

# Topic: Exploratory Data Analysis (EDA)

## Measures of Variability - Part A

School of Mathematics and Applied Statistics



UNIVERSITY  
OF WOLLONGONG  
AUSTRALIA

# Measuring Variability: Motivating Example

Consider the following data sets:  $n=7$

Data set 1: 55, 55, 55, 55, 55, 55, 55

Data set 2: 47, 51, 54, 55, 56, 59, 63

Data set 3: 39, 47, 53, 55, 57, 63, 71

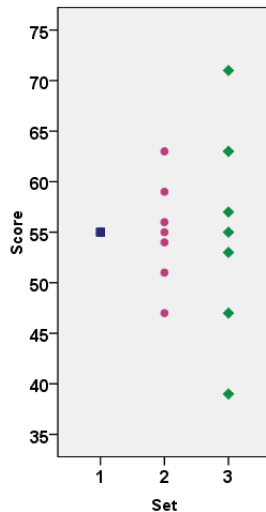
Total
385
385
385

For each data set

- Median = 55 =  $\theta_2$
- Mean = 385/7 = 55 =  $\bar{x}$

But the spread of the scores

vary.



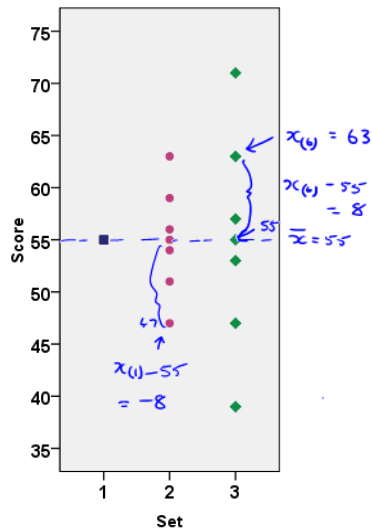
# How do we Measure Variability?

**Variability (spread)** can be measured by:

- **Variance**  $\sigma^2$  or  $s^2$ 
  - uses all data values but is inflated by outliers
- **Standard deviation**  $\sigma$  or  $s$ 
  - uses all data values but is inflated by outliers
- **Range** = maximum <sup>✓</sup> – minimum <sup>✓</sup> =  $x_{(n)} - x_{(1)}$ 
  - unreliable measure, depends on extreme values
- **Interquartile range**:  $IQR = Q_3 - Q_1$ 
  - spans middle 50% of data, <sup>upper</sup> <sup>lower quartile</sup>
  - unaffected by outliers, ignores variation in tails

# Variance and Standard Deviation

- The **mean** is used as a reference point.  $\bar{x} = 55$
- Consider the **deviation** from the  $i$ th point to the mean:  $x_i - \bar{x}$
- Deviations may be positive or negative
- **Variance** is based on the squared deviations.
- We can also describe the difference in spread using a notion of average distance from the mean.
- This measure of variability is called the **standard deviation**.
- $\text{sum of deviations} = 0$ .



# Variance

**Variance** is based on the squared distances of individual data points from the mean.

## Population Variance

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

*Handwritten notes:* "sigma" with an arrow pointing to  $\sigma^2$ ; "pop mean" with an arrow pointing to  $\mu$ ; "mv." with an arrow pointing to  $\mu$ ; the  $N$  in the denominator is circled.

## Sample Variance

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

*Handwritten notes:* "sample mean" with an arrow pointing to  $\bar{x}$ ; the  $n-1$  in the denominator is underlined.

- measurement is in squared units
- is never negative, and
- is only zero when all data values are identical
- $s^2$  is an unbiased estimator of  $\sigma^2$ . ✓✓

Σ

# Standard Deviation

## Standard deviation $\sigma$ or $s$

- is the square root of the variance
- is measured in same units of measurement as data

### Population standard deviation

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

$$\sigma = \sqrt{\sigma^2}$$

### Sample standard deviation

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

$$s = \sqrt{s^2}$$

- On calculator, enter data then use: *STAT mode*

*Pop sd.*  $\left\{ \begin{array}{l} \sigma_n \text{ or } x\sigma_n \\ \sigma_x \end{array} \right\}$  or  $\sigma_{n-1} \text{ or } x\sigma_{n-1} \text{ or } s_{n-1} \text{ or } s \text{ key}$   $\left. \vphantom{\left\{ \right\}} \right\}$  *Sample sd.*  
 $\sigma_x$   $s_x$

# Example

- Data set 1: 55, 55, 55, 55, 55, 55, 55
- Data set 2: 47, 51, 54, 55, 56, 59, 63
- Data set 3: 39, 47, 53, 55, 57, 63, 71

	Median	Mean	SD
Set 1	55	55	
Set 2	55	55	
Set 3	55	55	

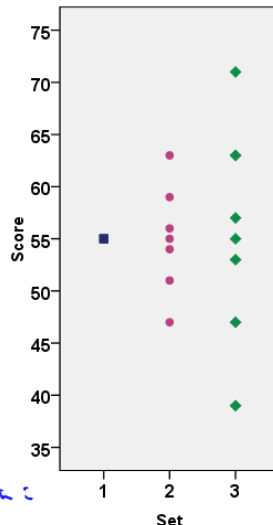
**Discuss:** What is the SD for Set 1? Why?

$$s_1 = 0 \quad (S)$$

$$x_i - \bar{x} = 0$$

$$x_i = 55 \text{ for each } i$$

$$\bar{x} = 55$$



# Activity: Calculate Variance and SD

**Exercise:** What is the sample variance and standard deviation for Set 2 and Set 3?  
Calculate  $s^2$  and  $s$ :

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

## Data Set 2 $\bar{x} = 55$

47, 51, 54, 55, 56, 59, 63

$x_i$	$x_i - \bar{x}$	$(x_i - \bar{x})^2$
47	-8	64
51	-4	16
54	-1	1
55	0	0
56	1	1
59	4	16
63	8	64
	<u>0</u>	<u>162</u>

$= \sum_{i=1}^7 (x_i - \bar{x})^2$

## Data Set 3

39, 47, 53, 55, 57, 63, 71

$$s^2 = \frac{1}{7-1} \times 162$$

$$= 27$$

$$s = \sqrt{27}$$

$$= 5.196 \dots$$



# Activity: Calculate Variance and SD cont.

Set 3.

$x_i$	$x_i - \bar{x}$	$(x_i - \bar{x})^2$
39	-16	256
47	-8	64
53	-2	4
55	0	0
57	2	4
63	8	64
71	16	256
		<u>648</u>

$$= \sum_{i=1}^7 (x_i - \bar{x})^2$$

$$s_3^2 = \frac{1}{7-1} \times 648$$

$$= \underline{108}$$

$$s_3 = \sqrt{108}$$

$$= 10.392.$$

	Set 1	Set 2	Set 3
s	0	5.196	<u>10.392</u>

# In R: Calculate mean, variance and sd

```
> Set2 <- c(47, 51, 54, 55, 56, 59, 63)
```

```
> mean(Set2)
```

```
[1] 55
```

```
> var(Set2)
```

```
[1] 27
```

```
> sd(Set2)
```

```
[1] 5.196152
```

```
> Set3 <- c(39, 47, 53, 55, 57, 63, 71)
```

```
> mean(Set3)
```

```
[1] 55
```

```
> var(Set3)
```

```
[1] 108
```

```
> sd(Set3)
```

```
[1] 10.3923
```