

Categorical data

classifies individuals or items into different group.

1. Married
2. Widowed
3. Divorced
4. Separated
5. Never married
6. Living w/ Partner
7. Refused
8. Don't know

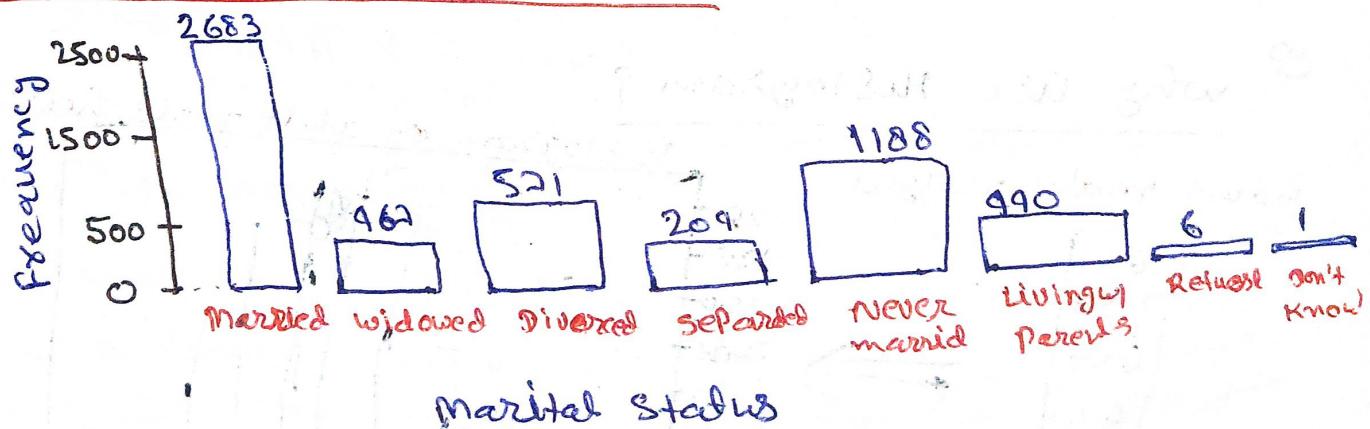
Id	marital status
62229	1
62230	3
62231	1
62232	1
62233	4
....

Frequency Table

- Counts
- Percentage

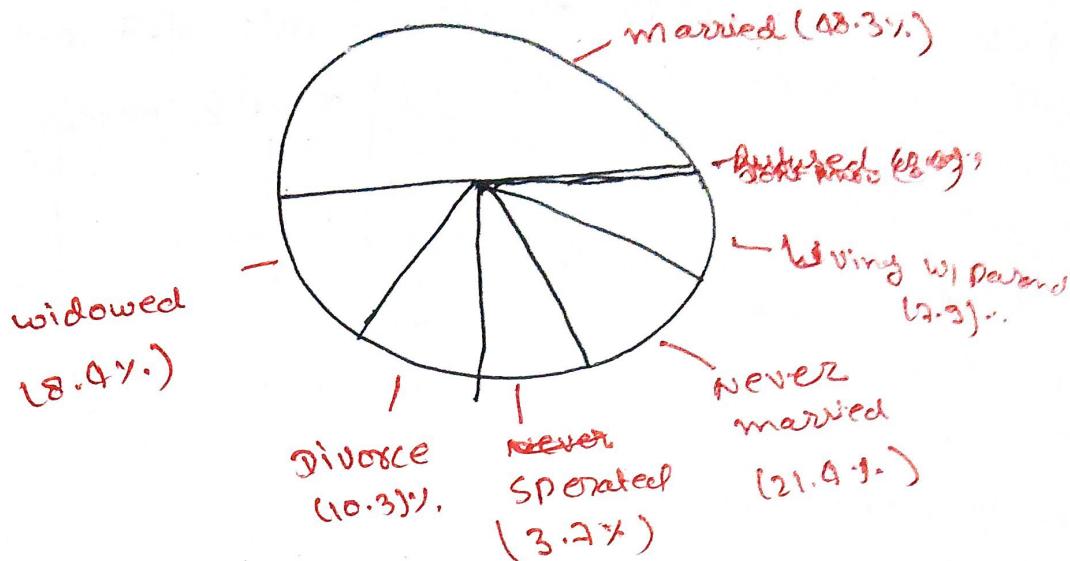
Marital Status	Count	Percent
Married	2683	48.3%
Widowed	462	8.4%
Divorced	571	10.3%
Separated	209	3.7%
Never married	1188	21.0%
Living w/ Partner	490	7.9%
Refused	6	0.1%
Don't know	1	0.0%
Total	5560	100%

Bar chart of marital status



we can also ~~use~~ use percentage instead of
Count.

Pie chart of marital status



→ Frequency Table

↳ Create for numerical summarized

→ Bar chart

↳ Create for visualization

Histogram

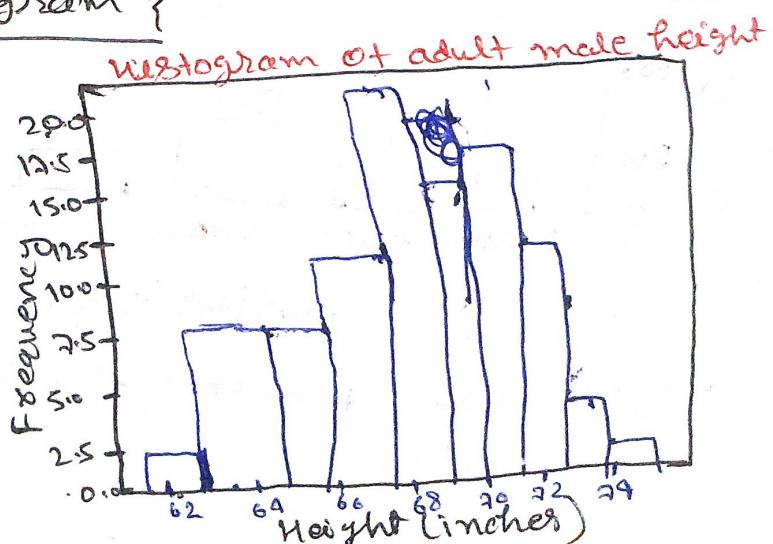
what are Quantitative variables?

Variables that have a numerical value that we can perform mathematical operations on
Ex → Height, weight, income, test scores.

④ why use histogram?

Adult male height

66.3
71.5.1
67.9
67.6
70.0
69.9
64.9.8
—



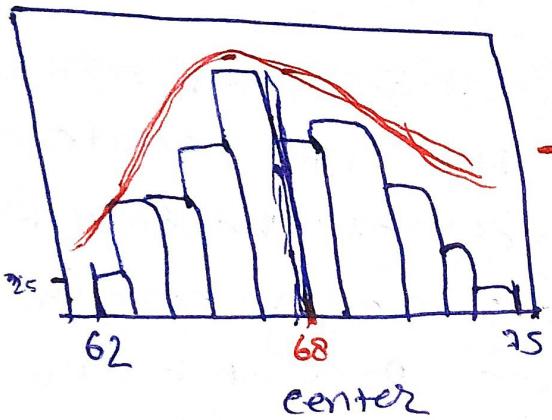
4 main aspect

Shape → overall appearance of histogram. can be symmetric, bell-shaped, left skewed, right skewed etc

center → mean or median

Spread → How far our data spreads. Range, interquartile range (IQR). Standard deviation, variance.

Outliers → Data points that fall far from the bulk of the data.



→ Bell-shaped / unimodal
center ≈ 68 = median
mean = 68

$$\begin{aligned} \text{spread} &= \text{range} = \text{max} - \text{min} \\ &= 75 - 62 \\ &= 13 \end{aligned}$$

Outlier → No apparent outliers

→ The distribution of adult male heights is roughly bell shaped with a center of about 68 inches, a range of 13 inches (62 to 75) and one apparent outlier.

→ Bar chart use for categorical data and Histogram use for quantitative data

Shape →

Right Skewed
Binodal

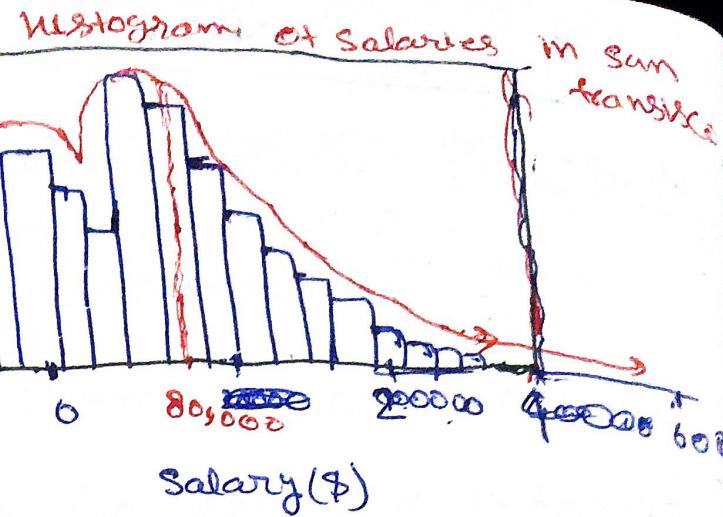
center →

median → 80,000

mean → 85000

Spread →

Range → $60,000 - 0$
= 600000



Outliers

High end

The distribution of salaries in San Francisco Cisco is bimodal and skewed to the right; centered at about \$80,000 with most of the data btw \$40,000 and \$120,000, a range of roughly \$60,000, and outliers are present on the higher end.

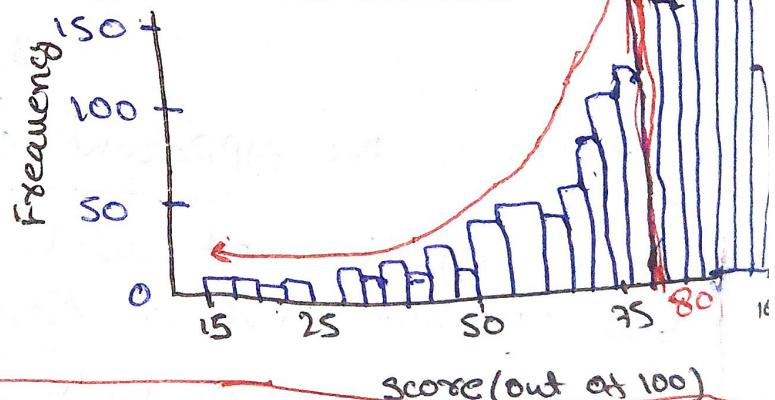
median = 80

Range = $100 - 100$
= 85

Outlier = lower end

Exam score

Histogram of exam score



The distribution of exam score is skewed left centered at about 80 points with most score being btw 65 and 90 points; a range of roughly 85 and some outlier are present below 50 points.

Numerical summaries

Adult male height

5 number summary

- min
- 1st quartile (25%)
- median (50%)
- 3rd quartile (75%)
- max

Height

min. : 61.7

1st Qu. : 66.5

median : 68.3

mean : 68.3

3rd Qua : 70.1

max. : 75.1

$$\begin{aligned} \text{IQR} &= Q_3 - Q_1 \\ &= 70.1 - 66.5 \\ &= 3.6 \end{aligned}$$

Salaries in San Francisco (2011-2017)

~~Number summary~~

min	25%	50%	75%	max	mean	SD	n
618.1	36169	71427	105839	567595	74768	50517	148654

Input error Q_1 meadia Q_3
because salaries \rightarrow IQR
 can't be negative

\downarrow Standard deviation \downarrow Sample size

~~mean + median~~
 boulders exist

Exam score

min	100.0	(10%).
25%	68.0	(68%).
med	78.0	(78%).
75%	87.0	(87%).
max	100.0	(100%).
mean	76.3	(76.3%).
SD	10.4	

$n = 1802.0$

$SD = 10.4$

$$\begin{aligned} \text{IQR} &= 87 - 68 \\ &= 19 \end{aligned}$$

mean < median

\rightarrow Right Skewed

LSI tells us where is most of data falls

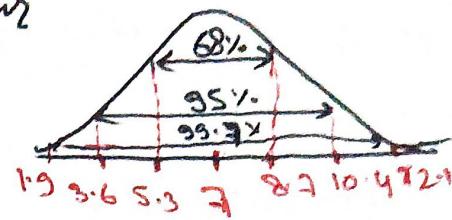
\rightarrow mean is sensitive to extreme observations and the few unusually low scores will tend to bring the mean a bit lower than the median

Amount of sleep (Z=2-2)

$$\rightarrow \mu = 7 \text{ hours}$$

$$\rightarrow \text{standard deviation} (\sigma) = 1.7 \text{ hours}$$

SD \Rightarrow roughly the avg distance of the values from their mean.



Empirical rule (68-95-99.7)

For normal distributions:

$$\rightarrow 1SD (H \pm \sigma) = 68\% \text{ of values}$$

\Rightarrow Range: 5.3 to 8.7 hours

$$\rightarrow 2SD (H \pm 2\sigma) = 95\% \text{ of values}$$

\Rightarrow Range: 3.6 to 10.4 hours

$$\rightarrow 3SD (H \pm 3\sigma) = 99.7\% \text{ of values}$$

Range: 1.9 to 12.1 hours

Very useful to judge how unusual/extreme a value is

Z-score (Standard Score)

\rightarrow Z tells us how unusual they are compared to the mean.

$$\rightarrow \text{Formula} \rightarrow Z = \frac{x - \mu}{\sigma}$$

$x = \text{observation}$
 $\mu = \text{mean}$
 $\sigma = SD$

Interpretation

\rightarrow Positive Z \rightarrow above mean

\rightarrow negative Z \rightarrow below mean

\rightarrow magnitude \rightarrow how far (in SD units) from mean.

$\frac{x - \mu}{\sigma}$

(1) Reed (10 hrs sleep)

$$Z = \frac{10 - 7}{1.7} = 1.76$$

→ above average, fairly unusual (close to upper 95% limit)

(2) Roommate (6 hrs sleep)

$$Z = \frac{6 - 7}{1.7} = -0.59$$

→ slightly below average, not unusual

(3) mark ($Z = -2.7$)

Observation

$$x = \mu + Z\sigma$$

$$= -2.7 \times 1.7 + 7$$

$$= 2.41 \text{ hrs}$$

→ very low sleep, rare situation

Box Plot

What is a Box Plot?

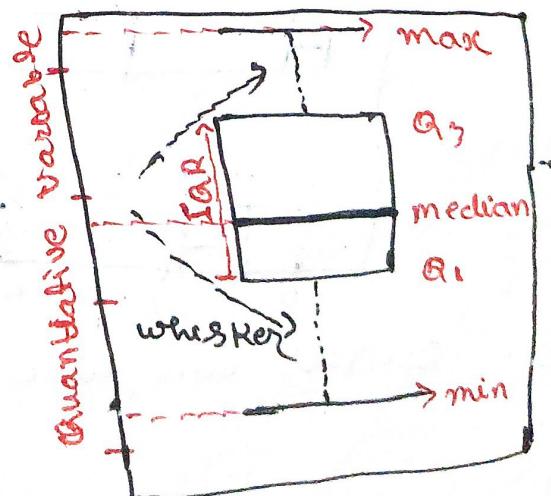
→ graphical representation of the five number summary

Five number summary

min	Q_1	median	Q_3	max
		center		

Range

whiskers: → line extending from box to min and max



Ex-heights of adult males

Height

min. : 61.7

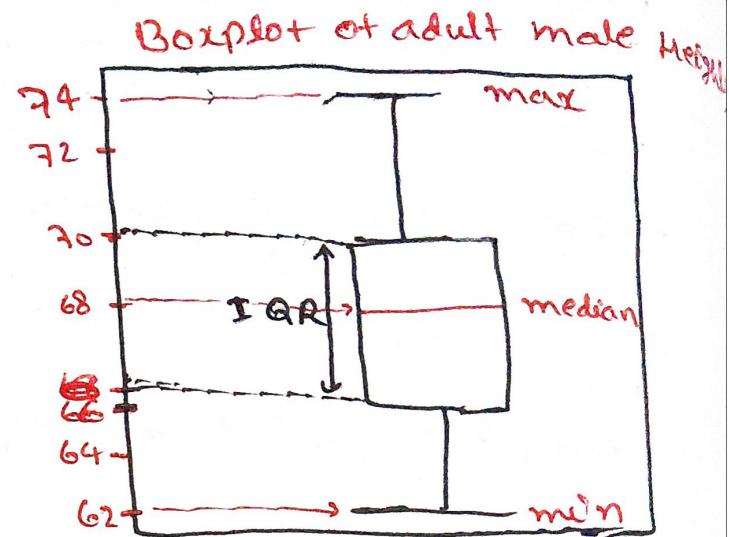
1st QI : 66.5

median : 68.3

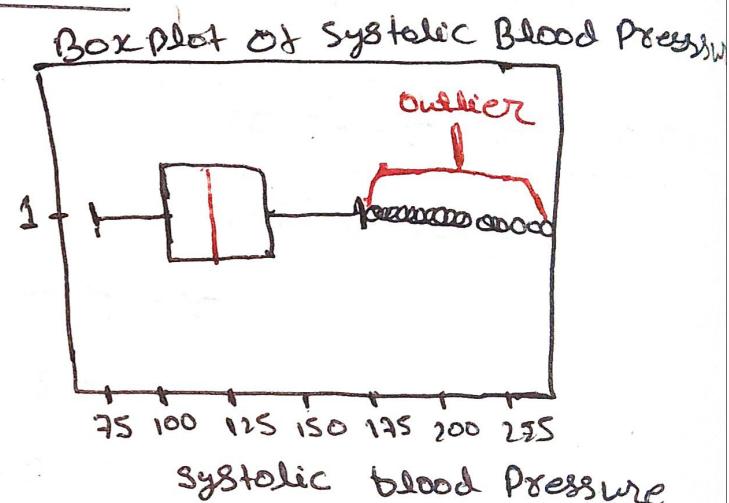
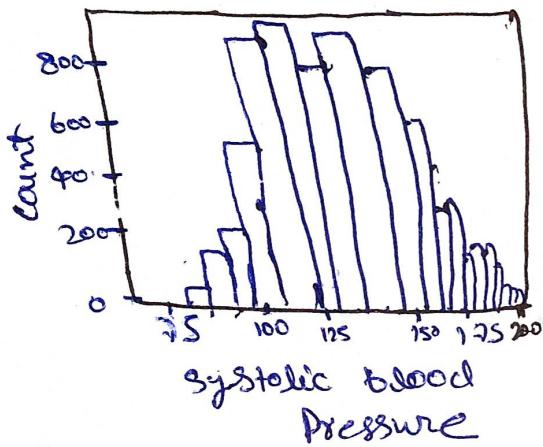
mean : 68.3

3rd QI : 70.1

max : 75.1



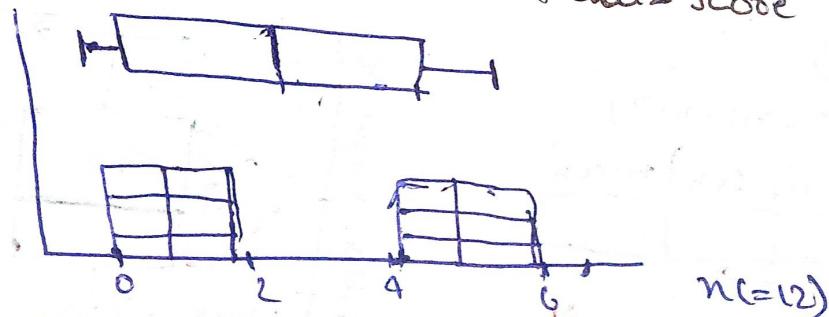
Ex-Systolic Blood Pressure



→ Boxplots can help identify outliers

Ex-Quiz score

Boxplot and Histogram of Quiz Score



Boxplots can hide gaps and clusters

Advantage of Boxplot

→ summarizes data quickly with five number summary.

→ show:

- center → median
- spread → IQR and overall range
- outliers → plotted separately
- excellent for comparison across groups.

Limitations

→ don't show the exact shape

→ histograms are better for shape

→ Boxplots focus on summary + comparisons