# Statistics Data Science

Estimates of Location / measure of central tendency

Mean : The sum of all values divided by the number of values.

Median : The value such that one-half of the data lies above and below.
Synonym : 50th percentile

aka → middle data point in Sorted data

for odd count => $m = \dfrac{m_{-1} + m_{+1}}{2}$

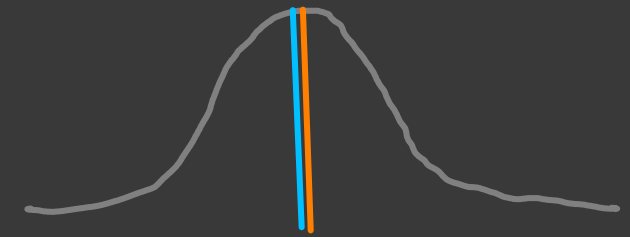Estimates of Variability / measure of dispersion or spread

Deviations : The difference between the observed values and the estimate of location.
Synonyms : errors, residuals

$= O - E$

$O = X \ , E = \bar{X}_{population} / \bar{X}_{mean}$

Variance : The sum of squared deviations from the mean divided by n – 1 where n is the number of data values.
Synonym : mean-squared-error

$S^2 = \dfrac{\sum\limits_{i=1}^{n}(O_i - E)^2}{n-1}$

in a normal distribution
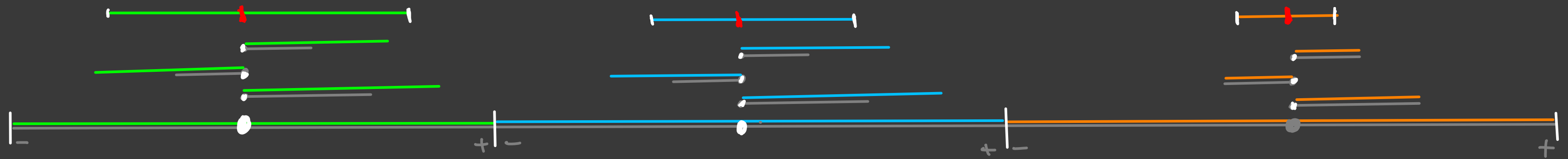


they overlap but in other distribution they differe

Standard deviation : The square root of the variance. $S = \sqrt{var}$

Mean absolute deviation : The mean of the absolute values of the deviations from the mean.
Synonyms : L1-norm, Manhattan norm

$$L1 = \sum_{i=1}^{n} \frac{abb(O-E)}{n}$$
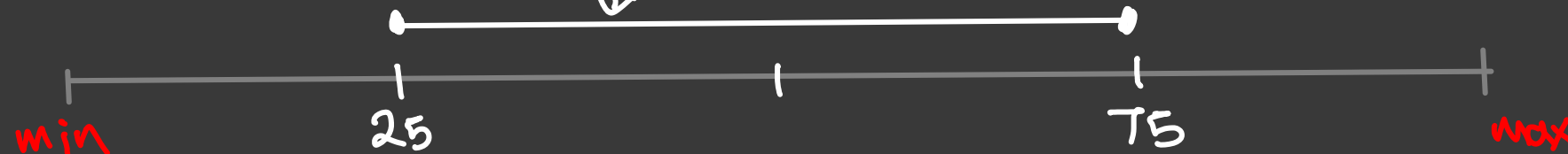
$$abb(-z/+z) \Rightarrow +z$$

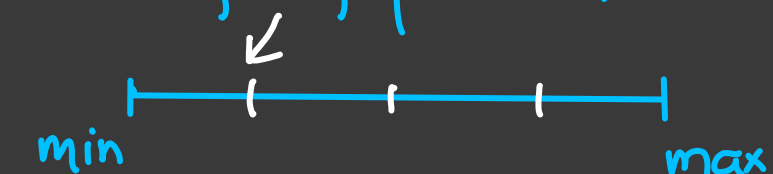Range : The difference between the largest and the smallest value in a data set.

Interquartile range : The difference between the 75th percentile and the 25th percentile.
Synonym : IQR

Range = max - min

min    25    75    max

* after sorting
data 25% lie
left of point X

min    max

*More robust metrics include mean absolute deviation, median absolute deviation from the median, and percentiles (quantiles).

Robust: Resistant toward outliers

mid: median of set/data of $O_{1...n}$ observation

Median absolute deviation from the median :
The median of the absolute values of the deviations from the median.

$$MAD = median(\,abs(O_1 - mid),\,abs(O_2 - mid),\,...\,)$$

Percentile : The value such that P percent of the values take on this value or less and (100-P) percent take on this value or more.

Synonym : quantile

let X=15

$100\% \rightarrow 1.0$
$30\% \rightarrow 0.3$
$28\% \rightarrow 0.28$

$$Percentile = \frac{count\ of\ O\ below\ X}{count\ of\ O\ total} = \frac{below\ count(X)}{50} = \frac{14}{50} = 0.28$$

50 percentile -> 50 % point in Observation are below some point X

$O_{1...50} = sorted(1,2,3,...,50)$

↓
0.5

data after sort = D

1. min of D
2. max of
3. First Quartile / Q1 / 25th percentile of
4. median / 50th percentile of  D
5. Third Quartile / Q3 / 75th percentile of

lower limit / fence ( left (<--) this point consider outlier )

higher limit / fence ( right (-->) this point consider outlier )

IRQ : is the difference in Q3 and Q1    $IRQ = Q3 - Q1$

$lower\ limit = Q1 - (1.5 \times IRQ)$

$higher\ limit = Q3 + (1.5 \times IRQ)$

outlier

outlier

min

max

$Q_2$

$Q_1$

$Q_3$

data good to use