Lec3

### 1.3.3 Population Mean

The average of all values in the population is called **pop mean** and is denoted by **$\mu$**. When there are $N$ values in the population, then $\mu = \dfrac{\sum \text{values in whole pop}}{N}$

One of our first tasks in statistical inference will be to present methods based on the **sample mean** for drawing conclusions about a **pop. mean**

For example, Ex. 1.14     $\bar{x} = 16.36$

$\uparrow$ guess for $\mu$

in whole pop.

### 1.3.4 Population Median

Analogous to $\widetilde{x}$ as the middle value in the sample, the **population median** is **middle value** denoted by **$\widetilde{\mu}$**. As with $\bar{x}$ and $\mu$, we can think of using **sample median** to make an inference about **pop median**.

**Summary:**

sample mean + median     V.S. pop. mean + median

### 1.3.5 Other Measures of Location

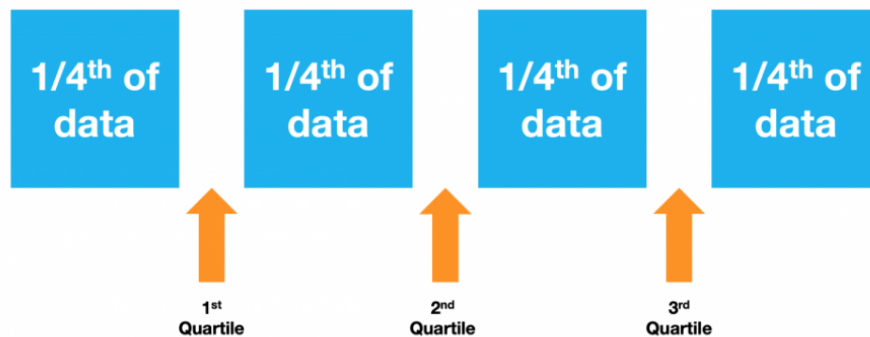**Maximum**

*Maximal value*

**Minimum**

*Minimal value*

An **outlier** is _an atypical data pt (too large or too small)_

Sometimes _Max and/or min_ are outliers in the data set.

**Quartiles and Percentiles**

Quartiles _divide data set into 4 equal parts_
_(25% )Q1, 25% )Q2=Median, 25% )Q3, 25%)_

with the observations above the third quartile constituting the upper quarter of the data set,

the second quartile being identical to the _Median_, and the first quartile separating

the lower quarter from the upper three-quarters.



|  |  |  |  |
|---|---|---|---|
| 1/4th of data | 1/4th of data | 1/4th of data | 1/4th of data |

1st Quartile    2nd Quartile    3rd Quartile

If the quantiles divide the data into _100 groups_, then they're called _percentiles_

10

### 1.3.6 The effect of skewness on the mean and median

$$\tilde{x} \approx \bar{x} \rightarrow \text{Symmetric}$$

$$\tilde{x} < \bar{x} \rightarrow \text{right skewed}$$

$$\tilde{x} > \bar{x} \rightarrow \text{left skewed}$$

Think about the graph in extreme cases.

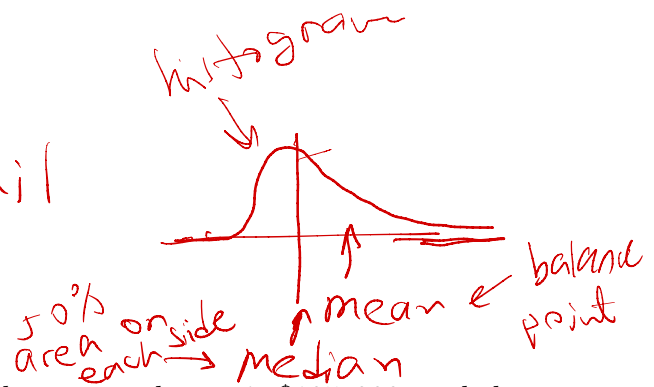Suppose we have 11 data in a set of students grades, in which ten of them are 60s and one of them is 100.

$$60 \quad 60 \quad 60 \quad 60 \cdots 60 \qquad 100$$

$$\tilde{x} = 60 \qquad \bar{x} > 60$$

On the other hand, suppose we have 11 data in a set of students grades, in which ten of them are 100s and one of them is 60.

$$60 \quad 100 \ 100 \ 100 --- 100$$

$$\tilde{x} = 100 \qquad \bar{x} < 100$$

11

**Note.** *mean follows the tail*

*histogram*

*50% area on each side → median*

*mean ← balance point*

**Example 4.** Suppose we have 10 people in a room, the mean salary $\bar{x}$ is $105,000$, and the median salary $\tilde{x}$ is $65,000$. What can we say about the distribution?

*Ans :* right skewed

There's likely an outlier

**What is the more appropriate measure of center when there are extreme valus in the data set and why?**
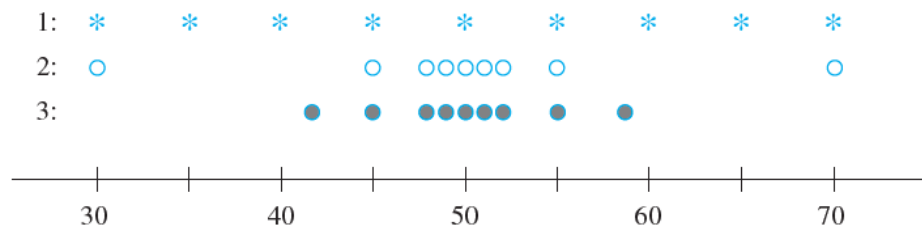
*Ans :* Median

## 1.4 Measures of Variability

Figure below shows dotplots of three samples with the same mean and median, but the extent of spread about the center is different for all three samples. The first sample has the largest amount of variability the third has the smallest amount, and the second is intermediate.

```
1:   *      *      *      *      *      *      *      *      *
2:   o                    o   ooooo   o                     o
3:                  ●      ●   ●●●●●   ●      ●
     ┼──────┼──────┼──────┼──────┼──────┼──────┼──────┼
        30            40            50            60            70
```

Samples with identical measures of center but different amounts of variability

## Measures of Variability for Sample data

The simplest measure of variability in a sample is the _range_, which is _span of data set_

- $max - min$

The value of the range for sample 1 in Figure above is _much larger_ than sample 3, reflecting _more variability_ in the first sample than in the third. A defect of the range, though, is _sensitive to outliers._

Samples 1 and 2 in Figure above have _the same range_ yet when the observations between the two extremes are taken into account, there is _much less variability_ in the second sample than in the first.

Our primary measures of variability involve _values − mean_

$$x_1 - \bar{x} \qquad x_2 - \bar{x} \qquad x_3 - \bar{x} \qquad ---$$

A deviation will be **positive** (+) if $\underline{x_i > \bar{x}}$ and **negative** (−) if $\underline{x_i < \bar{x}}$

If all the deviations are _small_, then all $x_i$'s are _close to $\bar{x}$_ and there is _little variability_. Alternatively, if some of the deviations are _large_, then some $x_i$'s _far from $\bar{x}$_, suggesting a _larger variability_

13

A simple way to combine the deviations into a **single quantity** is to **average**. Out

Unfortunately, this is a bad idea:

$$\text{sum of deviations} = (x_1 - \bar{x}) + (x_2 - \bar{x}) + \cdots + (x_n - \bar{x})$$

$$= x_1 + \cdots + x_n - n\bar{x} = n\bar{x} - n\bar{x}$$

so that the average deviation is _____.

$$= 0$$

To prevent **negative and positive values** from counteracting one

another when they are combined, we consider instead **squared deviations**

$$(x_1 - \bar{x})^2 \quad (x_2 - \bar{x})^2 \quad -- \quad \leftarrow \text{think}$$

Euclidean distance

To get the **averaged square deviation**, for several reasons we will not cover here, we

divide the sum of squared deviations by $n - 1$ instead of $n$.

---

**Definition 1.** The sample **variance**, denoted by $s^2$, is given by

$$s^2 = \frac{(x_1 - \bar{x})^2 + \cdots + (x_n - \bar{x})^2}{n - 1}$$

The sample **standard deviation(SD)**, denoted by $s$, is the $\sqrt{\phantom{xxx}}$ of the variance:

$$s = \sqrt{s^2} = \sqrt{\frac{(x_1 - \bar{x})^2 + \cdots + (x_n - \bar{x})^2}{n - 1}}$$

---

The **standard deviation** is our preferred measure of variability, because

**it has the same unit as data.**

---

A shortcut formula for $s$,

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

---

This is because we can write $s^2$ as

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

$$\downarrow 2n\bar{x}^2$$

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

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

$$(\sum x_i)^2 \leftarrow n\left(\frac{\sum x_i}{n}\right)^2$$

14

**Example 5.** Suppose we have a data set $x_1, x_2, \cdots, x_5$ where $\sum_{i=1}^{5} x_i = 10.9$ and $\sum_{i=1}^{5} x_i^2 = 29.97$. What is the SD?

*Ans* :

$$\sqrt{\dfrac{29.97 - \dfrac{(10.9)^2}{5}}{4}} = 1.246$$

$$\bar{x} = \dfrac{10.9}{5}$$

.check $\nearrow$ equals the s you get from the original formula!

**Properties of the mean and SD:**

Let $x_1, x_2, \cdots, x_n$ be a sample and $c$ be any nonzero constant.

1. If $y_1 = x_1 + c,\ y_2 = x_2 + c, \cdots, y_n = x_n + c,$

   then mean $\underline{\bar{y} = \bar{x} + c}$ $\leftarrow$ $\dfrac{(x_1+c) + \cdots + (x_n+c)}{n} = \dfrac{x_1 + \cdots + x_n + nc}{n}$

   the sample variance $\underline{S_y^2 = S_x^2}$ $= \bar{x} + c$

   the SD $\underline{S_y = S_x}$

   $\left( (x_i + c) - (\bar{x} + c) \right)^2$

2. If $y_1 = cx_1,\ y_2 = cx_2, \cdots, y_n = cx_n,$

   then mean $\underline{\bar{y} = c\bar{x}}$

   the sample variance $\underline{S_y^2 = c^2 S_x^2}$

   the SD $\underline{S_y = |c| S_x}$

   $\dfrac{cx_1 + \cdots + cx_n}{n} = c\bar{x}$

$\uparrow$

$\sqrt{c^2} = |c|$

$\left( \sqrt{(-5)^2} = 5 \text{ not } -5 \right)$

$(c x_i - c\bar{x})^2$

$= c^2 (x_i - \bar{x})^2$

$=$

15