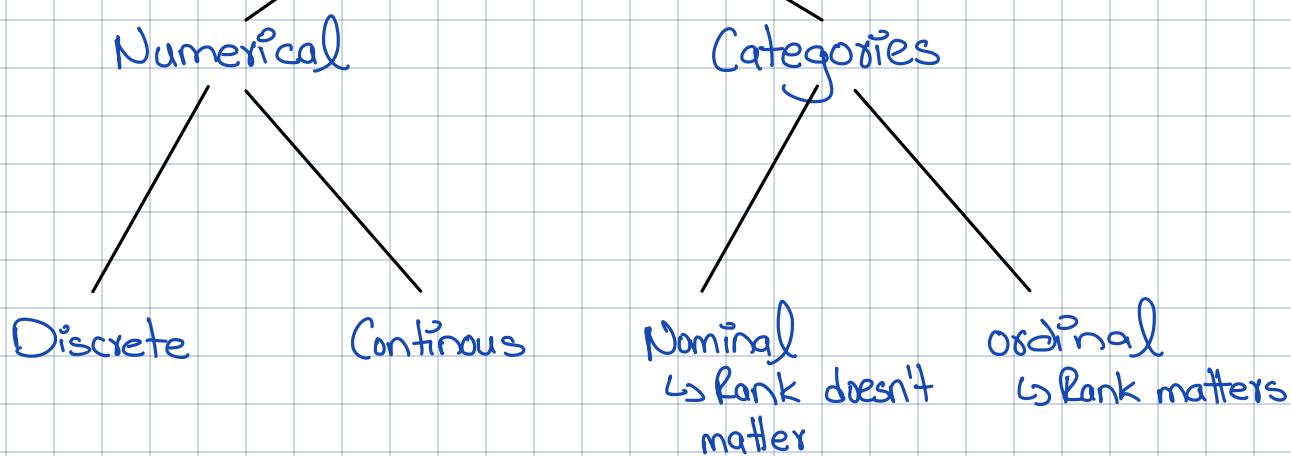


Data Types:



Descriptive Analysis:

↳ Organizes & summarizes data

Eg:

- ↳ Sorting
- ↳ Grouping
- ↳ Creating tables

- ↳ Mean
- ↳ Median
- ↳ Mode

} Central
tendency

- ↳ Variance
- ↳ Range
- ↳ Standard deviation

} Spread of
data

Graphs:

1. Bar chart:

- ↳ Represent categorical data
- ↳ Shows frequency

2. Histogram:

- ↳ Represents numerical data
- ↳ Shows distribution of data over bins.

Summarizing 1D continuous data:

Inferential Analysis:

↳ Make inferences & predictions

Eg:

- ↳ Predicting future sales

1. Mean:

$$\text{mean}(\{x_i\}) = \frac{1}{N} \sum_i^N x_i$$

- ↳ It's the centre of data.
- ↳ Describes central tendency of data.
- ↳ It's sensitive to outliers.

Properties of Mean:

1. Scale:

$$\text{mean}(\{k \cdot x_i\}) = k \cdot \text{mean}(\{x_i\})$$

- ↳ Changes position of points.
- ↳ Relative distance b/w points remain same.
- ↳ Spread of data changes.

2. Translate:

$$\text{mean}(\{x_i + c\}) = \text{mean}(\{x_i\}) + c$$

- ↳ Changes position of points.
- ↳ Spread of data remains same.

3. Signed distances sum to 0:

$$\sum_i^N (x_i - \text{mean}(\{x_i\})) = 0$$

4. Mean minimizes sum of squared distance:

$$\underset{\mu}{\operatorname{argmin}} \sum_i^N (x_i - \mu)^2 = \text{mean}(\{x_i\})$$

Q. Why is there a need to square differences?

- ↳ Emphasize/Penalize larger values.
- ↳ Allows positive contribution from all deviations.

Standard deviation:

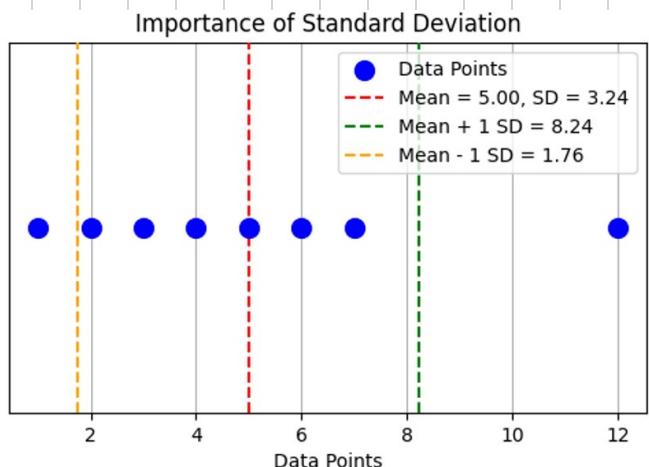
- ↳ Tells spread of w.r.t mean

$$\text{std}(\{x_i\}) = \sqrt{\frac{1}{N} \sum_i^N (x_i - \text{mean}(\{x_i\}))^2}$$

$$= \sqrt{\text{mean}(\{(x_i - \text{mean}(\{x_i\}))^2\})}$$

Q. How do identify outliers in data?

- ↳ Calculate mean
- ↳ Calculate std
- ↳ Range: $(\text{mean} - \text{std}, \text{mean} + \text{std})$
- ↳ Points outside are outliers



Properties of the standard deviation:

1. Scaling:

$$\text{std}(\{k \cdot x_i\}) = |k| \cdot \text{std}(\{x_i\})$$

- ↳ Changes position of points.
- ↳ Relative distance b/w points remain same.
- ↳ Spread of data changes. i.e std changes

2. Translate:

$$\text{std}(\{x_i + c\}) = \text{std}(\{x_i\})$$

- ↳ Changes position of points.
- ↳ Spread of data remains same. i.e std remains same.

Variance:

$$\text{Var}(\{x_i\}) = \frac{1}{N} \sum_{i=1}^N (x_i - \text{mean}(\{x_i\}))^2$$

Properties of Variance:

1. Scaling:

$$\text{Var}(k \cdot \{x_i\}) = k^2 \cdot \text{Var}(\{x_i\})$$

2. Translating:

$$\text{Var}(\{x_i + c\}) = \text{Var}(\{x_i\})$$

Standard Coordinates/Normalize:

↳ Allows us to compare two datasets at different ranges.

$$\hat{x}_i = \frac{x_i - \text{mean}(\{x_i\})}{\text{std}(\{x_i\})}$$

↳ Standard coordinates always have:

↳ Mean: 0

↳ StdDev: 1

↳ Var: 1

↳ It is unitless.

Median:

↳ Sort data

↳ If odd:

Middle value

$$\text{Median} = \left(\frac{n+1}{2} \right)^{\text{th}}$$

↳ If even:

Mean of two middle values

$$\text{Median} = \frac{\left(\frac{n}{2} \right)^{\text{th}} + \left(\frac{n}{2} + 1 \right)^{\text{th}}}{2}$$

Properties of Median:

1. Scale:

$$\text{median}(\{k \cdot x_i\}) = k \cdot \text{median}(\{x_i\})$$

2. Translate:

$$\text{median}(\{x_i + c\}) = \text{median}(\{x_i\}) + c$$

Percentile:

↳ k^{th} percentile means $k\%$ of data is smaller than or equal to it.

Total data points

↑

$$P = \frac{n \cdot k}{100} \rightarrow \text{Percentile}$$

If P is in points:

$$P = X_A + (X_B - X_A) \times \text{point}$$

↓ ↓ ↓
 P^{th} value $(P+1)^{\text{th}}$ value Value after decimal point

Interquartile Range (IQR):

- ↳ $\text{IQR} = (75^{\text{th}} \text{ percentile} - 25^{\text{th}} \text{ percentile})$.
- ↳ Tells spread of data between the 75th & 25th percentile.
- ↳ Less sensitive to outliers.
- ↳ Can identify outliers.

Q. What indicates a large OR small IQR in a dataset?

Relative Size

Context:

- ↳ Compare IQR to range of data.
- ↳ IQR of 17 with range 60-95 indicates moderate variation.
- ↳ Range: $95 - 60 = 35$
- ↳ Proportion: $\frac{17}{35} = 0.49$
- ↳ How big is 50% values compared to entire data.
- ↳ 50% of scores take up half of total spread, which indicates moderate variability

Benchmarking:

- ↳ Compare IQR to similar datasets.
- ↳ In education, 10-15 IQR is normal
- ↳ So, 17 seems too large.

Scale of measurement:

1. Unit matters:

- ↳ IQR of 17 is significant spread in middle 50% of the students
- ↳ IQR of 5 may be small in a dataset of heights in cm.

2. Domain knowledge:

- ↳ Use your experience to assess the IQR.
- ↳ IQR of 10 is normal for tests

Comparison with standard deviation:

- ↪ Small IQR, data is clustered
- ↪ std dev: 11, IQR: 17
- ↪ Significant spread of students in middle 50% of the students

Visualization:

- ↪ Use box plots

Properties of iqf:

1. Scale:

$$iqf(\{k \cdot x_i\}) = k \cdot iqf(\{x_i\})$$

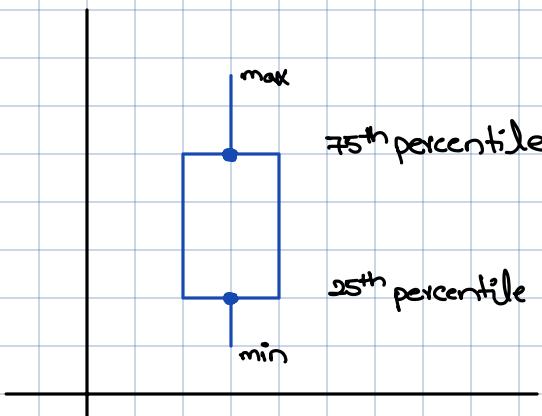
2. Translate:

$$iqf(\{x_i + c\}) = iqf(\{x_i\})$$

Box plot:

↪ Max: 75th percentile + (iqf × 1.5)

↪ Min: 25th percentile - (iqf × 1.5)



↪ Above max & below min are outliers.

↪ If max of data is NOT in range of max calculated, then we use max as the max of data. Similarly, for min.

Example,

Data: [15, 8, 18, 12, 11, 20, 2, 3, 5, 7, 1, 18, 15]

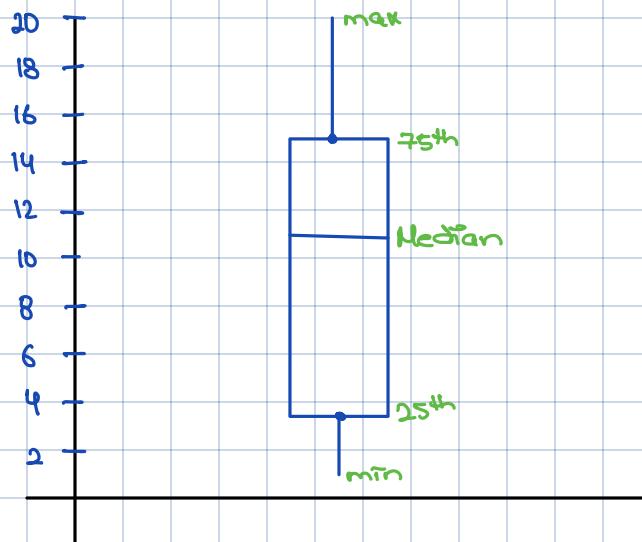
↳ 1 2 3 5 7 8 11 12 15 15 18 18 20
 1 2 3 4 5 6 7 8 9 10 11 12 13

$$P_{25} = \frac{13 \times 25}{100} = 3.25$$

$$P'_{25} = 3 + (5 - 3) \times 0.25 = 3.5$$

Range: 15 - 3.5
11.5

$$\text{Max} = 15 + (11.5 \times 1.5) = 32.25 = 32$$



$$P_{75} = \frac{13 \times 75}{100} = 9.75$$

$$P'_{75} = 15 + (15 - 15) \times 0.75 = 15$$

$$\text{Min} = 3.5 - (11.5 \times 1.5) = -13.75$$

↳ As this is far too less than 1
so we say
= 1

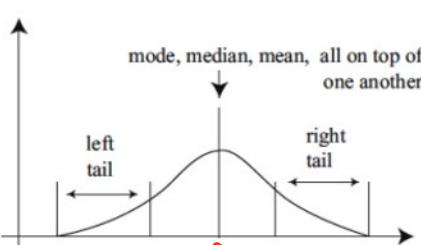
$$\text{Median} = \frac{13 + 1}{2} = 7^{\text{th}} = 11$$

Modes:

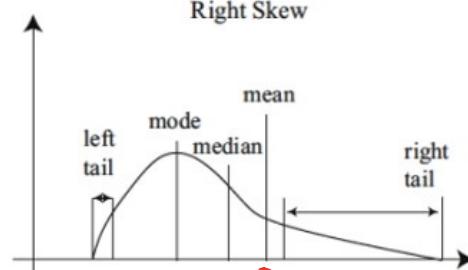
- ↳ The most frequent value.
- ↳ Peaks in histogram.
- ↳ If there are multiple peaks then the data is bimodal & there are sub groups in the population.

Tails and Skews:

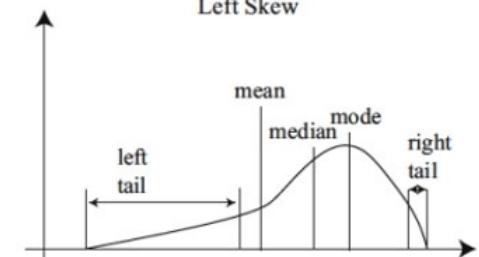
Symmetric Histogram



Right Skew



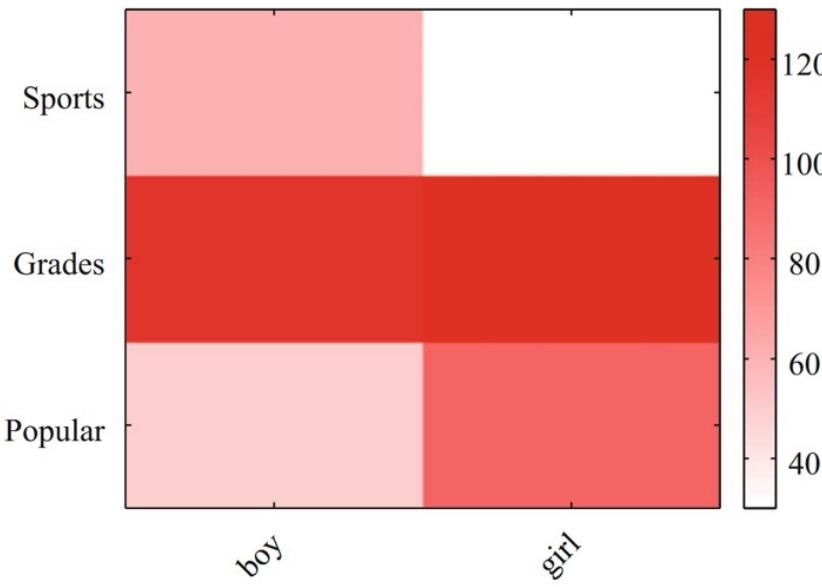
Left Skew



- ↳ Mean is always in the middle.
- ↳ Mode is the peak.
- ↳ Median's in b/w.

Visualizing relationships with a heatmap.

- ↳ Individual values are represented by colors.

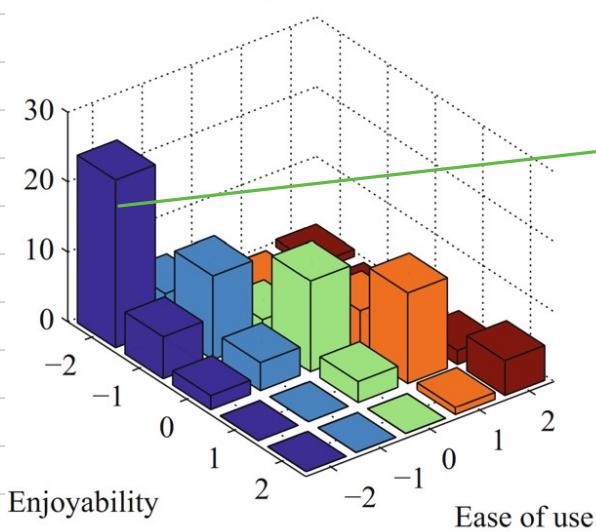


- ↳ Each cell corresponds to count of numbers of elements of that type
- ↳ This map tells:

- ↳ Both boys & girls like grades
- ↳ Boys like sports more
- ↳ Girls like popularity more.

Visualising with a 3D bar chart.

- ↳ Consider a dataset of user rating on enjoyability, ease of use of a software.



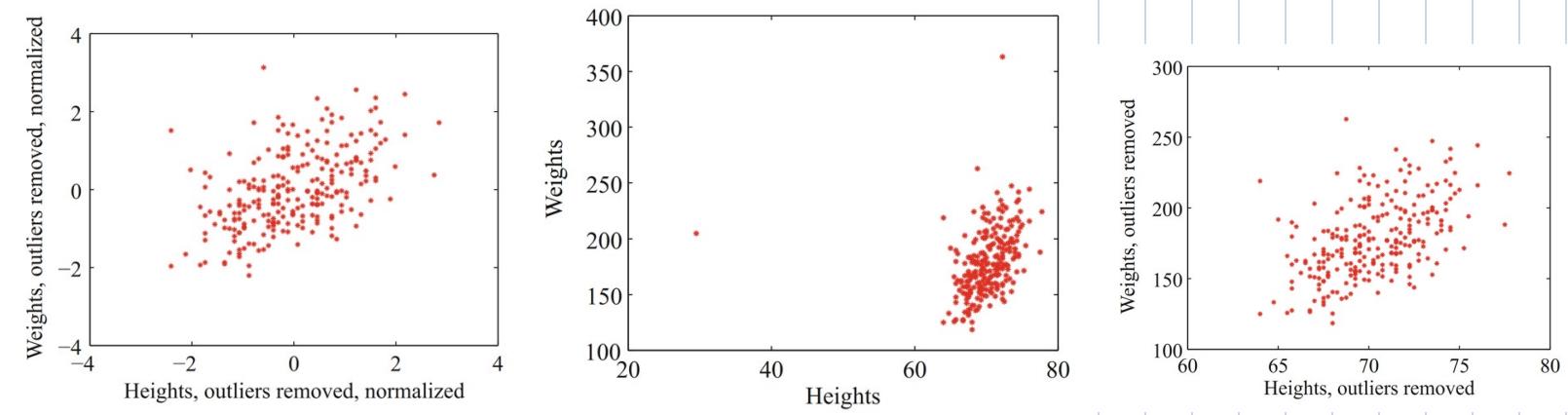
↳ Part of the plot is hidden

→ Those who found it hard to use also didn't enjoy it.

Disadvantages for heatmap & 3D bar plot:

1. For discrete categorical data.
2. Assumes direct relation

Scatter plot:



Correlation

1. Positive correlation:

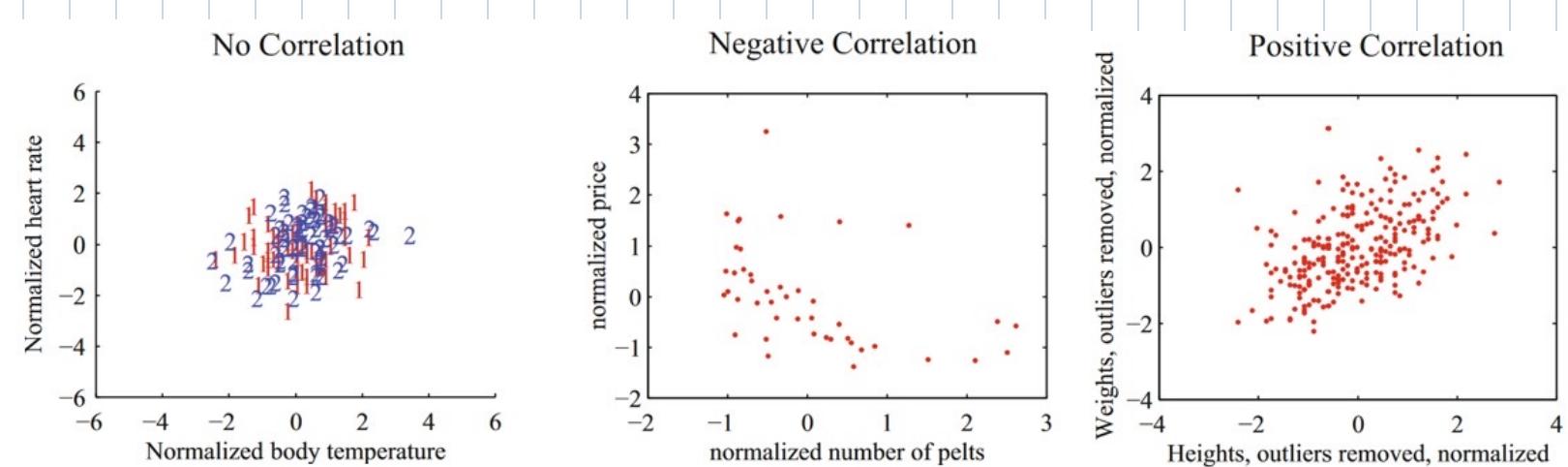
↳ When one variable increases, the other also increases.

2. Negative correlation:

↳ When one variable increases, the other decreases.

3. Zero correlation:

↳ No relation between variables.



↳ Correlation does not mean causation.

Assignment 01:

Task 01:

$$\hat{x} = \frac{x_i - \text{mean}\{x_i\}}{\text{std}\{x_i\}}$$

$$\begin{aligned}\text{mean}\{\hat{x}_i\} &= \sum_{i=1}^n \frac{1}{n} \hat{x}_i \\ &= \frac{1}{n} \sum_{i=1}^n \frac{x_i - \text{mean}\{x_i\}}{\text{std}\{x_i\}} \\ &= \frac{1}{n} \cdot \frac{1}{\text{std}\{x_i\}} \sum_{i=1}^n (x_i - \text{mean}\{x_i\})\end{aligned}$$

$$\begin{aligned}\sum_{i=1}^n (x_i - \text{mean}\{x_i\}) &= \sum_{i=1}^n x_i - \sum_{i=1}^n \text{mean}\{x_i\} \\ &= \sum_{i=1}^n x_i - n \text{mean}\{x_i\} \\ &= \sum_{i=1}^n x_i - \cancel{n} \cdot \cancel{\frac{1}{n} \sum_{i=1}^n x_i} \\ &= \sum_{i=1}^n x_i - \sum_{i=1}^n x_i \\ &= 0 \\ &= \frac{1}{n} \cdot \frac{1}{\text{std}\{x_i\}} \times 0 \\ &= 0\end{aligned}$$

Task 02:

$$d = \{572, 572, 573, 568, 569, 575, 565, 570\}$$

$$1. \text{ mean} = \underline{\underline{572 + 572 + 573 + 568 + 569 + 575 + 565 + 570}}$$

$$= 570.5$$

1 2 3 4 5 6 7 8
 565, 568, 569, 570, 572, 572, 573, 575

$$\text{median} = \frac{\frac{8^{\text{th}}}{2} + \left(\frac{8}{2} + 1\right)^{\text{th}}}{2}$$

$$= \frac{570 + 572}{2} = 571$$

2. $\text{Var} = 8.75$

$$\text{std} = \sqrt{8.75} = 2.95$$

Range: 575 - 565

$$= 10$$

4. The tires seem to be of good quality. Most of the tires lie around 570 with small std of 2.95
3. The data relatively balanced but median > mean shows presence of smaller values \rightarrow slightly skewed.

Task 3:

~~919 1196 785 1126 936 918 1156 920 948 1027 1052 1117 1120 929 905 921 1035 1045 835 1195 1195 1240~~
~~1112 938 970 1227 956 1182 1117 978 882 1049 1117 1111 1049 765 938 902 1022 1333 811 1227 1045 836 958~~
~~1311 1037 702 929 33~~

1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16
 702, 765, 785, 811, 832, 855, 896, 902, 905, 918, 919, 920, 923, 929, 936, 938,

17 18 19 20 21 22 23 24 25 26 27 28 29
 948, 950, 956, 958, 958, 970, 972, 973, 1009, 1009, 1022, 1035, 1037,

30 31 32 33 34 35 36 37 38 39 40 41
 1045, 1067, 1085, 1092, 1102, 1122, 1126, 1151, 1156, 1157, 1162, 1170,

42 43 45 46 47 48 49 50
 1195, 1195, 1196, 1217, 1237, 1311, 1333, 1340

$$P_{25} = \frac{25 \times 50}{100}$$

$$= 12.5$$

$$P'_{25} = 920 + (923 - 920) \times 0.5$$

$$= 921.5$$

$$\text{Range: } 1153.5 - 921.5$$

$$= 232$$

$$\min: 921.5 - 232 \times 1.5$$

$$= 573$$

$$= 702$$

$$P_{75} = \frac{75 \times 50}{100}$$

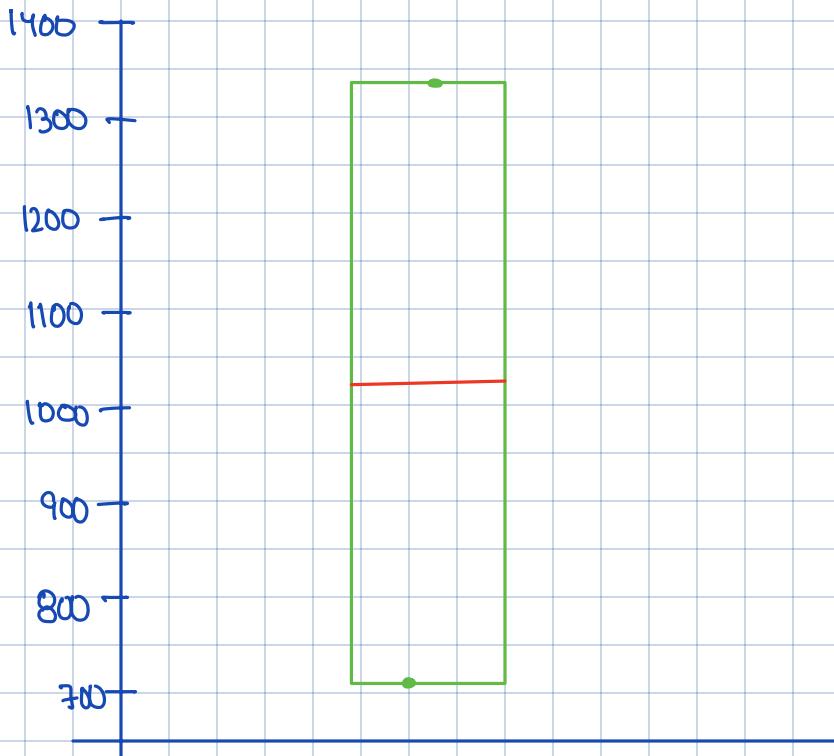
$$= 37.5$$

$$P'_{75} = 1151 + (1156 - 1151) \times 0.5$$

$$= 1153.5$$

$$\max = 1153.5 + (232 \times 1.5)$$

$$= 1340$$



$$\text{Median} = \frac{\frac{50^{\text{th}}}{2} + \left[\frac{50}{2} + 1 \right]^{\text{th}}}{2}$$

$$= \frac{1009 + 1009}{2}$$

↪ The boxplot tells that there is quite a lot of variation in the data.

Task 4:

6.72 6.87 6.82 6.70 6.88 6.70 6.82 6.85 6.86 6.66 6.64 6.76 6.73 6.80 6.72 6.76 6.76 6.88 6.86 6.62 6.72 6.76 6.70 6.78
 6.76 6.77 6.70 6.72 6.74 6.81 6.79 6.78 6.86 6.76 6.76 6.72

1
 6.64, 6.62, 6.62, 6.66, 6.66, 6.66, 6.66, 6.67, 6.68, 6.70, 6.70, 6.70, 6.70,
 6.72, 6.72, 6.72, 6.72, 6.72, 6.73, 6.74, 6.75, 6.76, 6.76, 6.76, 6.76,
 6.76, 6.76, 6.76, 6.77, 6.78, 6.78, 6.78, 6.79, 6.80, 6.81, 6.82

$n = 36$

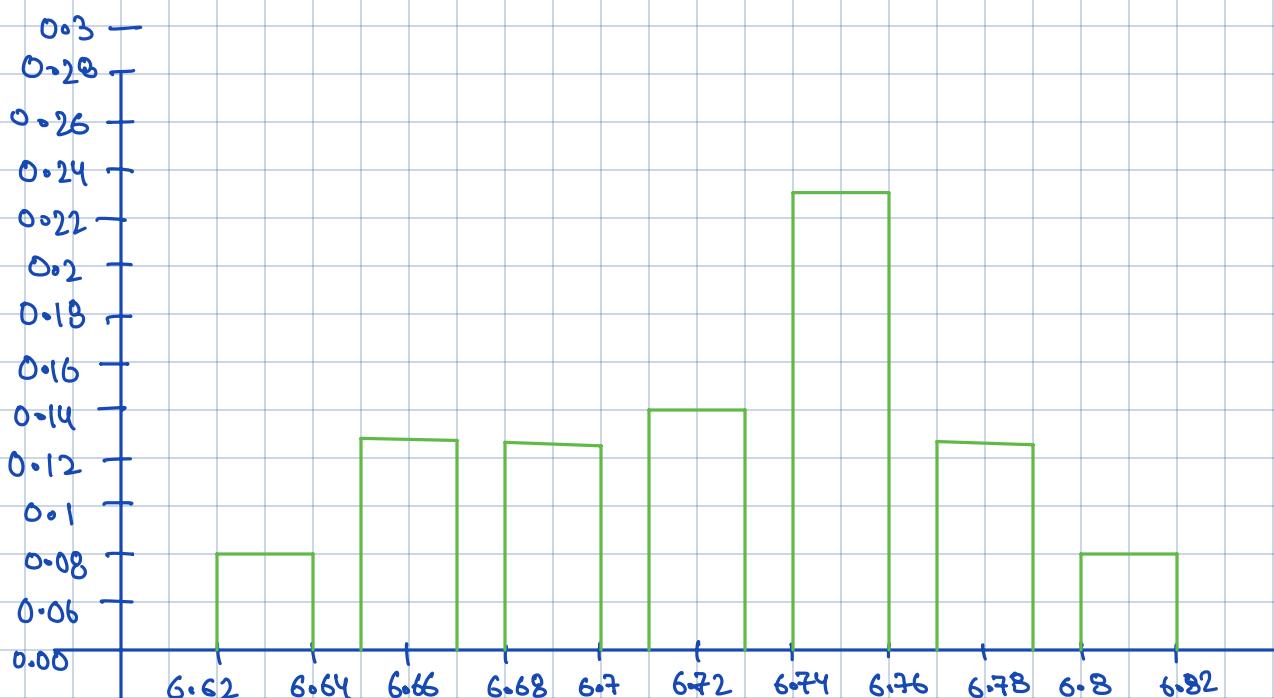
mean : 6.726

std : 0.0536

$$\text{width} = \frac{6.82 - 6.64}{6}$$

$$= 0.03$$

3
 5
 5
 6
 9
 6.62 → 6.64, 6.65 - 6.67, 6.68 - 6.70, 6.71 - 6.73, 6.74 - 6.76,
 6.77 - 6.79, 6.80 - 6.82



3. Almost symmetric but with some skewness.

Properties of correlation:

↳ Translating the data does not effect the correlation

↳ Scaling the data may change the sign of correlation

$$\text{corr}(\{(ax_i + b, cy_i + d)\}) = \text{sign}(a/c) \text{corr}(\{(x_i, y_i)\})$$

↳ Correlation is bounded between [-1, 1]

$$\text{corr}(\{(x_i, y_i)\}) = 1 \text{ if } \hat{x}_i = \hat{y}$$

$$\text{corr}(\{(x_i, y_i)\}) = -1 \text{ if } \hat{x}_i = -\hat{y}$$

Predictions using correlation:

1. $\text{mean}(\{\hat{u}\}) = 0$

$$\begin{aligned} \text{mean}(\{\hat{u}\}) &= \text{mean}(\{\hat{y} - \hat{y}_p\}) \\ &= \text{mean}(\{\hat{y}\}) - \text{mean}(\{\hat{y}_p\}) \\ &= \text{mean}(\{\hat{y}\}) - \text{mean}(\{a\hat{x} + b\}) \\ &= \text{mean}(\{\hat{y}\}) - a \text{mean}(\{\hat{x}\}) - b \\ &= 0 - a(0) + b \end{aligned}$$

$\therefore b = 0$

2. $\text{Var}(\{\hat{u}\}) = \text{Var}(\{\hat{y} - \hat{y}_p\})$

$$\begin{aligned} &= \text{Var}(\{\hat{y} - a\hat{x} + b\}) \\ &= \text{Var}(\{\hat{y} - a\hat{x}\}) \\ &= \text{mean}(\{\hat{y} - a\hat{x}\}^2) \\ &= \text{mean}(\{\hat{y}^2 - 2a\hat{x}\hat{y} + a^2\hat{x}^2\}) \\ &= \text{mean}(\{\hat{y}^2\}) - 2a\text{mean}(\{\hat{x}\hat{y}\}) + a^2\text{mean}(\{\hat{x}\}^2) \\ &= 1 - 2ax + a^2 \end{aligned}$$

$$\frac{d \text{var}(\{y_i\})}{da} = -2s + 2a$$

$$-2s + 2a = 0$$

$$a = s$$

$$\hat{y}_p = a\hat{x} + b \\ = s\hat{x}$$

Procedure:

1. Find \hat{x} & \hat{y}

2. Find $\text{corr}(\{\hat{x}_i, \hat{y}_i\}) = \frac{1}{N} \sum_{i=1}^N \hat{x}_i \hat{y}_i$

3. $y_i^p = \text{std}(y) \cdot s \cdot \hat{x} + \text{mean}(\{y_i\})$