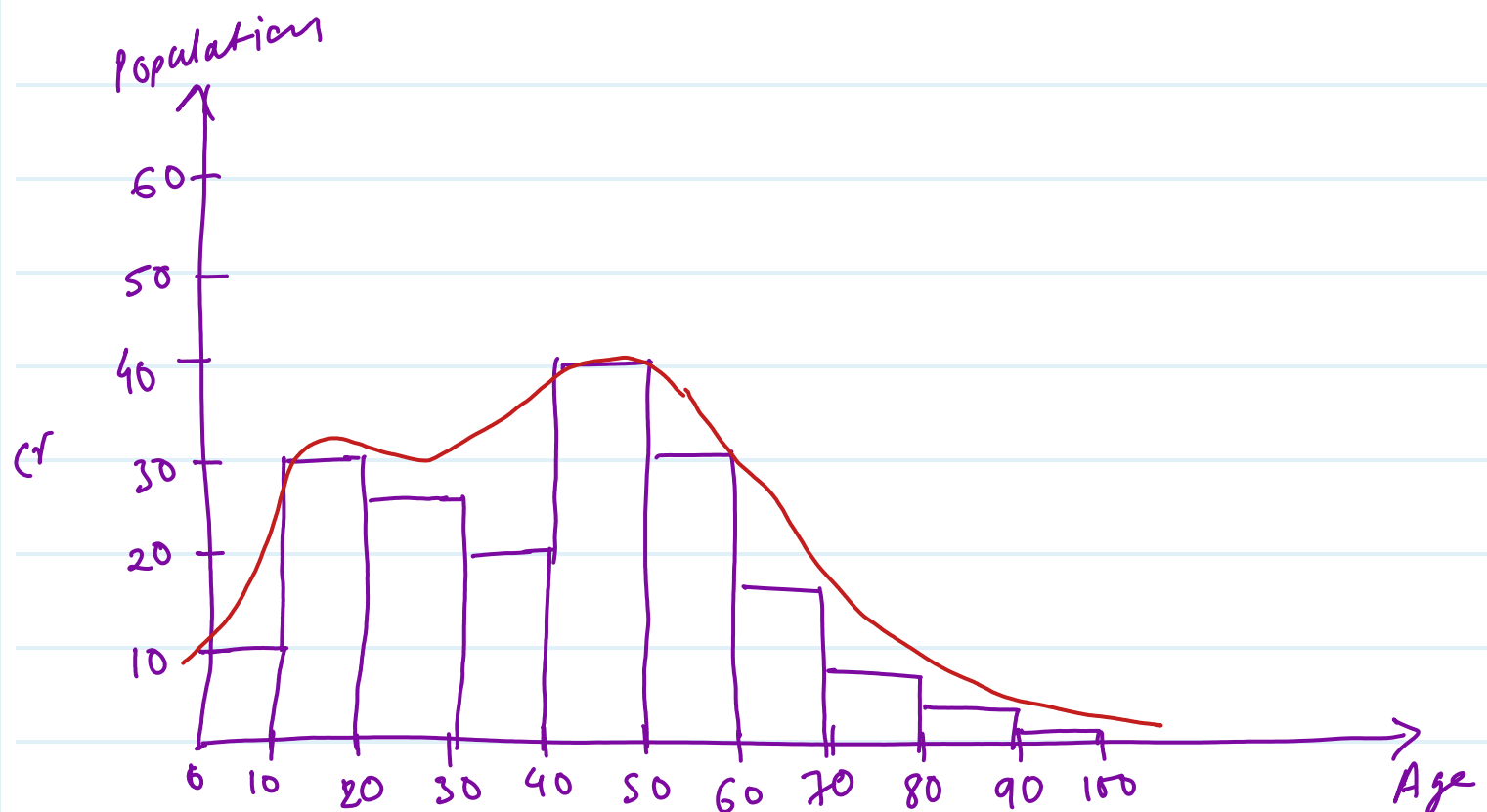Descrite random variable

Continous Random variable

$X_1 = $ Gender
$= m | F | T$

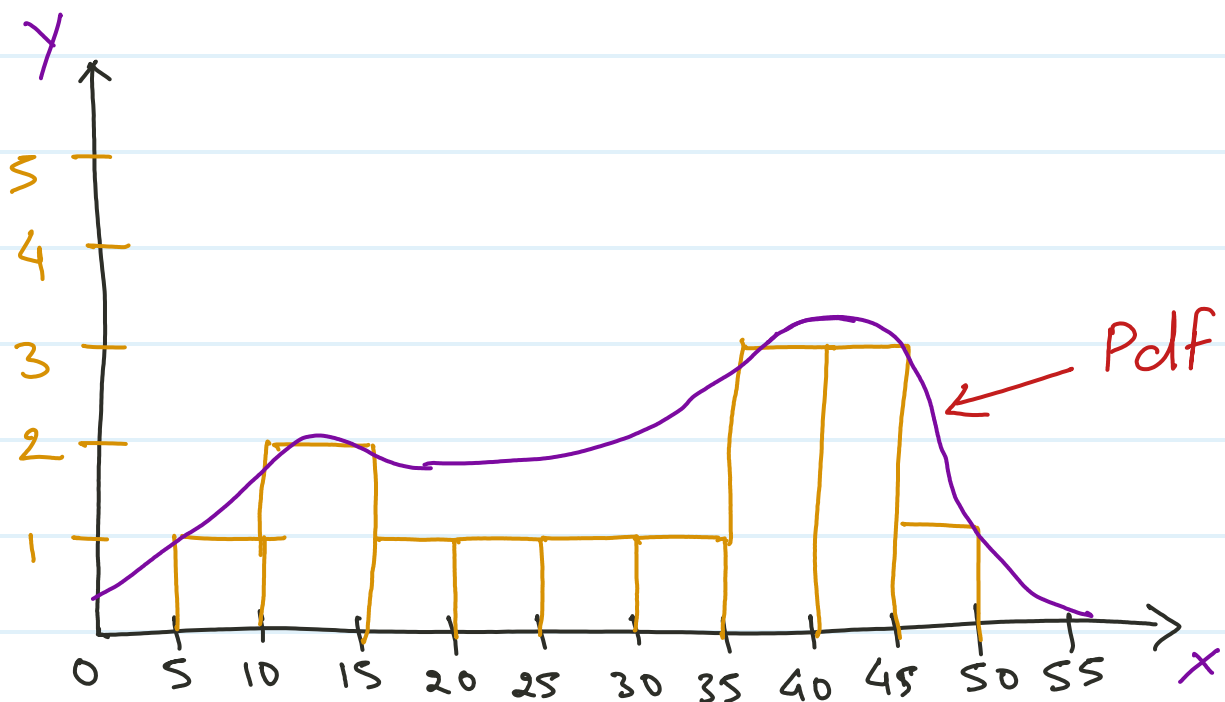$X_2 = $ Result
$(P/F)$

$X_1 = $ Rain
$(1mm, 6.5mm)$

# ✳ Histogram

Dataset = [10, 12, 14, 18, 24, 30, 35, 36, 37, 40, 41, 42, 43, 50, 51]
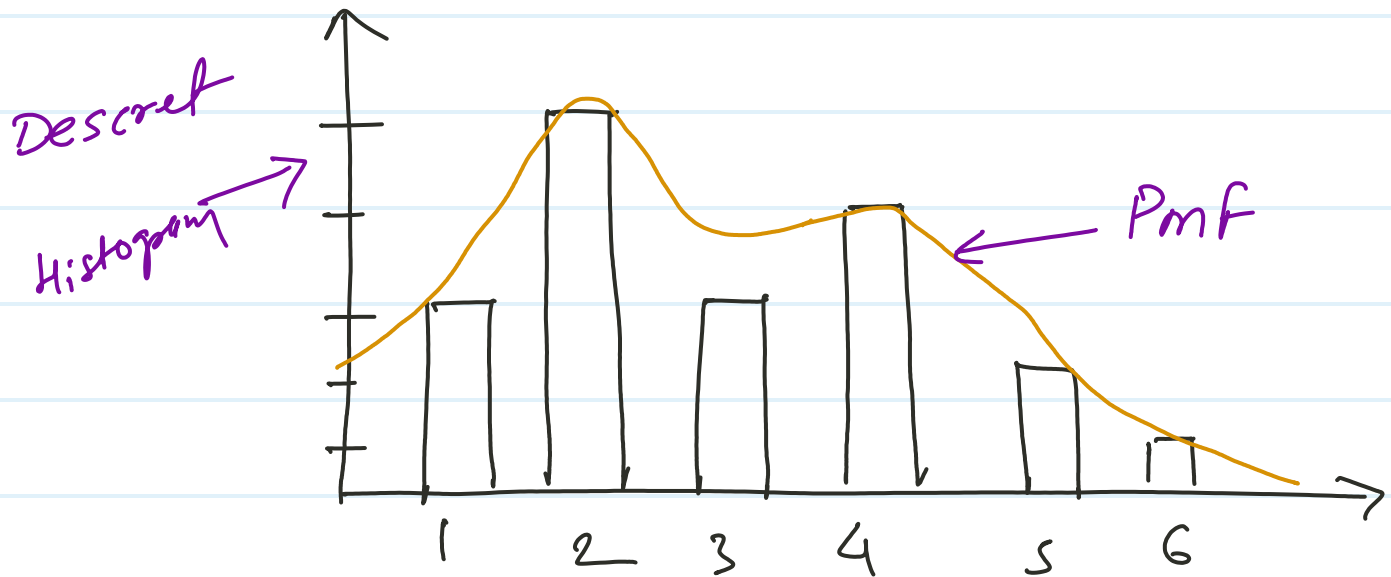
Bin / Bin Size    <u>Assume bin size = 5</u>

$\Rightarrow$ No of Bin = $\dfrac{50}{5}$ = 10
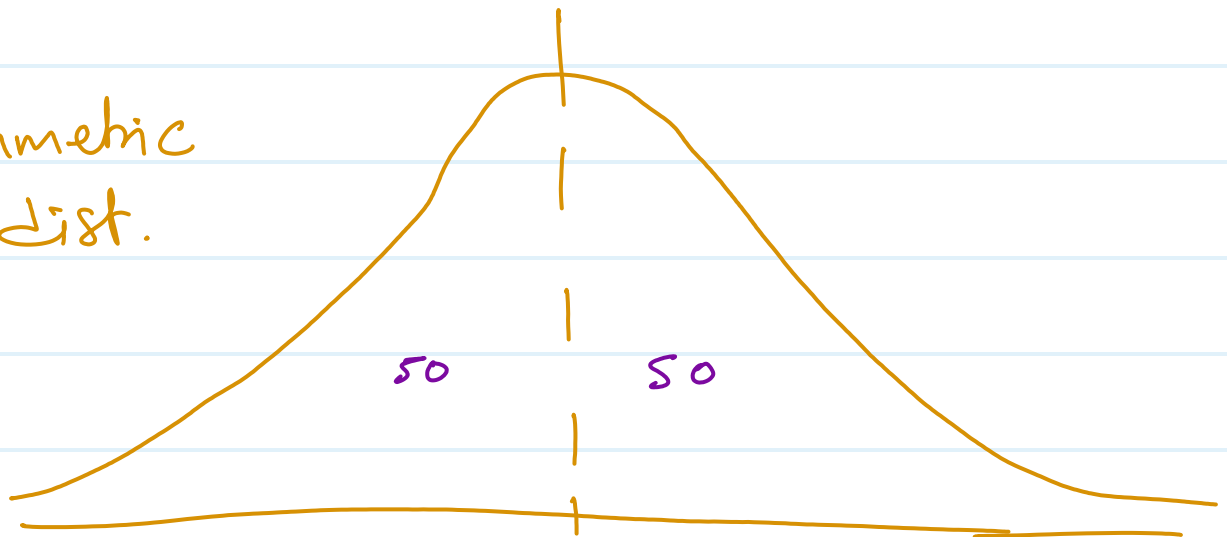


Pdf = Probability density function.

continuous histogram
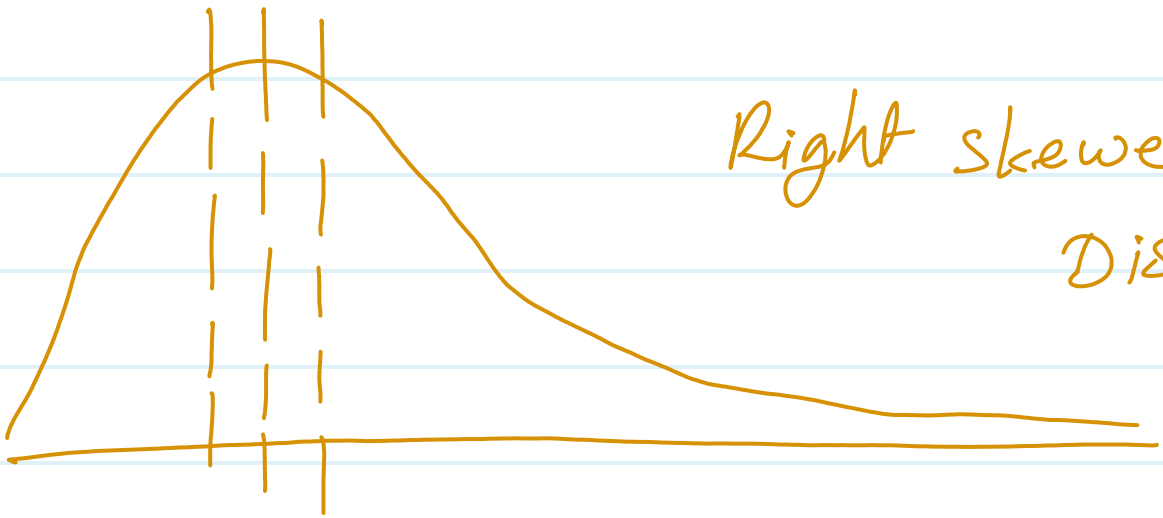
Descret
Histogram

Pmf

$$Pmf = Probability\ mass\ function$$

① Symmetric
dist.

50    50

mean = median = mode

② 

Right skewed Dist.

mean > median > mode

③ 

Left skewed Dist.

mean < median < mode

# Sampling methods

↓

| Probability sampling | Non-probability sampling |
| --- | --- |
| ① Simple Random Samp | ① Convenience Sample. |
| ② clustered samp. | ② consicutive samp. |
| ③ Systematic samp. | ③ Quato sample. |
| ④ Stratified Random samp | ④ Purposive/Judgement Sample. |
| | ⑤ Snowball Sampling. |

## Quartes

# Percentile and Quartile

100%.

$25\%$ percentile $=$ $Q_1$

$50\%$ percentile $=$ $Q_2$/median

$75\%$ —||— $=$ $Q_3$

$100\%$ —||— $=$ $Q_4$

| Score | Rank |
|-------|------|
| 30 | 1 |
| 33 | 2 |
| 43 | 3 |
| 53 | 4 |
| 56 | 5 |
| 67 | 6 |
| 68 | 7 |
| 72 | 8 |

find out where is the 25th percentile is in the above list.

Rank at 25th percentile

$$\# \quad Rank = \frac{Percentile}{100} \times (n+1)$$

$$= \frac{25}{100} \times 9$$

$$= 0.25 \times 9$$

$$= 2.25$$

when rounding up/down the closest value will be Rank.

So Rank is $= 2$

## Rank 75%

$$= \frac{75}{100} \times (n+1)$$

$$= 0.75 \times 9$$

$$= 6.75$$

Rank $= 7$

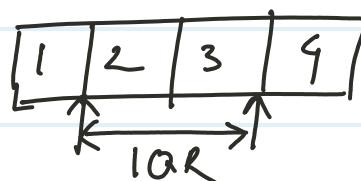\#   $25\% = Q_1$

     $50\% = Q_2 / median$

     $75\% = Q_3$

     $100\% = Q_4$

\# IQR (Inter Quartile Range)

$$IQR = Q_3 - Q_1$$



☆   .5 - Number Summary

①  min

②  $Q_1$

③  median.      IQR

④  $Q_3$.

⑤  max.

$= [1, 2, 3, 4, 6, 8, 11, 14, 18, 19, 5, 21, 82, 95, 7]$

1, 2, 3, 4, 5, 6, 8, 11, 14, 18, 19, 21, 82, 95, 140

To find outlier we use 5 number summery to display values in Box - whisker plot

<u>formula</u> upper limit $= Q_3 + 1.5(IQR)$
lower limit $= Q_1 - 1.5(IQR)$

$$Q_1 = \frac{Q_1}{100} \times (n+1)$$

$$= \frac{25}{100} \times (15+1)$$

$$= \frac{25}{100} \times 16 \qquad =$$

$$Q_1 = 4$$

$$Q_3 = \frac{\overset{3}{75}}{\underset{25}{\cancel{100}}} \times \overset{4}{\cancel{16}}$$

$$Q_3 = 12$$

$$IQR = Q_3 - Q_1$$
$$= 12 - 4$$
$$= 8$$

$$\text{lower limit} = 4 - 1.5 \times 8$$
$$= -8$$

$$\text{upper limit} = 12 + 1.5 \times 8$$
$$= 24$$

$$\text{min} = -8$$
$$\text{max} = 24$$
$$Q_1 = 4$$
$$Q_3 = 12$$

Box and whisker plot

Box



-8   4   12   24
min   Q₁   median   Q₃   max
whisker

82, 95,
140

# To treat outlier, we can use median



① mean
② median
③ mode

1
2
3
4
5
6 → 6.5
7
8
9

Outlier $\dfrac{45+75+80+89}{12} = \dfrac{289}{12} =) \underline{\underline{24}}$

24 | 75 — 6.5
24 | 80 — 6.5
24 | 89 — 6.5

7

1
2

7.5 NAN 15.2

5 →
6
9

7.5 NAN 15.2

11
7.5 NAN 15.2
12

1
2
5
6 → 75
9
11

12
77

15.2

$\dfrac{122}{8}$

$\dfrac{54}{8} =) \underline{\underline{8.8}}$

missing — → mean ← numeric data
         → median mode

P
F
P
F
NA — P
F
NA — P
P
NA — P
P

mode