



INTRODUCTION TO PROBABILITY AND STATISTICS FOURTEENTH EDITION

Chapter 2

Describing Data

with Numerical Measures

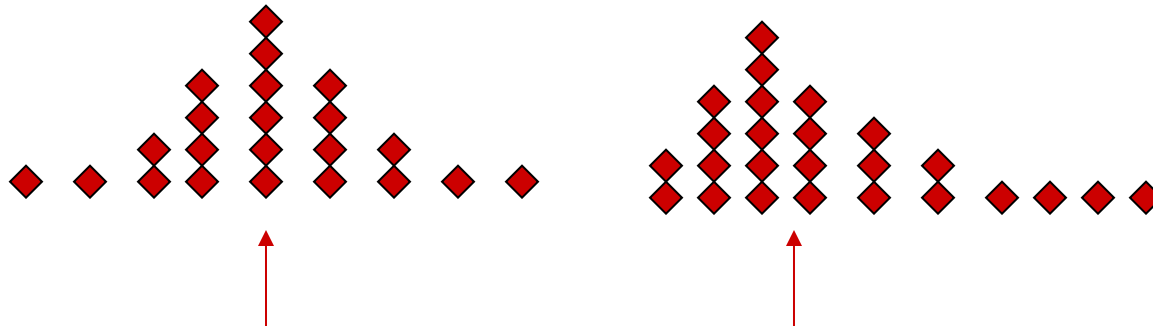
DESCRIBING DATA WITH NUMERICAL MEASURES

- Graphical methods may not always be sufficient for describing data.
- **Numerical measures** can be created for both **populations** and **samples**.
 - A **parameter** is a numerical descriptive measure calculated for a population.
 - A **statistic** is a numerical descriptive measure calculated for a sample.



MEASURES OF CENTER

- A measure along the horizontal axis of the data distribution that locates the center of the distribution.



ARITHMETIC MEAN OR AVERAGE

- The **mean** of a set of measurements is the sum of the measurements divided by the total number of measurements.

$$\bar{x} = \frac{\sum x_i}{n}$$

where n = number of
measurements

$\sum x_i$ = sum of all the measurements



EXAMPLE

- The set: 2, 9, 1, 5, 6

$$\bar{x} = \frac{\sum x_i}{n} = \frac{2 + 9 + 1 + 5 + 6}{5} = \frac{23}{5} = 4.6$$

If we were able to enumerate the whole population, the **population mean** would be called μ (the Greek letter “mu”).



MEDIAN

- The **median** of a set of measurements is the middle measurement when the measurements are ranked from smallest to largest.
- The **position of the median** is

$$.5(n + 1)$$

once the measurements have been ordered.



EXAMPLE

- The set: 2, 4, 9, 8, 6, 5, 3 $n = 7$
- Sort: 2, 3, 4, 5, 6, 8, 9
- Position: $.5(n + 1) = .5(7 + 1) = 4^{\text{th}}$

Median = 4th largest measurement

- The set: 2, 4, 9, 8, 6, 5 $n = 6$
- Sort: 2, 4, 5, 6, 8, 9
- Position: $.5(n + 1) = .5(6 + 1) = 3.5^{\text{th}}$

Median = $(5 + 6)/2 = 5.5$ — average of the 3rd and 4th measurements



MODE

- The **mode** is the measurement which occurs most frequently.
- The set: 2, 4, 9, 8, 8, 5, 3
 - The mode is **8**, which occurs twice
- The set: 2, 2, 9, 8, 8, 5, 3
 - There are two modes—**8** and **2** (bimodal)
- The set: 2, 4, 9, 8, 5, 3
 - There is **no mode** (each value is unique).



EXAMPLE

The number of quarts of milk purchased by 25 households:



0 0 1 1 1 1 1 2 2 2 2 2 2
2 2 2 3 3 3 3 3 4 4 4 5

○ Mean?

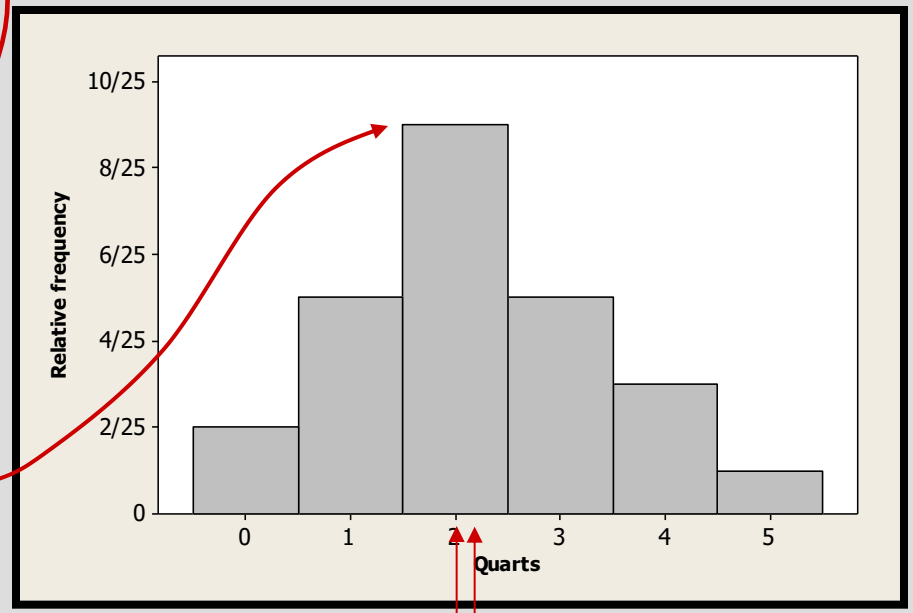
$$\bar{x} = \frac{\sum x_i}{n} = \frac{55}{25} = 2.2$$

○ Median?

$$m = 2$$

○ Mode? (Highest peak)

$$\text{mode} = 2$$



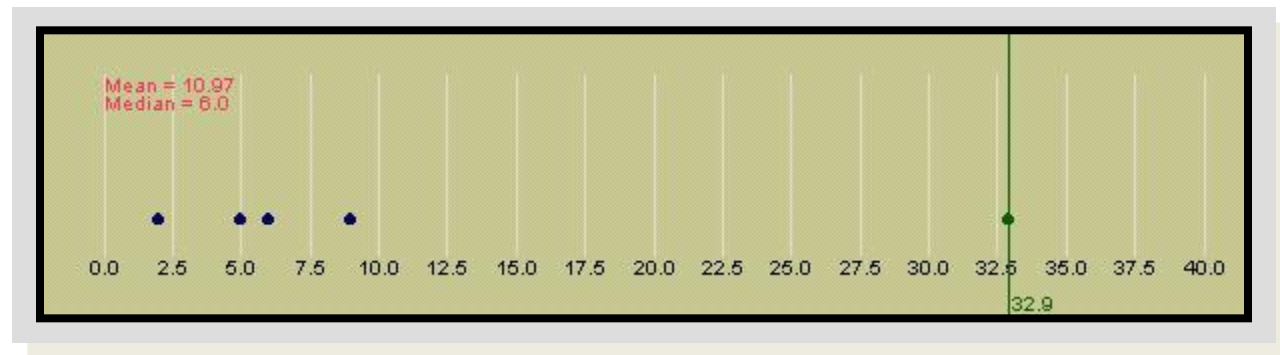
FREQUENCY TABLE FOR QUARTS OF MILK PURCHASED

Quarts of Milk	Freq.
0	2
1	5
2	9
3	5
4	3
5	1
Total	25



EXTREME VALUES

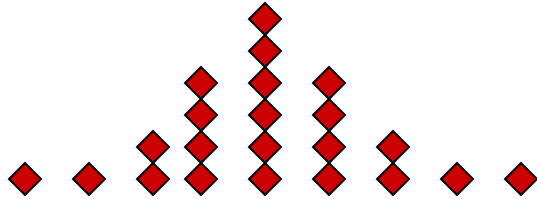
- The mean is more easily affected by extremely large or small values than the median.



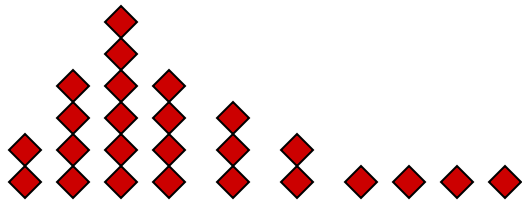
- The median is often used as a measure of center when the distribution is skewed.



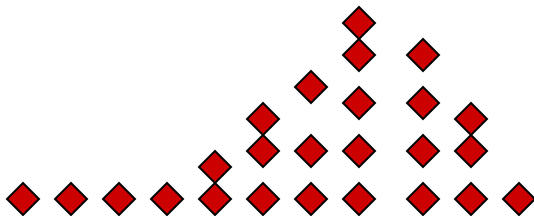
EXTREME VALUES



Symmetric: Mean = Median



Skewed right: Mean > Median

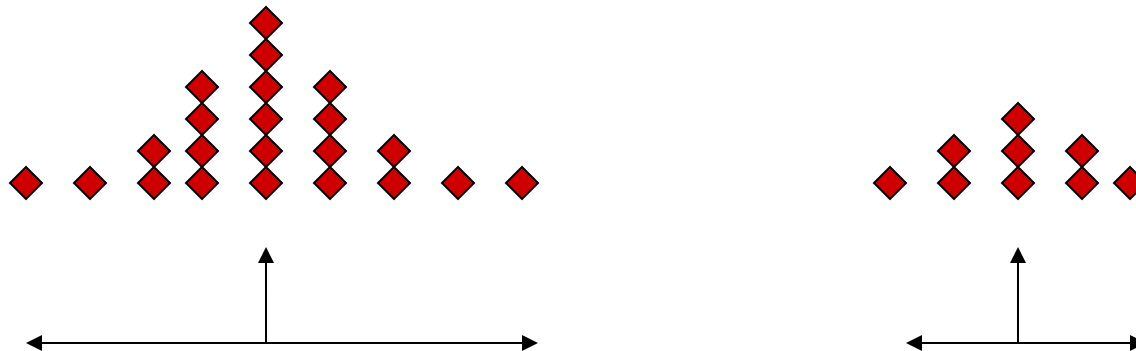


Skewed left: Mean < Median

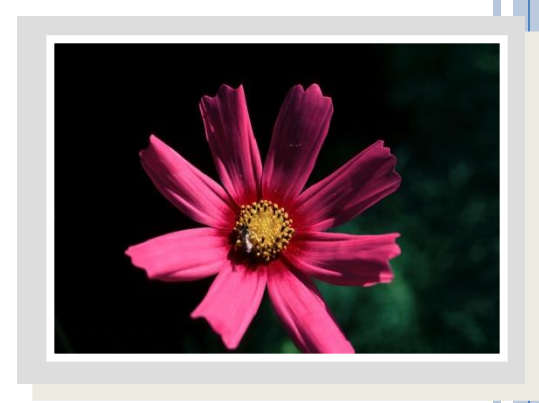


MEASURES OF VARIABILITY

- A measure along the horizontal axis of the data distribution that describes the **spread** of the distribution from the center.



THE RANGE



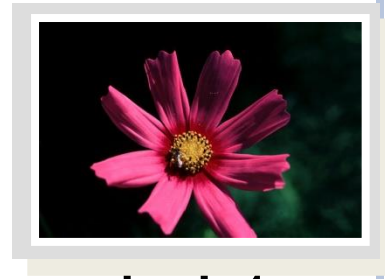
- The **range, R** , of a set of n measurements is the difference between the largest and smallest measurements.
- **Example:** A botanist records the number of petals on 5 flowers:

5, 12, 6, 8, 14

- The range is **$R = 14 - 5 = 9$.**

• Quick and easy, but only uses 2 of the 5 measurements.

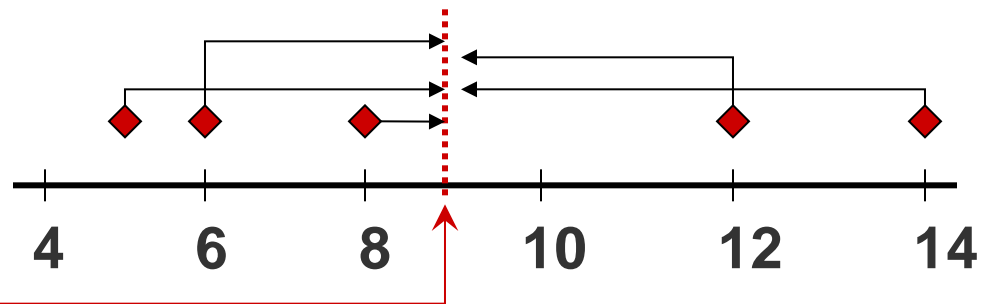




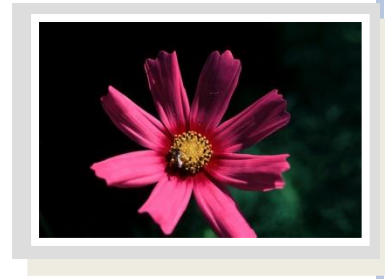
THE VARIANCE

- The **variance** is a measure of variability that uses all the measurements. It measures the average deviation of the measurements about their mean.
- **Flower petals: 5, 12, 6, 8, 14**

$$\bar{x} = \frac{45}{5} = 9$$



THE VARIANCE



- The **variance of a population** of N measurements is the average of the squared deviations of the measurements about their mean μ .

$$\sigma^2 = \frac{\sum (x_i - \mu)^2}{N}$$

- The **variance of a sample** of n measurements is the sum of the squared deviations of the measurements about their mean, divided by $(n - 1)$.

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



THE STANDARD DEVIATION



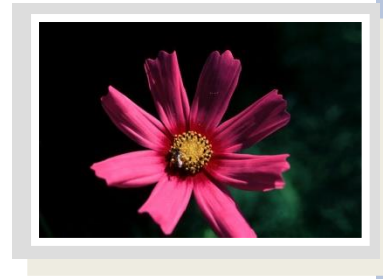
- In calculating the variance, we squared all of the deviations, and in doing so changed the scale of the measurements.
- To return this measure of variability to the original units of measure, we calculate the **standard deviation**, the positive square root of the variance.

Population standard deviation : $\sigma = \sqrt{\sigma^2}$

Sample standard deviation : $s = \sqrt{s^2}$



TWO WAYS TO CALCULATE THE SAMPLE VARIANCE



Use the Definition Formula:

	x_i	$x_i - \bar{x}$	$(x_i - \bar{x})^2$
	5	-4	16
	12	3	9
	6	-3	9
	8	-1	1
	14	5	25
Sum	45	0	60

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

$$= \frac{60}{4} = 15$$

$$s = \sqrt{s^2} = \sqrt{15} = 3.87$$

TWO WAYS TO CALCULATE THE SAMPLE VARIANCE



Use the Computational Formula:

	x_i	x_i^2
	5	25
	12	144
	6	36
	8	64
	14	196
Sum	45	465

$$s^2 = \frac{\sum x_i^2 - \frac{(\sum x_i)^2}{n}}{n-1}$$

$$= \frac{465 - \frac{45^2}{5}}{4} = 15$$

$$s = \sqrt{s^2} = \sqrt{15} = 3.87$$

SOME NOTES

- The value of s is **ALWAYS** positive.
- The larger the value of s^2 or s , the larger the variability of the data set.
- **Why divide by $n - 1$?**
 - The sample standard deviation s is often used to estimate the population standard deviation σ . Dividing by $n - 1$ gives us a better estimate of σ .

