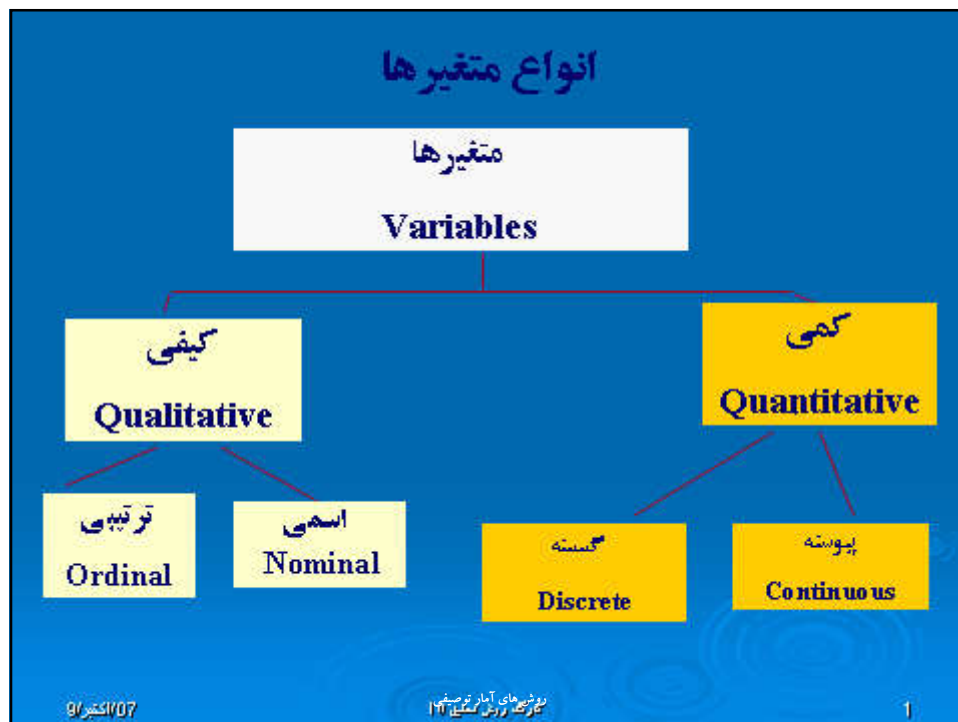


# مروری بر روش های آمار توصیفی

دکتر یددا... محرابی

دانشکده بهداشت - دانشگاه علوم پزشکی شهید بهشتی

[ymehrabi@gmail.com](mailto:ymehrabi@gmail.com)



## Numerical Summaries

- **Center of the data**
  - mean
  - Median
- **Variation**
  - range
  - quartiles (interquartile range)
  - variance
  - standard deviation

## Mean or Average

- Traditional measure of center
- Sum the values and divide by the number of values

$$\bar{x} = \frac{1}{n}(x_1 + x_2 + \cdots + x_n) = \frac{1}{n} \sum_{i=1}^n x_i$$

## Median ( $M$ )

- A resistant measure of the data's center
- At least half of the **ordered** values are less than or equal to the median value
- At least half of the **ordered** values are greater than or equal to the median value
- If  $n$  is odd, the median is the middle ordered value
- If  $n$  is even, the median is the average of the two middle ordered values

## Median ( $M$ )

Location of the median:  $L(M) = (n+1)/2$  ,  
where  $n$  = sample size.

**Example:** If 25 data values are recorded, the Median would be the  $(25+1)/2 = 13^{\text{th}}$  ordered value.

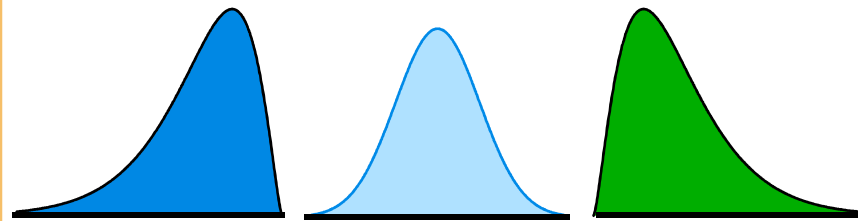
## Median

- Example 1 data: 2 4 6  
Median ( $M$ ) = 4
- Example 2 data: 2 4 6 8  
Median = 5 (ave. of 4 and 6)
- Example 3 data: 6 2 4  
Median  $\neq$  2  
(**order** the values: 2 4 6 , so Median = 4)

## Comparing the Mean & Median

- The mean and median of data from a **symmetric** distribution should be close together. The actual (true) mean and median of a symmetric distribution are exactly the same.
- In a **skewed** distribution, the mean is farther out in the long tail than is the median [the mean is 'pulled' in the direction of the possible outlier(s)].

# Skewness

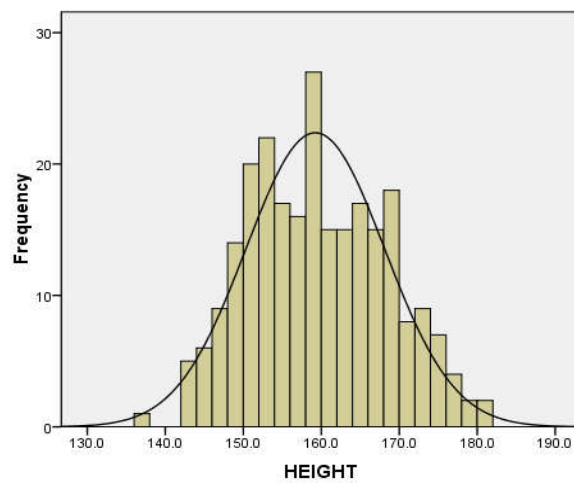


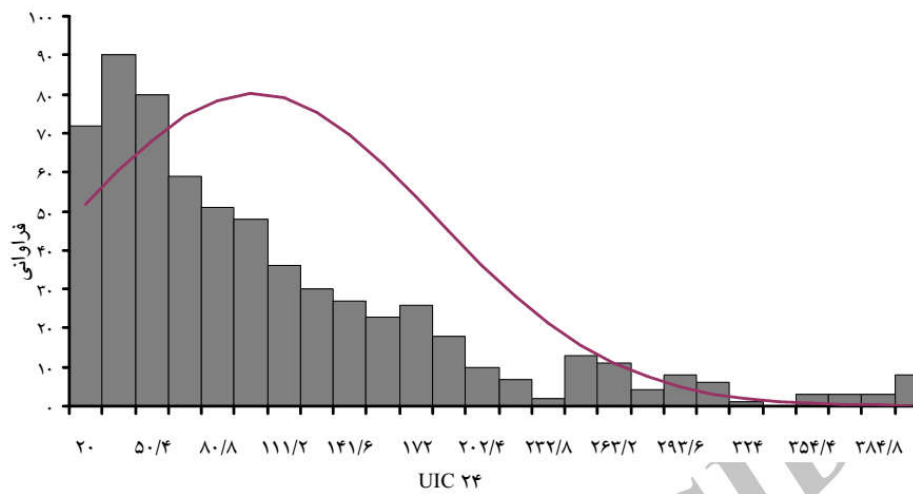
**Negatively  
Skewed**

**Symmetric  
(Not Skewed)**

**Positively  
Skewed**

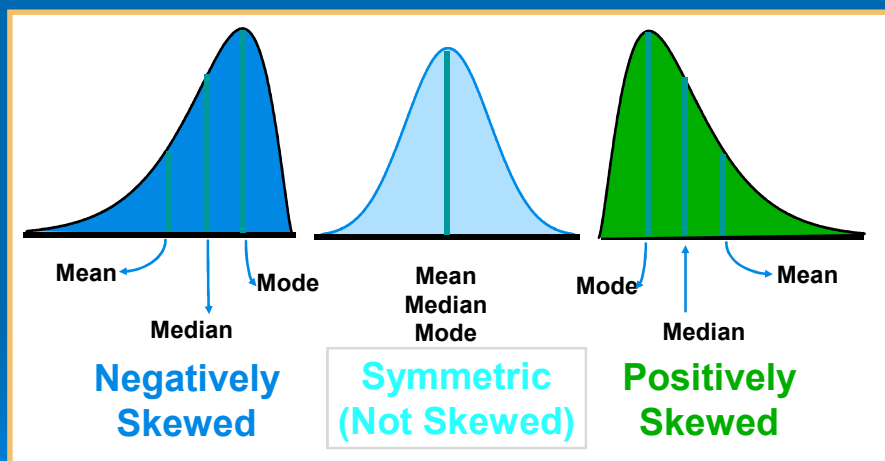
HEIGHT

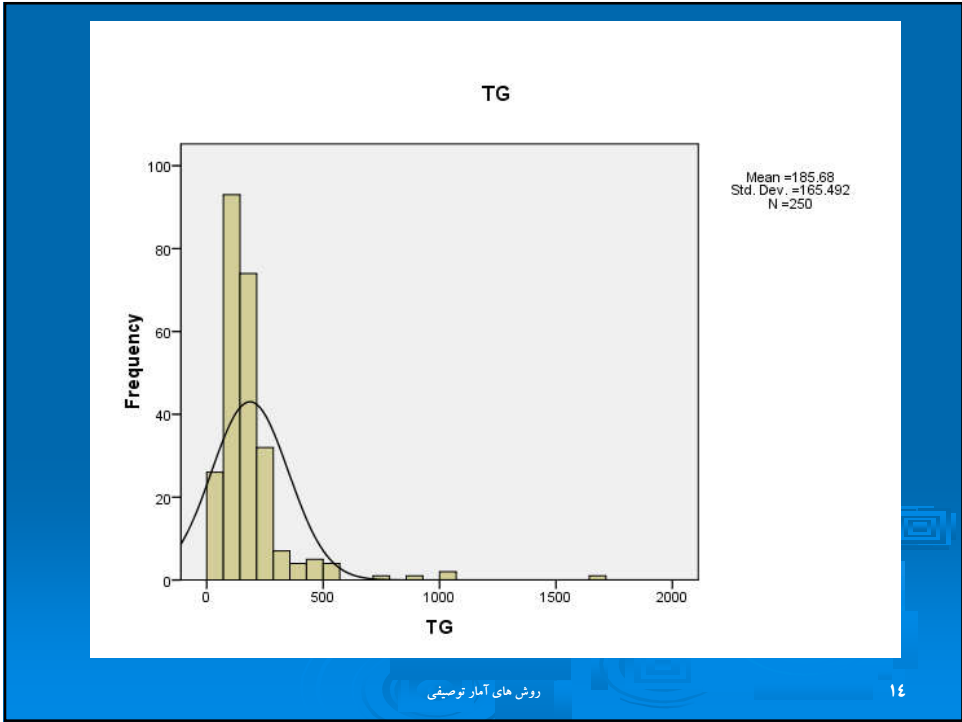
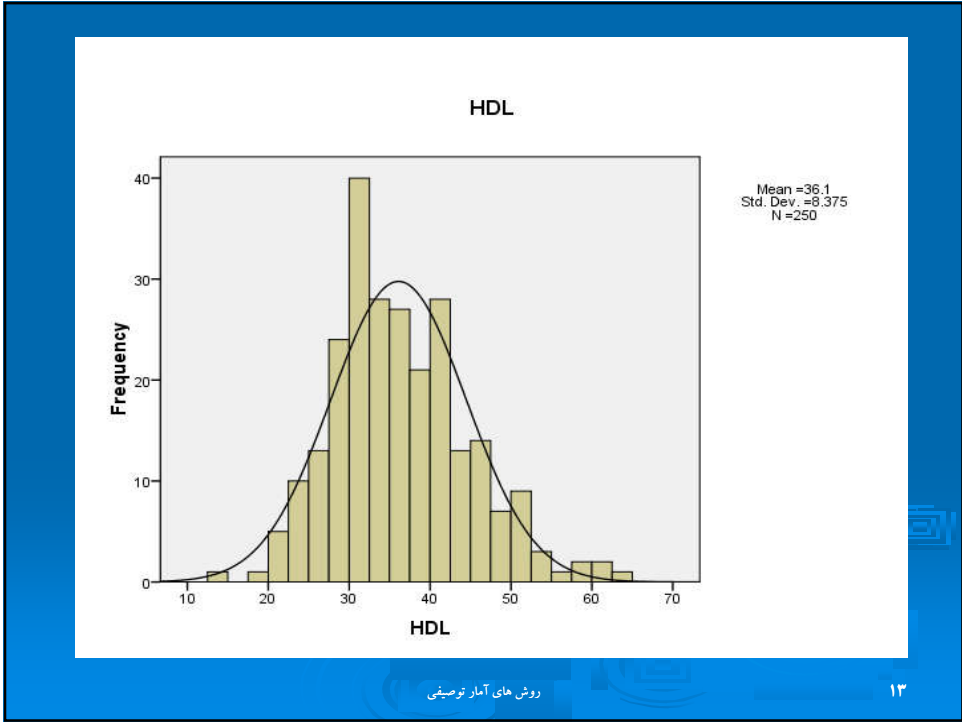


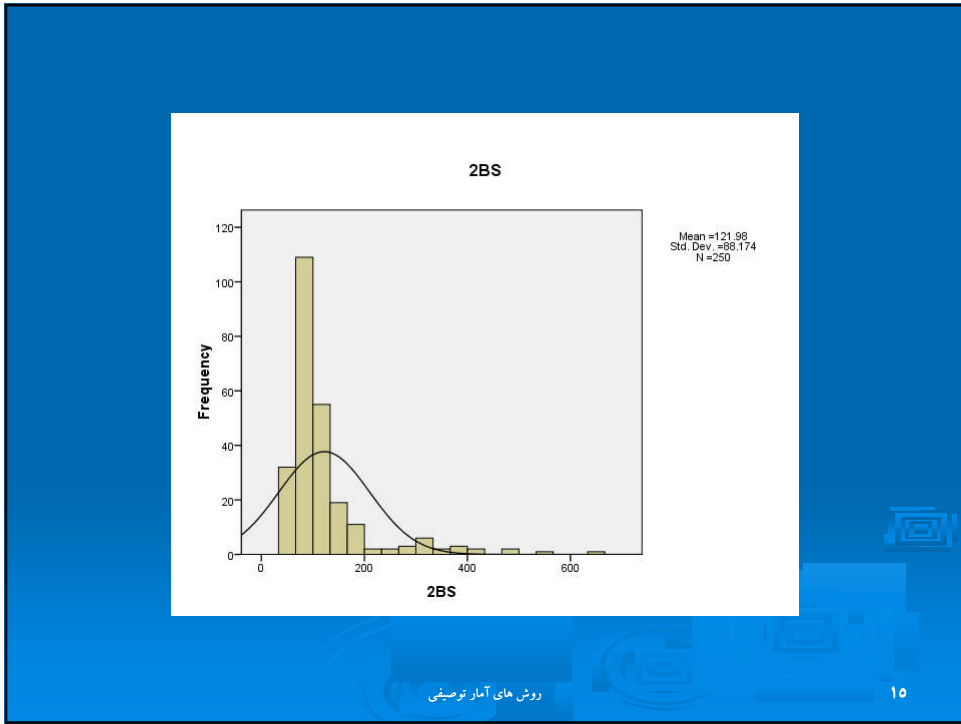


نمودار ۱- هیستوگرام غلظت ید ادران ۲۴ ساعته

## Skewness



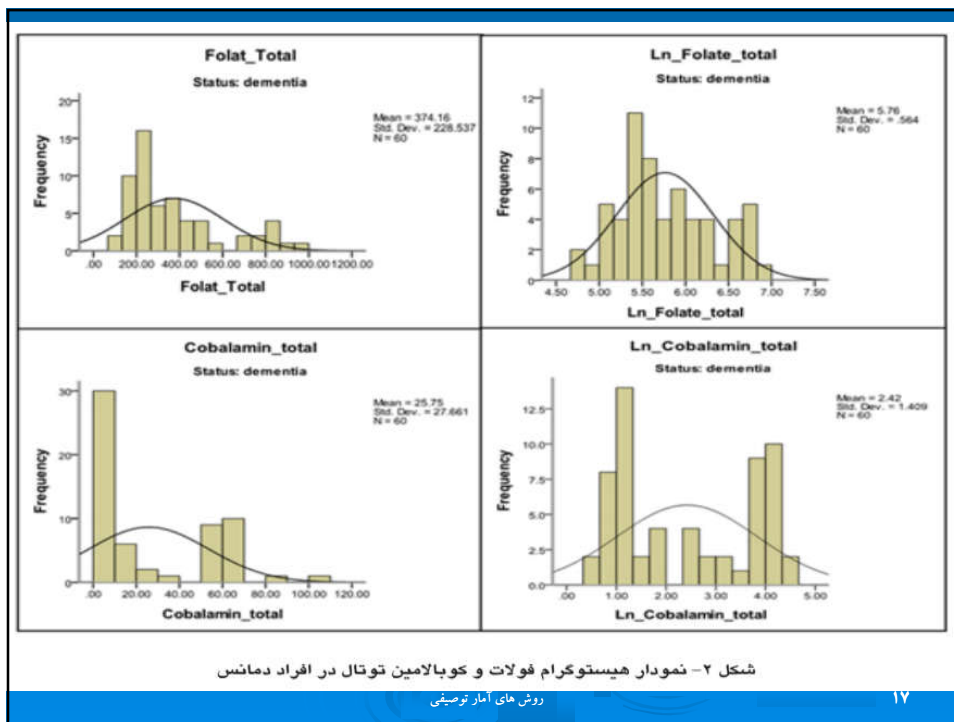




**مثال**

	WEIGHT	HEIGHT	HDL	TG	2hrBS
N Valid	250	249	250	250	250
<b>Mean</b>	<b>67.5</b>	<b>159.2</b>	36.1	<b>185.6</b>	<b>121.9</b>
<b>Median</b>	<b>67.2</b>	<b>158.5</b>	35.0	<b>148.5</b>	<b>97.0</b>
Std. Deviation	12.40	8.8	8.3	165.5	88.2
Range	63.5	44.0	49	1645	619
Minimum	38.5	137.0	14	43	34
Maximum	102.0	181.0	63	1688	653





چندک ها Quantiles

صدک ها Centiles

دهک ها Tentiles

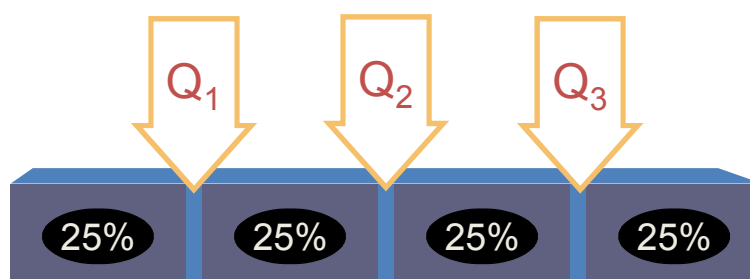
پنجک ها Quintiles

چارک ها Quartiles

## Quartiles

- Three numbers which divide the ordered data into four equal sized groups.
- $Q_1$  has 25% of the data below it.
- $Q_2$  has 50% of the data below it. (Median)
- $Q_3$  has 75% of the data below it.

## Quartiles



## Cholesterol

$$L(M) = (53+1)/2 = 27$$

$$L(Q_1) = (26+1)/2 = 13.5$$

100	124	148	170	185	215
101	125	150	170	185	220
106	127	150	172	186	260
106	128	152	175	187	
110	130	155	175	192	
110	130	157	180	194	
119	133	165	180	195	
120	135	165	180	203	
120	139	165	180	210	
123	140	170	185	212	

## Cholesterol Quartiles

- $Q_1 = 127.5$
- $Q_2 = 165$  (Median)
- $Q_3 = 185$

## Five-Number Summary

- minimum = 100
- $Q_1 = 127.5$
- $M = 165$
- $Q_3 = 185$
- maximum = 260

$$\left. \begin{array}{l} \\ \\ \\ \end{array} \right\} \begin{array}{l} \text{Interquartile} \\ \text{Range (IQR)} \\ = Q_3 - Q_1 \\ = 57.5 \end{array}$$

IQR gives spread of middle 50% of the data

## ارتباط مصرف فرآورده های لبنی با سندرم متابولیک و اجزای آن در

نوجوانان: مطالعه قند و لیپید تهران

**لطفا مفهوم اعداد گزارش شده برای تری گلیسرید را بیان کنید.**

جدول ۱- ویژگی های آمارنگاری و اجزای سندرم متابولیک افراد شرکت کننده در مطالعه در دو گروه سالم و مبتلا به سندرم متابولیک

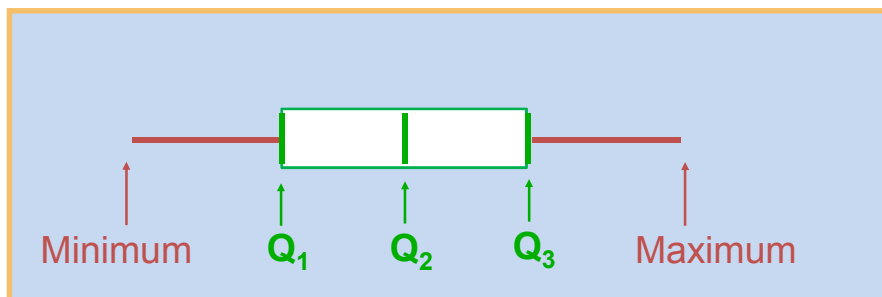
P <sup>†</sup>	مبتلا به سندرم متابولیک		سالم	
	۱۷۴ نفر			
<۰/۰۰۱	۱۴/۲±۲/۸	۱۴/۹±۲/۹		سن (سال)
<۰/۰۰۱	۲۵/۲±۴/۴	۲۱/۱±۴/۲		نمایه ی توده ی بدن (کیلوگرم بر متر مربع)
<۰/۰۰۱	۹۴/۶±۷/۸	۹۱/۸±۸/۱		گلوکز ناشتای خون (میلی گرم در صد میلی لیتر)
<۰/۰۰۱	۱۲۱ (۱۱۳-۱۶۶)	۷۴ (۶۰-۹۴)		تری گلیسرید خون (میلی گرم در صد میلی لیتر) <sup>‡</sup>
<۰/۰۰۱	۴۰/۹±۶/۹	۵۲/۵±۱۰/۳		کلسترول - HDL خون (میلی گرم در صد میلی لیتر)
<۰/۰۰۱	۸۶/۳±۱۰/۹	۷۴/۲±۱۱/۵		دورکمر (سانتی متر)
<۰/۰۰۱	۱۰۷±۱۳/۰	۱۰۰±۱۱/۶		فشار خون سیستولی (میلی متر جیوه)
<۰/۰۰۱	۷۰/۳±۱۱/۶	۶۶/۲±۹/۴		فشار خون دیاستولی (میلی متر جیوه)

\* داده ها به استثنای تری گلیسرید صورت میانگین±انحراف معیار است. † مقادیر P به استثنای تری گلیسرید با استفاده از آزمون آماری تی به دست آمد.  
‡ برای تری گلیسرید میانه و دامنه میان چارکی گزارش شده است.

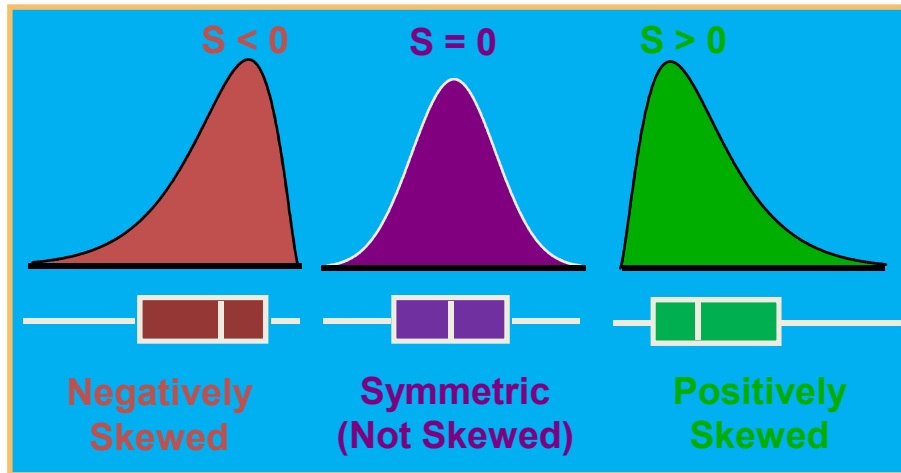
## Boxplot

- Central box spans  $Q_1$  and  $Q_3$ .
- A line in the box marks the median  $M$ .
- Lines extend from the box out to the minimum and maximum.

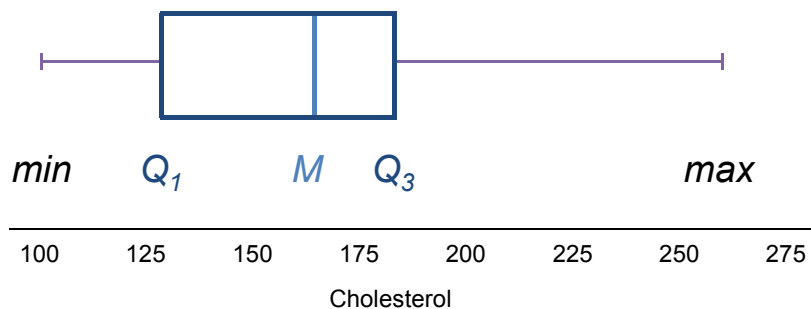
## Box Plot



## Skewness: Box Plots, and Coefficient of Skewness



## Cholesterol: Boxplot



## Identifying Outliers

- The central box of a boxplot spans  $Q_1$  and  $Q_3$ ; recall that this distance is the Interquartile Range ( $IQR$ ).
- We call an observation a suspected **outlier** if it falls more than  $1.5 \times IQR$  above the third quartile or below the first quartile.

## چندک ها

Quantiles					
	WEIGHT	HEIGHT	SYSTOLIC BP	TG	CHOLESTEROL
Median	67.250	158.500	120.00	148.50	192.50
Percentiles					
5	47.000	145.750	95.00	63.00	126.10
10	51.550	148.000	100.00	71.00	146.10
20	56.000	151.000	105.00	94.00	157.00
25	58.000	152.000	110.00	102.00	163.00
30	60.000	153.000	110.00	110.30	168.00
40	64.000	156.000	110.00	129.40	179.40
50	67.250	158.500	120.00	148.50	192.50
60	70.000	161.000	120.00	169.00	200.60
70	74.000	164.000	130.00	198.00	219.70
75	76.125	166.000	136.25	209.00	226.00
80	78.000	168.000	140.00	223.00	231.00
90	83.900	171.000	150.00	302.10	251.90
95	90.000	174.250	170.00	434.50	271.00

## لطفا مفهوم اعداد گزارش شده برای کلسیم را بیان کنید.

جدول ۳- میانگین دریافت مواد مغذی بر اساس چارک‌های دریافت کل لبنیات\*

p*	چارک‌های دریافت لبنیات <sup>†</sup>				۱	۲	۳	۴
	۱	۲	۳	۴				
<./۰.۰۱	۲۸۹۵ ± ۵۸/۰	۲۵۴۵ ± ۵۷/۸	۲۵۶۱ ± ۵۷/۶	۲۹۲۳ ± ۵۷/۹	انرژی دریافتی (کیلوکالری در روز)			
<./۰.۰۱	۹۶/۵ ± ۱/۳	۹۲/۰ ± ۱/۳	۸۸/۳ ± ۱/۳	۸۲/۴ ± ۱/۳	دانسیتتهی انرژی			
-./۰.۴۶	۵۶/۹ ± ۱/۰	۵۷/۲ ± ۱/۰	۵۸/۲ ± ۱/۰	۶۰/۵ ± ۱/۰	کربوهیدرات (درصد از انرژی دریافتی)			
-./۰.۱۲	۱۷/۴ ± ۱/۰	۱۴/۴ ± ۱/۰	۱۳/۹ ± ۱/۰	۱۳/۰ ± ۱/۰	پروتئین (درصد از انرژی دریافتی)			
-./۰.۷۹	۳۷/۴ ± ۲/۳	۳۱/۲ ± ۲/۳	۳۱/۳ ± ۲/۳	۲۹/۶ ± ۲/۳	چربی (درصد از انرژی دریافتی)			
-./۰.۳۷	۱۷/۲ ± ۲/۳	۱۰/۸ ± ۲/۳	۹/۹ ± ۲/۳	۸/۶ ± ۲/۳	چربی اشباع (درصد از انرژی دریافتی)			
-./۰.۷۱۵	۲۰/۸ ± ۲/۵	۱۶/۸ ± ۲/۵	۱۸/۴ ± ۲/۵	۱۹/۱ ± ۲/۵	فیبر غذایی (گرم/۱۰۰۰ کیلوکالری)			
<./۰.۰۱	۱۱۲ ± ۳/۹	۹۸/۸ ± ۳/۹	۹۱/۵ ± ۳/۹	۸۱/۵ ± ۳/۹	کلسترول (میلی‌گرم/۱۰۰۰ کیلوکالری)			
<./۰.۰۱	۷۴۳ ± ۱۰/۶	۵۹۵ ± ۱۰/۵	۵۳۹ ± ۱۰/۵	۴۳۰ ± ۱۰/۵	کلسیم (میلی‌گرم/۱۰۰۰ کیلوکالری)			
<./۰.۰۱	۷۸۹ ± ۶/۲	۶۸۸ ± ۶/۱	۶۳۱ ± ۶/۱	۵۷۱ ± ۶/۲	فسفر (میلی‌گرم/۱۰۰۰ کیلوکالری)			

\* داده‌ها به صورت میانگین ± انحراف معیار است. مقادیر با استفاده از آنالیز کوواریانس (ANCOVA) پس از تعدیل برای سن و جنس به دست آمد. † چارک‌های دریافت کل لبنیات به ترتیب از کمترین به بیشترین عبارت بود از: <۲۱۱/۳، ۲۲۲/۸، ۲۳۳/۵، ۲۴۴/۸، ۲۵۵/۱۲، ۲۶۶/۱۵، ۲۷۷/۱۸، ۲۸۸/۲۱، ۳۰۰/۲۴، ۳۱۱/۲۷، ۳۲۲/۳۰، ۳۳۳/۳۳، ۳۴۴/۳۶، ۳۵۵/۳۹، ۳۶۶/۴۲، ۳۷۷/۴۵، ۳۸۸/۴۸، ۴۰۰/۵۱ از نظر آماری معنی‌دار است.

روش‌های آمار توصیفی

## خلاصه سازی داده‌ها

### شاخصهای پراکندگی

دامنه  
تغییرات  
Range

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

واریانس  
Variance

انحراف معیار  
Standard  
Deviation

ضریب تغییرات  
Coefficient of  
Variation

$$CV = \left( \frac{S}{\bar{X}} \right) \cdot 100\%$$

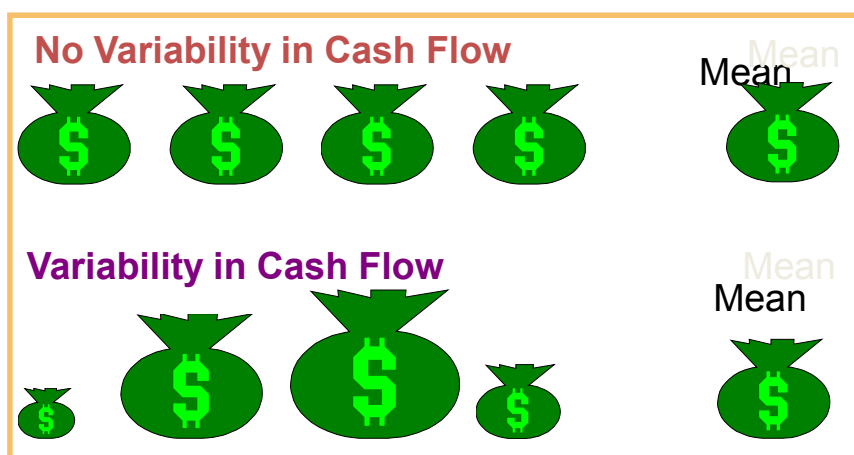
روش‌های آمار توصیفی

۳۲



# Variance and Standard Deviation

## Variability



## Variance: a measure of how data points differ from the mean

- Data Set 1: 3, 5, 7, 10, 10  
Data Set 2: 7, 7, 7, 7, 7

What is the mean and median of the above data set?

Data Set 1: mean = 7, median = 7

Data Set 2: mean = 7, median = 7

But we know that the two data sets are not identical! The **variance** shows how they are different.

We want to find a way to represent these two data set numerically.

## How to Calculate sd?

	Score X	$X - \bar{X}$	$(X - \bar{X})^2$
1	3		
2	5		
3	7		
4	10		
5	10		
Totals	35		

The mean is  $35/5=7$ .

## How to Calculate sd?

	Score $X$	$X - \bar{X}$	$(X - \bar{X})^2$
1	3	$3-7=-4$	
2	5	$5-7=-2$	
3	7	$7-7=0$	
4	10	$10-7=3$	
5	10	$10-7=3$	
Totals	35		

## How to Calculate sd?

	Score $X$	$X - \bar{X}$	$(X - \bar{X})^2$
1	3	$3-7=-4$	16
2	5	$5-7=-2$	4
3	7	$7-7=0$	0
4	10	$10-7=3$	9
5	10	$10-7=3$	9
Totals	35		38

## How to Calculate sd?

	Score X	$x - \bar{x}$	$(x - \bar{x})^2$
1	3	3-7=-4	16
2	5	5-7=-2	4
3	7	7-7=0	0
4	10	10-7=3	9
5	10	10-7=3	9
Totals	35		38

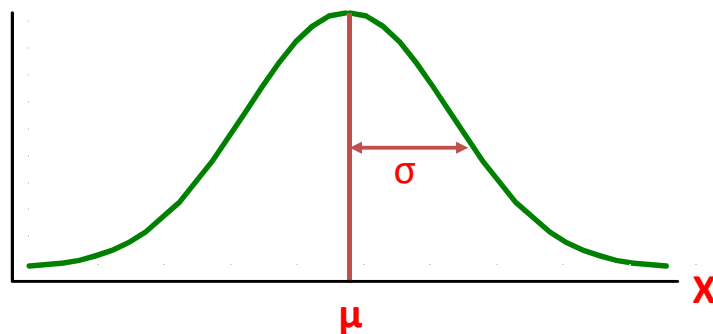
$$sd = \sqrt{\frac{\sum_{i=1}^n (x_i - \bar{x})^2}{(n-1)}} = \sqrt{\frac{38}{5-1}} = \sqrt{9.5} = 3.1$$

چهارشنبه، سپتامبر ۱۴، ۲۰۱۶

روش های آمار توصیفی

۳۹

## The Population Mean and Standard Deviation



چهارشنبه، سپتامبر ۱۴، ۲۰۱۶

روش های آمار توصیفی

40

$$\sigma^2 = \frac{\sum(x - \bar{X})^2}{N} \quad \text{For population variance}$$

$$s^2 = \frac{\sum(x - \bar{X})^2}{n - 1} \quad \text{For sample variance}$$

### Example: Age of Patients

Male	Female
28	27
22	27
21	28
26	6
18	27

Find the mean, median, range?

<b>mean</b>	<b>23</b>	<b>23</b>
<b>median</b>	<b>22</b>	<b>27</b>
<b>range</b>	<b>10</b>	<b>22</b>

What can be said about this data?

**Due to the outlier, the median is more typical of overall performance.**

Which diver was more consistent?

Dive	Male Age X	$X - \bar{X}$	$(X - \bar{X})^2$
1	28	5	25
2	22	-1	1
3	21	-2	4
4	26	3	9
5	18	-5	25
Totals	115	0	64

$$\text{Male age sd} = \sqrt{\frac{64}{5-1}} = \sqrt{16} = 4$$

$$\text{Female age sd} = \sqrt{\frac{362}{5-1}} = \sqrt{90.5} = 9.5$$

## Example: Sample Standard Deviation

- Square root of the sample variance

$X$	$X - \bar{X}$	$(X - \bar{X})^2$
2,398	625	390,625
1,844	71	5,041
1,539	-234	54,756
<u>1,311</u>	<u>-462</u>	<u>213,444</u>
7,092	0	663,866

$$\begin{aligned}
 S^2 &= \frac{\sum (X - \bar{X})^2}{n - 1} \\
 &= \frac{663,866}{3} \\
 &= 221,288.67 \\
 S &= \sqrt{S^2} \\
 &= \sqrt{221,288.67} \\
 &= 470.41
 \end{aligned}$$

## Example

Metabolic rates of 7 men (cal./24hr.) :

1792 1666 1362 1614 1460 1867 1439

$$\begin{aligned}\bar{x} &= \frac{1792+1666+1362+1614+1460+1867+1439}{7} \\ &= \frac{11,200}{7} \\ &= 1600\end{aligned}$$

## Example

Observations $x_i$	Deviations $x_i - \bar{x}$	Squared deviations $(x_i - \bar{x})^2$
1792	1792-1600 = 192	(192) <sup>2</sup> = 36,864
1666	1666 -1600 = 66	(66) <sup>2</sup> = 4,356
1362	1362 -1600 = -238	(-238) <sup>2</sup> = 56,644
1614	1614 -1600 = 14	(14) <sup>2</sup> = 196
1460	1460 -1600 = -140	(-140) <sup>2</sup> = 19,600
1867	1867 -1600 = 267	(267) <sup>2</sup> = 71,289
1439	1439 -1600 = -161	(-161) <sup>2</sup> = 25,921
	sum = 0	sum = 214,870

## Variance and Standard Deviation

Example from Text

$$s^2 = \frac{214,870}{7-1} = 35,811.67$$

$$s = \sqrt{35,811.67} = 189.24 \text{ calories}$$

## Standard Deviation      انحراف معیار

• مهم ترین شاخص پراکندگی

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



## مقایسه انحراف معیار داده ها

Data A



Mean = 15.5

Data B



Mean = 15.5

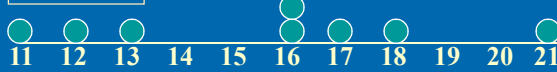
Data C



Mean = 15.5

## مقایسه انحراف معیار داده ها

Data A



Mean = 15.5  
sd = 3.338

Data B



Mean = 15.5  
sd = .9258

Data C



Mean = 15.5  
sd = 4.57

## SYSTOLIC BP

Age Groups	N	Mean	Std. Deviation	Minimum	Maximum	Range
<35	54	109.54	11.421	90	140	50
35-44	60	112.08	17.303	85	180	95
45-54	52	125.87	20.430	100	185	85
55-64	54	136.30	22.554	100	180	80
65+	30	134.33	23.034	100	210	110
Total	250	122.30	21.848	85	210	125

## ضریب تغییرات Coefficient of Variation

- پراکندگی نسبت به میانگین را نشان می دهد
- مقایسه پراکندگی دو یا چند گروه
- مقایسه پراکندگی دو متغیر با واحد اندازه گیری متفاوت

$$CV = \left( \frac{S}{\bar{X}} \right) \cdot 100\%$$

## Coefficient of Variation ضریب تغییرات

$$\bar{X}=3000 \quad Sd=300$$

وزن نوزادان

$$\bar{X} = 40Kg \quad Sd = 2$$

وزن افراد ۱۲ ساله

## Coefficient of Variation ضریب تغییرات

$$\bar{X} = 3000 \quad Sd = 300$$

$$\bar{X} = 3000 \text{ gr} \quad Sd = 300 \text{ gr} \quad CV = \frac{300}{3000} \times 100 = 10\%$$

وزن نوزادان

$$\bar{X} = 40Kg \quad Sd = 2$$

وزن افراد ۱۲ ساله

$$\bar{X} = 40kg \quad Sd = 2kg \quad CV = \frac{2}{40} \times 100 = 5\%$$

ضریب تغییرات داده های زیر را محاسبه کنید. کدام یک پراکندگی بیشتری دارد؟

SYSTOLIC BP

Age Groups	N	Mean	Std. Deviation
<35	54	109.54	11.421
35-44	60	112.08	17.303
45-54	52	125.87	20.430
55-64	54	136.30	22.554
.=65	30	134.33	23.034
Total	250	122.30	21.848

۵۵

ضریب تغییرات داده های زیر را محاسبه کنید. کدام یک پراکندگی بیشتری دارد؟

SYSTOLIC BP

Age Groups	N	Mean	Std. Deviation	CV%
<35	54	109.54	11.421	10.43
35-44	60	112.08	17.303	15.44
45-54	52	125.87	20.430	16.23
55-64	54	136.30	22.554	16.55
.=65	30	134.33	23.034	17.15
Total	250	122.30	21.848	17.86

## Weighted Mean میانگین وزنی

$n_1$	$n_2$	.....	$n_k$	تعداد در گروه $j$
$\bar{X}_1$	$\bar{X}_2$	.....	$\bar{X}_k$	میانگین گروه $j$

$$\bar{X} = \frac{1}{n} \sum_{j=1}^K n_j \bar{X}_j$$

$$n = \sum_{j=1}^k n_j$$

## نسبت وزنی

$n_1$	$n_2$	.....	$n_k$	تعداد در گروه $j$
$p_1$	$p_2$	.....	$p_k$	نسبت در گروه $j$

$$P = \frac{1}{n} \sum_{j=1}^K n_j p_j$$

$$n = \sum_{j=1}^k n_j$$