

---

---

---

---

---



# VARIATIONAL FORMULATION OF MEAN & MEDIAN

---

$$\text{mean} = \arg \min_a g(a)$$

$$\text{median} = \arg \min_a f(a)$$

$$g(a) = \frac{1}{n} \sum_{i=1}^n (x_i - a)^2$$

(L2 norm)

$$f(a) = \frac{1}{n} \sum_{i=1}^n |x_i - a|$$

(L1 norm)

$$g(a) = \frac{1}{n} \sum_i (x_i - a)^2$$

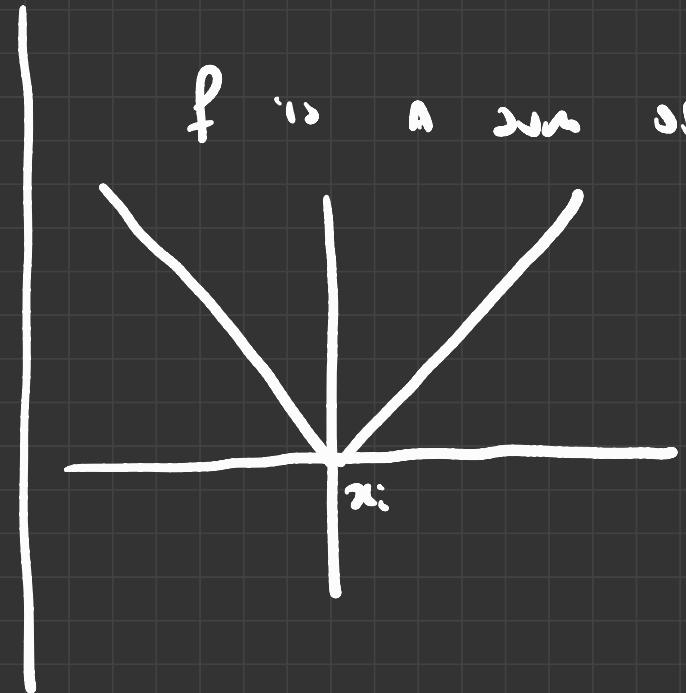
$$g'(a) = \cancel{\frac{2}{n} \sum_{i=1}^n (-)(x_i - a)} = 0$$

$$\Rightarrow \sum_i x_i = n a \Rightarrow a = \underline{\underline{\frac{1}{n} \sum_i x_i}}$$

WHAT ABOUT

$$f(a) = \frac{1}{n} \sum_i |x_i - a|$$

$f$  is a sum of





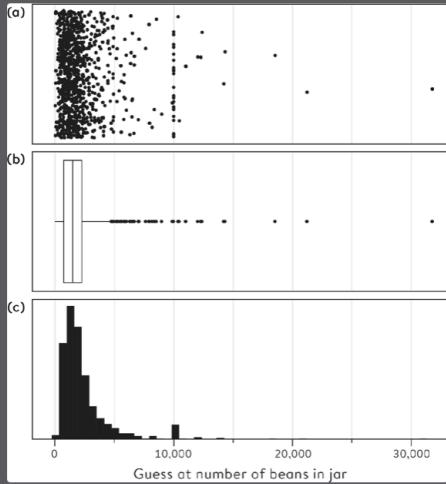
# MEAN

vs

# MEDIAN

# (INFLUENCE vs

# ROBUSTNESS)



MEAN : 2408

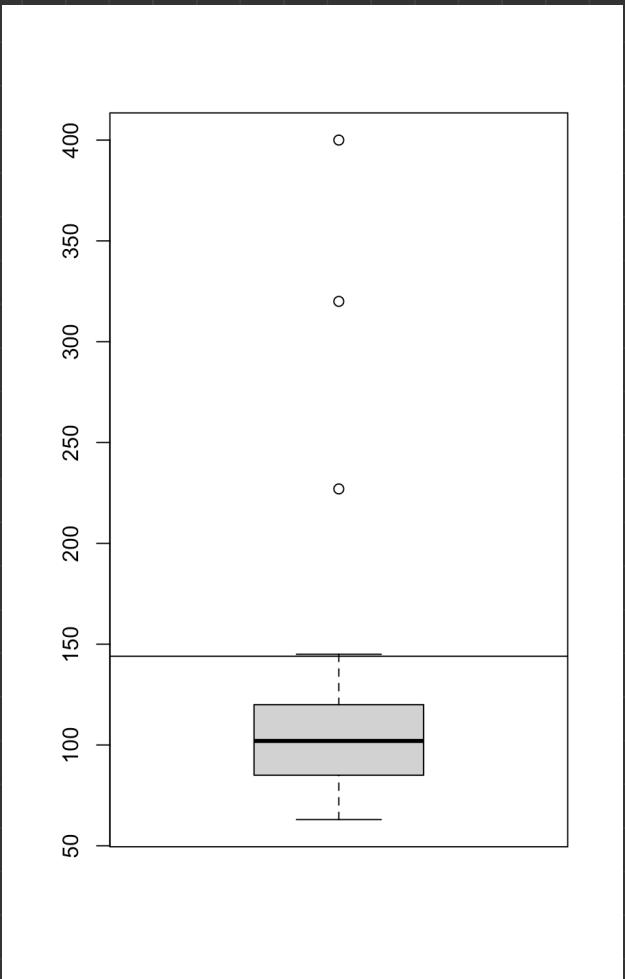
MEDIAN : 1775

TRIM : 1616

Figure 2.2

Different ways of showing the pattern of 915 guesses of the number of jelly beans in the jar. (a) A strip-chart or dot-diagram, with a jitter to prevent points lying on top of each other; (b) a box-and-whisker plot; (c) a histogram

# Pi, r6 - Pov6 RESULTS



# MEASURES / SUMMARIES OF SPREAD / DISPERSION / VARIABILITY

---

- RANGE : MAX - MIN  
↳ VERY UNINFORMATIVE
- INTERQUARTILE RANGE : IQR Q - IQR Q

## - STANDARD DEVIATION

$$SD = \sqrt{VAF}$$

$$VAF = \frac{1}{n} \sum_{i=1}^n (x_i - \text{mean})^2 = \min_a \frac{1}{n} \sum_i (x_i - a)^2$$

## - MEAN ABSOLUTE DEVIATION

$$MAD = \frac{1}{n} \sum_i |x_i - \text{median}| = \min_a \frac{1}{n} \sum_i |x_i - a|$$

## HISTOGRAMS

THIS IS NOT A SUMMARY OR VISUALIZATION  
OF THE DATA.

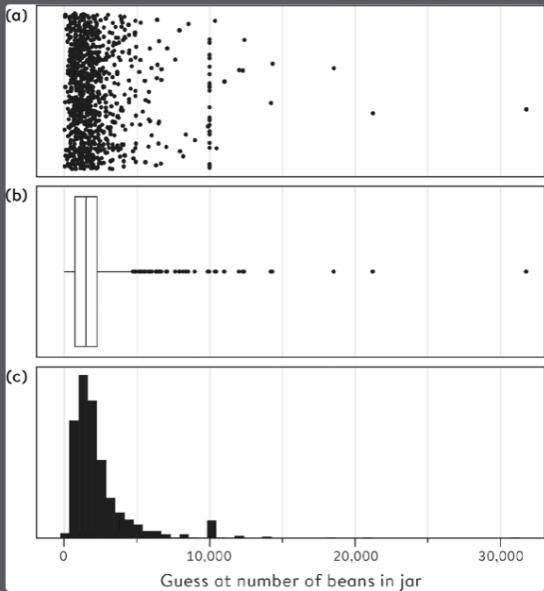
↳ IT IS AN INFERENCE ABOUT DATA

## CREATING A HISTOGRAM

1st Step : DEFINE DATA INTERVALS  
↳ BINS

2nd Step : FIND THE NUMBER OF DATA  
IN EACH BIN

3rd Step : BAR PLOT

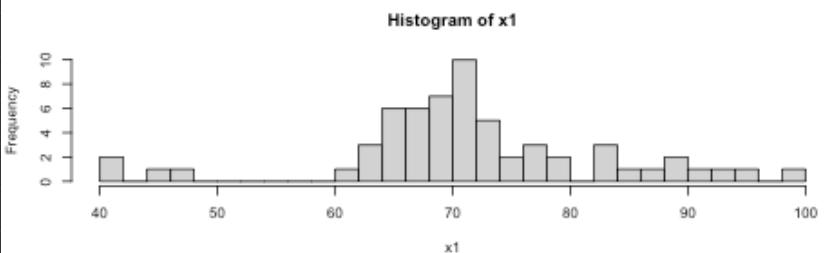
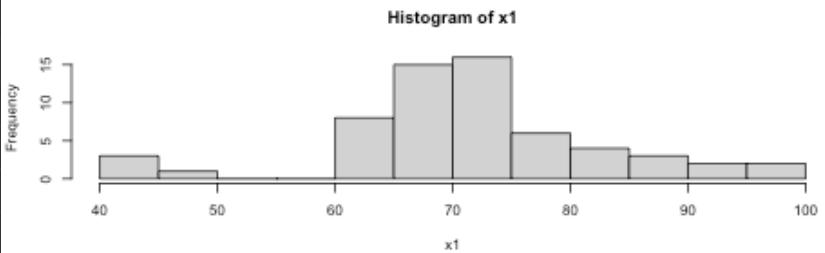
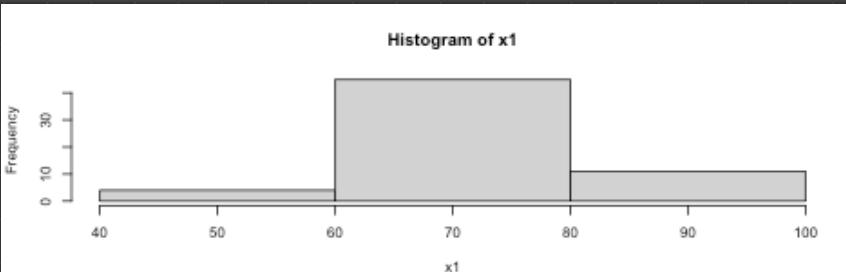


Histogram

**Figure 2.2**

Different ways of showing the pattern of 915 guesses of the number of jelly beans in the jar. (a) A strip-chart or dot-diagram, with a jitter to prevent points lying on top of each other; (b) a box-and-whisker plot; (c) a histogram

4 MODES  
BIMODAL  
(ATTENDING  
CLASS)



## BIVARIATE DATA

( $x_i, y_i$ )       $i = 1, \dots, n$

$x, y$  NUMERICAL

AoS

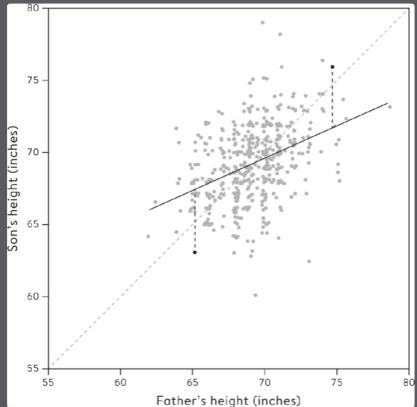


Figure 5.1  
Scatter of heights of 465 fathers and sons from Galton's data (many fathers are repeated since they have multiple sons). A jitter has been added to separate the points, and the diagonal dashed line represents exact equality between son and father's heights. The solid line is the standard 'best-fit' line. Each point gives rise to a 'residual' (dashed line), which is the size of the error were we to use the line to predict a son's height from his father's.

SCATTER PLOT

IGNORE SOLID BLACK

LINE

"REGRESSION TO THE MEAN"

## CORRELATION COEFFICIENT

MEASURE IF HOW "TIGHTLY" THE SCATTER PLOT  
IS AROUND A STRAIGHT LINE

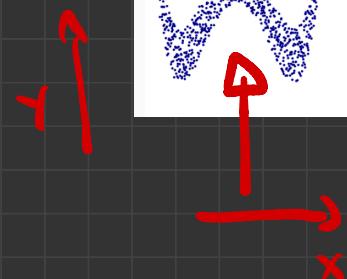
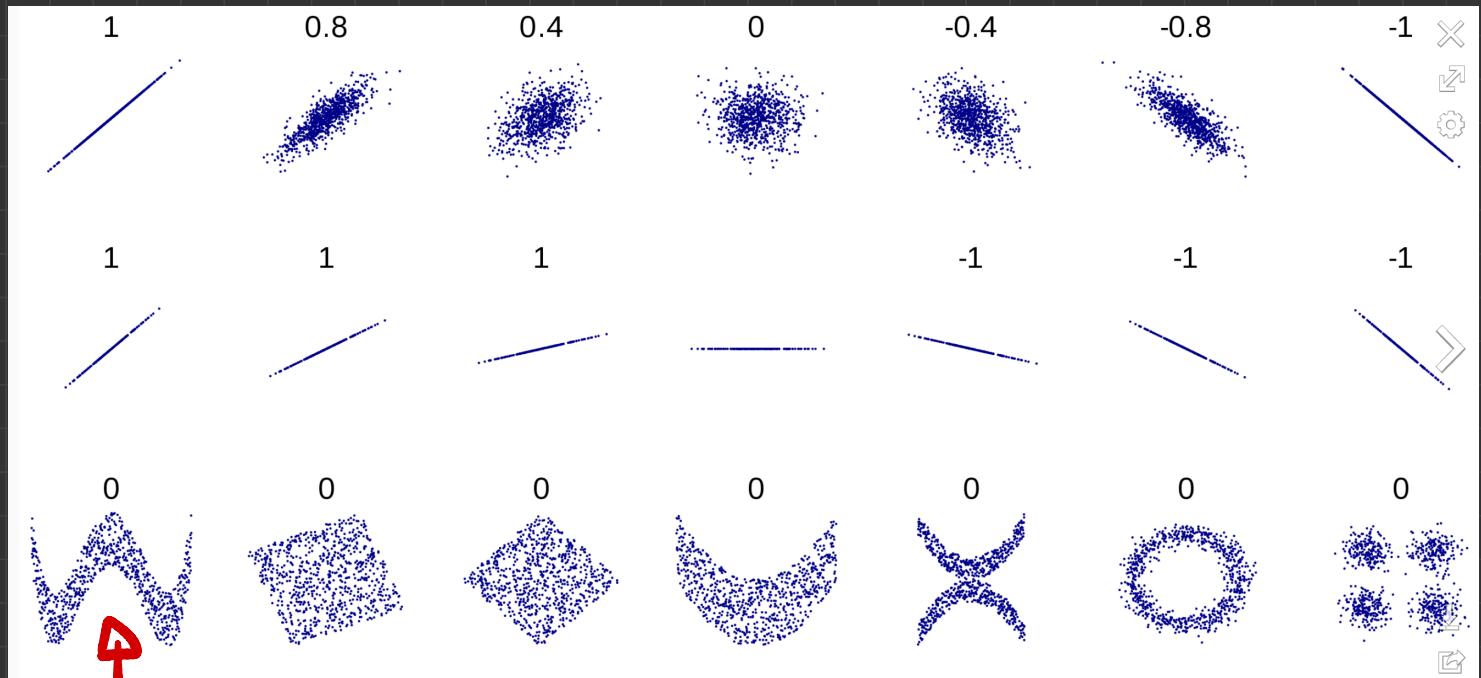
$$(x_i, y_i) \quad i = 1, \dots, n$$

$$\text{mean}(x) = \frac{1}{n} \sum_i x_i \quad \text{mean}(y)$$

$$sd(x) = \sqrt{\frac{1}{n} \sum_i (x_i - \text{mean}(x))^2} \quad sd(y)$$

$$\text{corr}(x,y) = \frac{1}{n} \sum_{i=1}^n \left( \frac{x_i - \text{mean}(x)}{\text{sd}(x)} \times \frac{y_i - \text{mean}(y)}{\text{sd}(y)} \right)$$

- $\frac{x_i - \text{mean}(x)}{\text{sd}(x)}$  : Z-score on the standardization of data points
- Definition implies :  $-1 \leq \text{corr}(x,y) \leq 1$   
↳  $\text{corr}(x,y)$  is a standardized measure of dependence

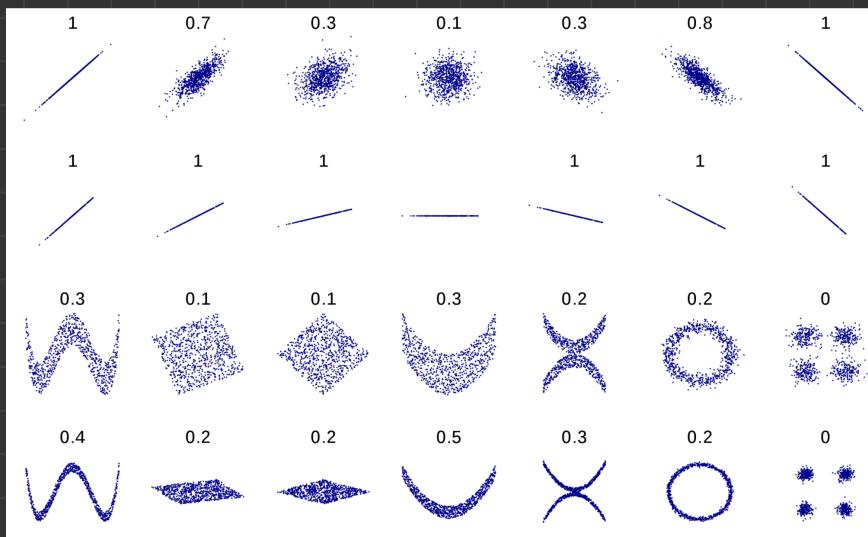


# DIFFERENT MEASURES OF (NON-LINEAR) DEPENDENCE

---

## DISTANCE CORRELATION

---



AoS



Color = 0

