

Chebyshev's rule applies to any sample of measurements, regardless of the shape of the frequency distribution:

At least  $1 - \frac{1}{k^2}$  of the measurements will fall within  $k$  std deviations of the mean  $(\bar{x} - ks, \bar{x} + ks)$  for any <sup>whole</sup> number  $k > 1$ .

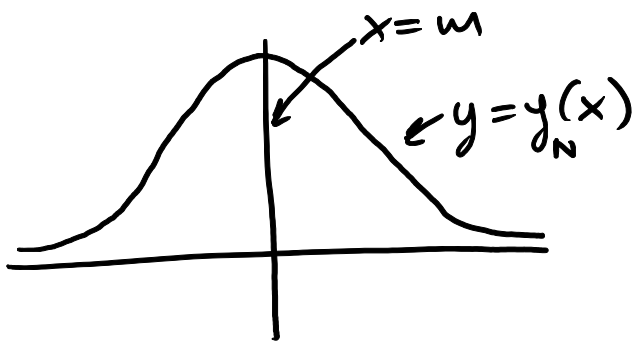
Empirical rule is a rule of thumb that applies to samples with freq distributions which are bell shaped (or "mound shaped" or "approximately normal").

68% of measurements in  $(\bar{x} - s, \bar{x} + s)$

95% of measurements in  $(\bar{x} - 2s, \bar{x} + 2s)$

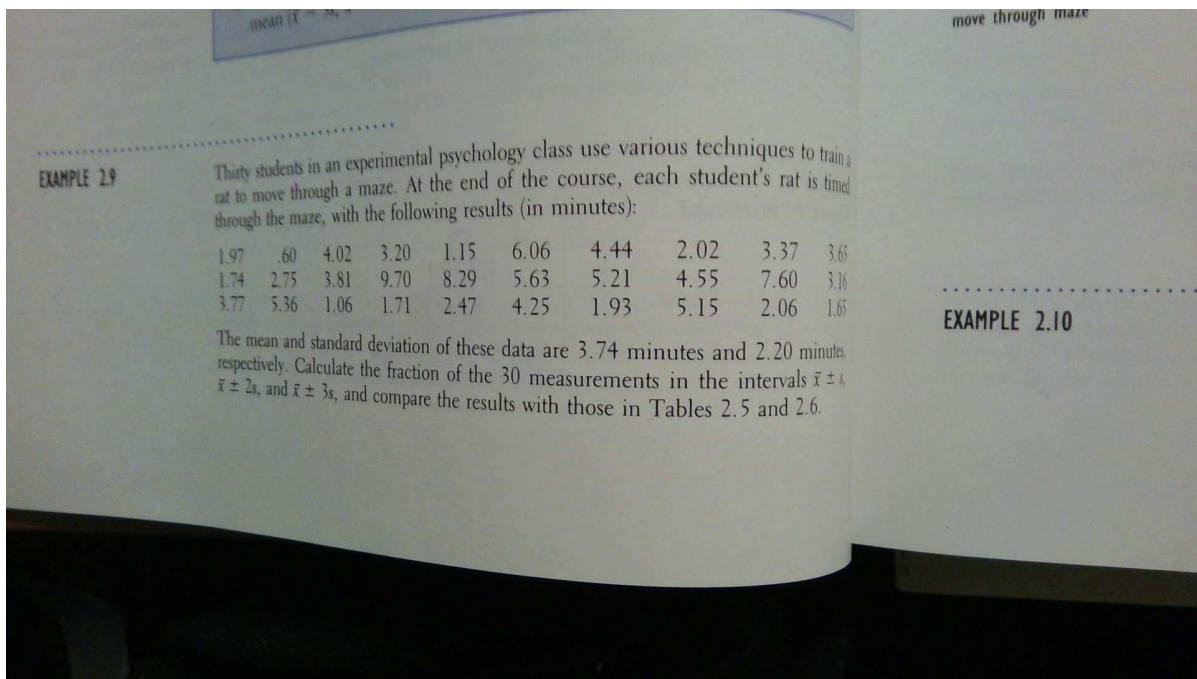
99.7% of measurements in  $(\bar{x} - 3s, \bar{x} + 3s)$

Derived from integral of  $y(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[-\frac{(x-\mu)^2}{2\sigma^2}\right]$



$$\int_{-\infty}^{+\infty} y_N(x) dx = 1$$

Consider an example with the data below:



Step 1: construct freq distribution and histogram.

$$\text{min} = 0.60$$

$$\text{max} = 9.70$$

Decide on class width  
say 1.5

$$\text{class 1: } [0.60, 2.09]$$

$$\text{class 2: } [2.10, 3.59]$$

$$\text{class 3: } [3.60, 5.09]$$

class 4 : [ 5.10, 6.59 ]

class 5 : [ 6.60, 8.09 ]

class 6 : [ 8.10, 9.59 ]

class 7 : [ 9.60, 11.09 ]

Use R to get counts

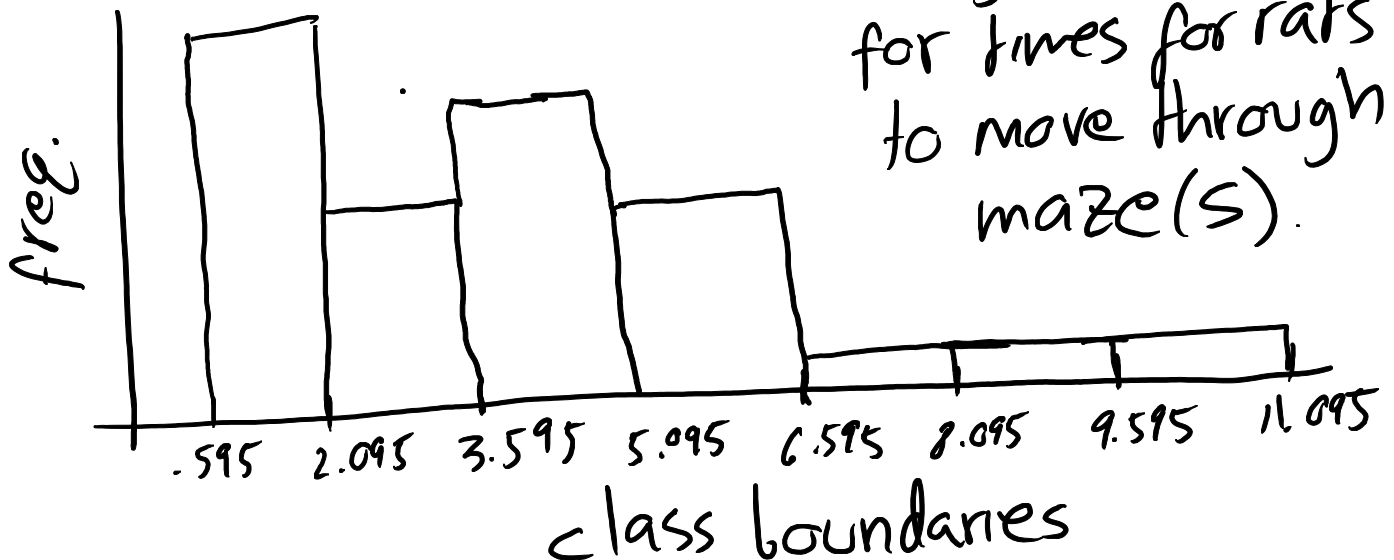
class	freq
class 1	10
class 2	5
class 3	7
class 4	5
class 6	1
class 7	1
class 8	1

Notice that this necessarily covers all the data as data is given to two decimal figures.

construct class boundaries:

0.595, 2.095, 3.595,  
5.095, 6.595, 8.095,  
9.595, 11.095

diffs are 1.5 = class width



$$\begin{array}{ll} \text{I: } (\bar{x}-s, \bar{x}+s) = (1.54, 5.94) & 23/30 \approx 77\% \\ \text{II: } (\bar{x}-2s, \bar{x}+2s) = (-.66, 8.14) & 28/30 \approx 93\% \\ \text{III: } (\bar{x}-3s, \bar{x}+3s) = (-2.86, 10.34) & 30/30 = 100\% \end{array}$$

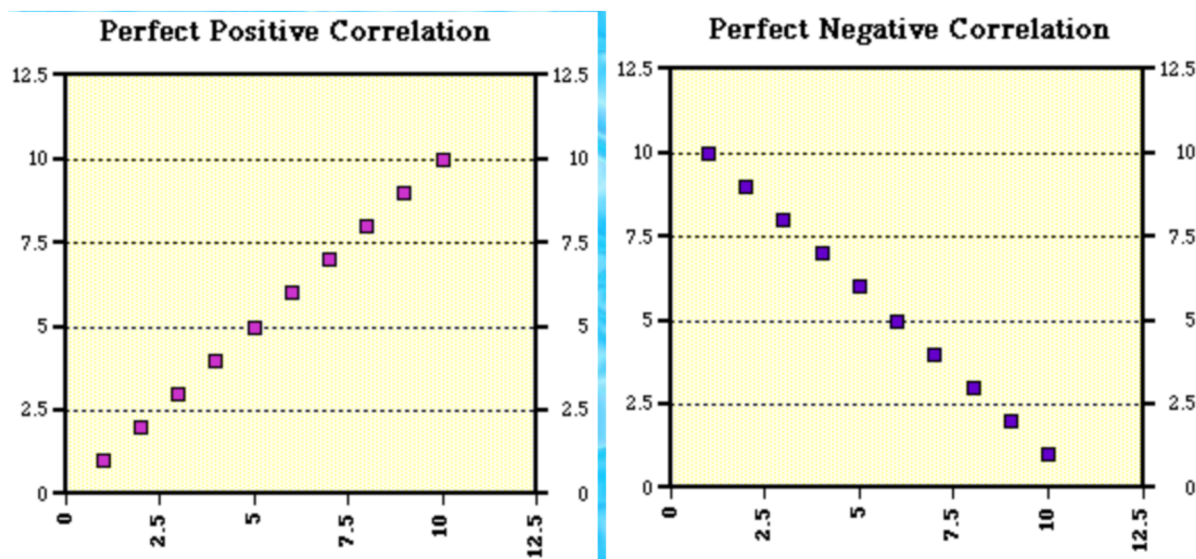
Empirical rule: (68, 95, 99.7)

Good estimate using the empirical rule.

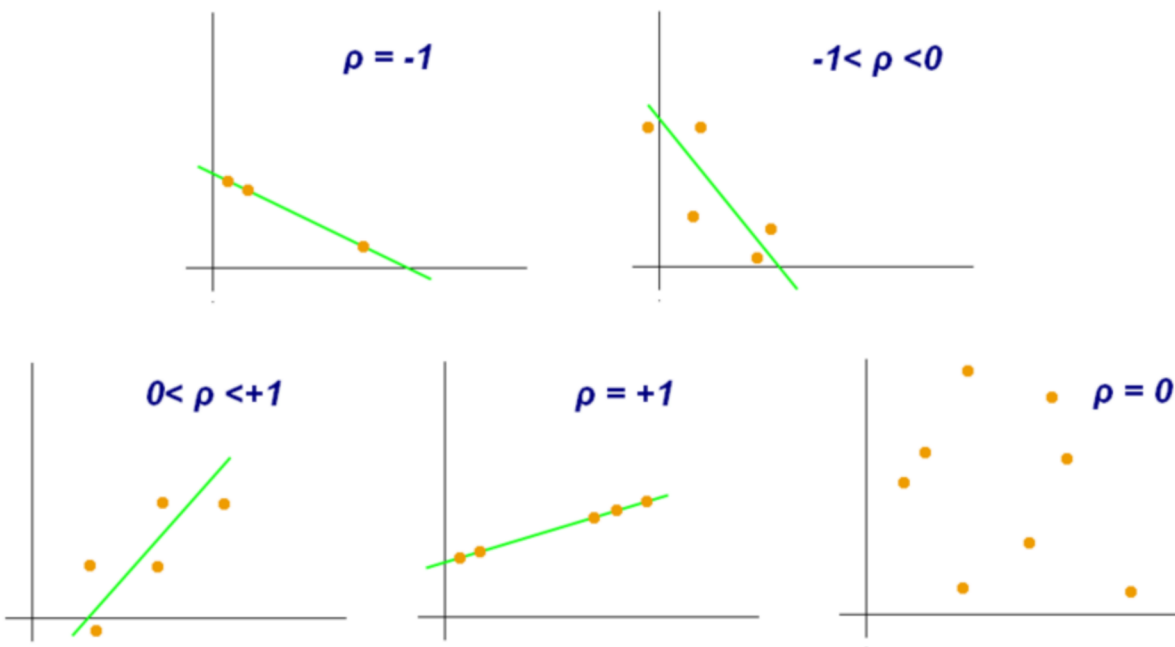
By Chebyshev's rule we expect at least 75%, 89% of measurements to lie within II and III, which holds above. Note that for interval I, Chebyshev's theorem can't be used.

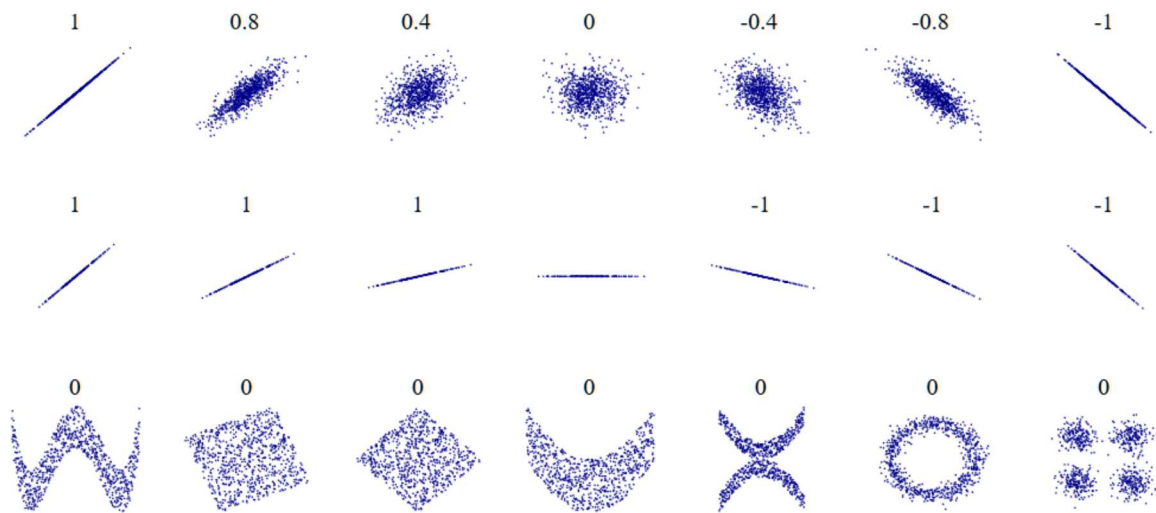
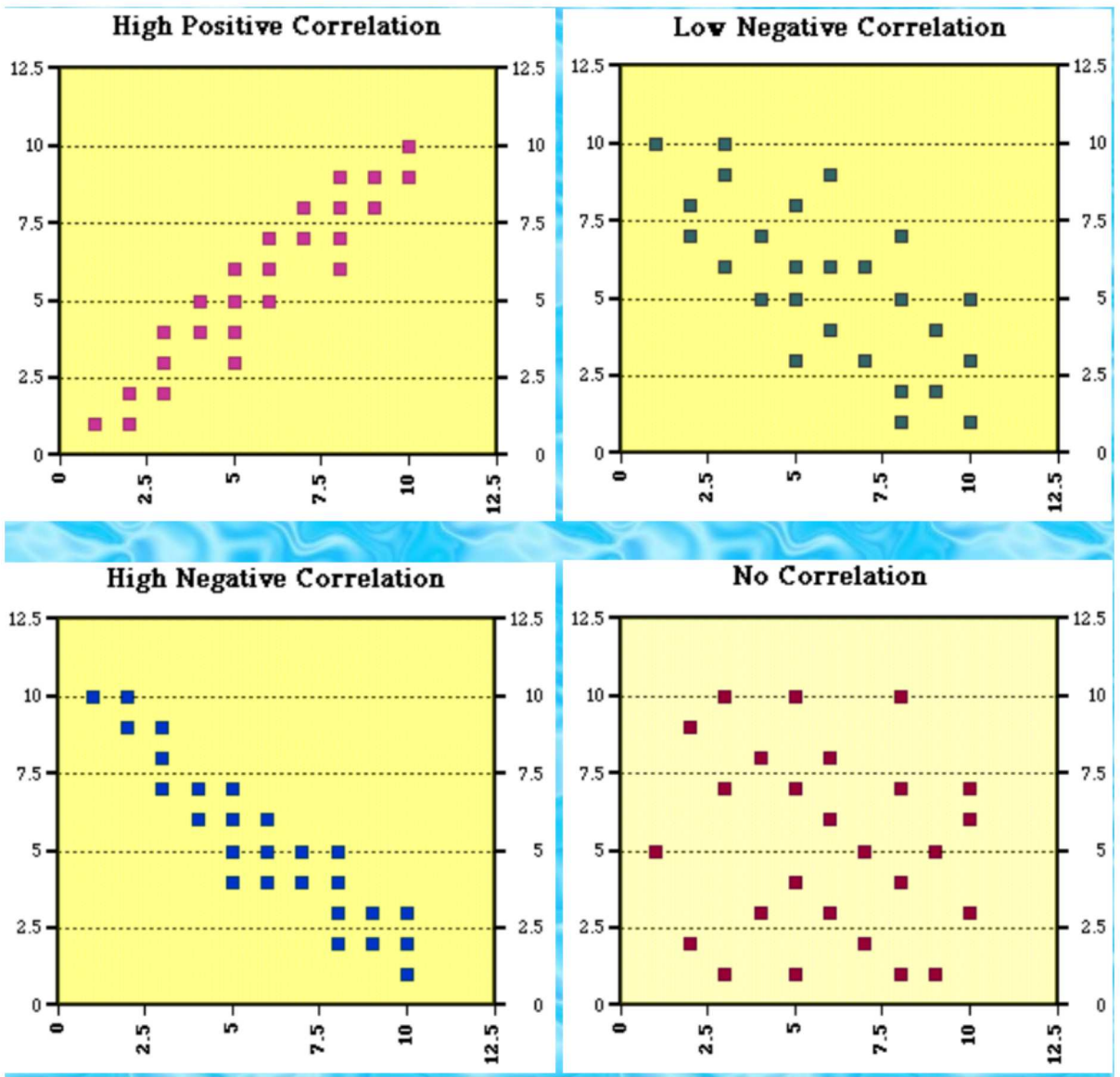
# Scatter plots and correlation

Allow us to compare two sets of data  $x = \{x_1, \dots, x_n\}$  and  $y = \{y_1, \dots, y_n\}$  and observe if any correlation (not causation) exists between them.



Correlation coefficient  $-1 \leq \rho \leq 1$





$$\text{Let } s_x = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2}$$

$$s_y = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2}$$

Then the Pearson's correlation coefficient may be written as:

$$r = \frac{1}{n-1} \sum_{i=1}^n \left( \frac{x_i - \bar{x}}{s_x} \right) \left( \frac{y_i - \bar{y}}{s_y} \right)$$

$$= \frac{\sum_{i=1}^n x_i y_i - n \bar{x} \bar{y}}{(n-1) s_x s_y}$$

If the correlation coefficient  $r$  is close to 1 or -1, then a straight line can be drawn through the data pts' region; which in "some measure" is the best fit line.

$$y_{LS} = \underbrace{mx}_{\text{slope}} + \underbrace{b}_{\text{y-intercept}}$$

For each data pair  $(x_i, y_i)$ , we define the residual as:  $r_i = mx_i + b_i - y_i$

Then let:

$$R(m, b) = \sum_{i=1}^n r_i^2 = \sum_{i=1}^n (mx_i + b_i - y_i)^2$$

We want to determine the constants  $m$  and  $b$  to make the total residual measure  $R(m, b)$  (sometimes called "cost function") as small as possible. This can be accomplished using calculus by setting:

$$\frac{\partial R}{\partial m} = 0 = \frac{\partial R}{\partial b}$$

The result is a system of equations from which  $m$  and  $b$  can be determined.



$$m = \frac{n \sum_{i=1}^n x_i y_i - \sum_{i=1}^n x_i \sum_{i=1}^n y_i}{n \sum_{i=1}^n x_i^2 - \left( \sum_{i=1}^n x_i \right)^2}$$

$$b = \frac{\sum_{i=1}^n x_i^2 \sum_{i=1}^n y_i - \sum_{i=1}^n x_i y_i \sum_{i=1}^n x_i}{n \sum_{i=1}^n x_i^2 - \left( \sum_{i=1}^n x_i \right)^2}$$

In R, the command `abline` can be used to plot a straight line with a specified slope and y-intercept.