

Standard Deviation and Variance

Measures of dispersion

Dataset 1

{1, 2, 3, 4, 5}

$$\text{Mean} = \frac{1+2+3+4+5}{5} = \frac{15}{5} = 3$$

Dataset 2

{0, 1, 2, 5, 7}

$$\text{Mean} = \frac{0+1+2+5+7}{5} = \frac{15}{5} = 3$$

Range = Max - Min

Dataset 1: Range = 5 - 1 = 4

Dataset 2: Range = 7 - 0 = 7

Variance

Definition:- Variance is basically the spread of the data from the mean

2 types:-

Population Variance

$$\sigma^2 = \sum_{i=1}^N \frac{(x_i - \mu)^2}{N}$$

Sample Variance

$$s^2 = \sum_{i=1}^N \frac{(x_i - \bar{x})^2}{n-1}$$

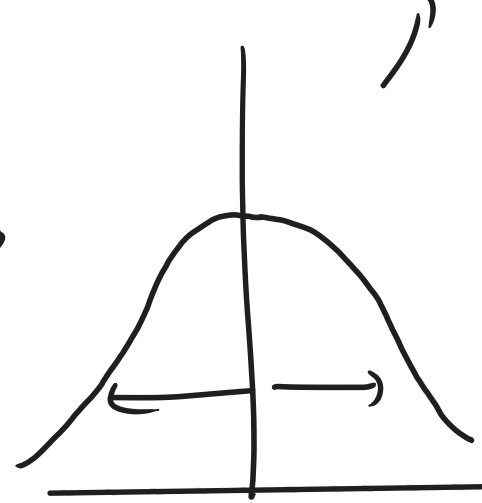
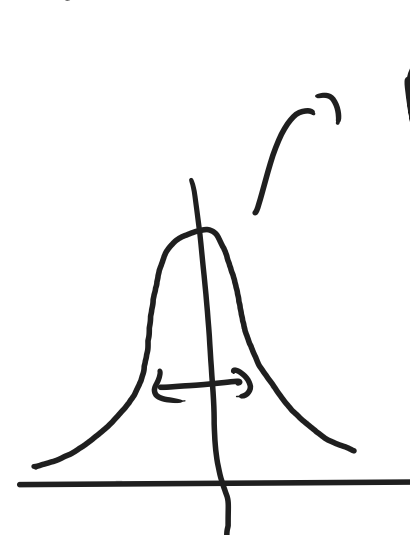
x	μ	$x - \mu$	$(x - \mu)^2$
1	3	(1-3)	4
2	3	(2-3)	1
3	3	(3-3)	0
4	3	(4-3)	1
5	3	(5-3)	4

Sum of $(x - \mu)^2 = 10$

Population Variance = $\frac{10}{5} = 2$

Dataset 1 Variance = 2

Dataset 2 Variance = 6.8



Variance is giving the overall spread

D1 variance = 2

D2 variance = 6.8

{1, 2, 3, 4, 5}

{0, 1, 2, 5, 7}

Standard Deviation

just the square root of variance

gives us the spread of a certain data with respect to the mean

D1 variance = 2 \rightarrow std = $\sqrt{2} = 1.414$

D2 variance = 6.8 \rightarrow std = $\sqrt{6.8} = 2.608$

D1 = {1, 2, 3, 4, 5}

D2 = {0, 1, 2, 5, 7}

Variance = Overall spread

Standard deviation = average spread or range of values around the mean

Dataset 1

mean

3

std

1.414

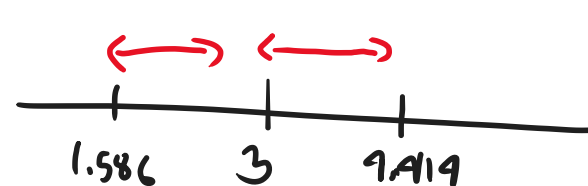
{0, 1, 2, 3, 4, 5}

3 + 1.414

= 4.414

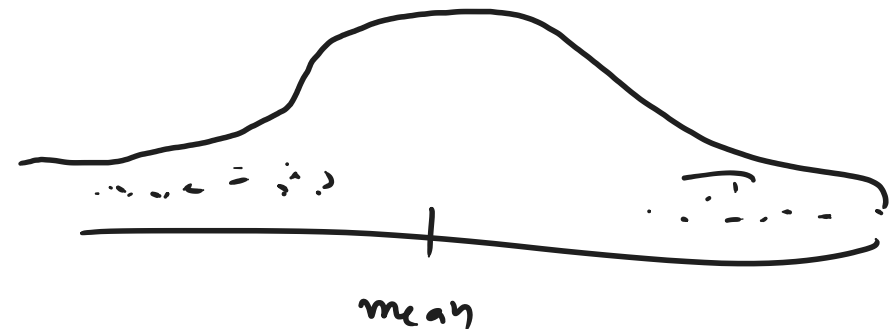
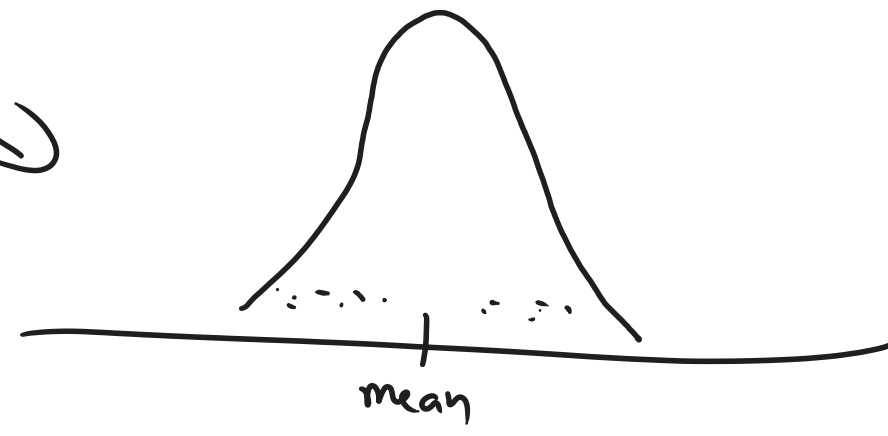
3 - 1.414

= 1.586



std = 10

std = 50



Why did we use n-1 in sample variance?

instead of N

based on degrees of freedom

we use sample mean to find the sample variance

there is uncertainty to prevent this we used (n-1)

gives us unbiased estimate