



and frequencies  $\{f_1, f_2, \dots, f_m\}$  respectively, the **population variance** is given by :

$$\sigma^2 = (O_1 - \mu)^2 \frac{f_1}{n} + (O_2 - \mu)^2 \frac{f_2}{n} + \dots + (O_m - \mu)^2 \frac{f_m}{n}.$$

Observations	Deviation	Squared Deviation	Sq. Dev. $\times$ Rel. Freq.
$O_i$	$O_i - \mu$	$(O_i - \mu)^2$	$(O_i - \mu)^2 \frac{f_i}{n}$
$O_1$	$O_1 - \mu$	$(O_1 - \mu)^2$	$(O_1 - \mu)^2 \frac{f_1}{n}$
$O_2$	$O_2 - \mu$	$(O_2 - \mu)^2$	$(O_2 - \mu)^2 \frac{f_2}{n}$
$\vdots$	$\vdots$	$\vdots$	$\vdots$
$O_m$	$O_m - \mu$	$(O_m - \mu)^2$	$(O_m - \mu)^2 \frac{f_m}{n}$
			$\sigma^2 = \text{Sum}$

**Population Standard Deviation:** The **population standard deviation** for the data is the square root of the population variance,

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

**Example 1:** Find the variance,  $\sigma^2$  and standard deviation,  $\sqrt{\sigma^2}$ , for the number of completions for each Quarterback above. The value of  $\mu$  is 28 for both players.

Quarterback A

$O_i$ # Completions	$f_i$ Frequency
19	1
22	1
24	1
25	2
27	2
28	2
29	1
30	2
32	2
33	1
37	1

Quarterback B

$O_i$ # Completions	$f_i$ Frequency
12	2
17	2
19	1
25	1
30	2
32	1
33	2
35	2
39	2
40	1

**Quarterback A:**

Observations	Deviation	Squared Deviation	Sq. Dev. $\times$ Rel. Freq.
$O_i$	$O_i - 28$	$(O_i - 28)^2$	$(O_i - 28)^2 \frac{f_i}{16}$
19			
22			
24			
25			
27			
28			
29			
30			
32			
33			
37			
			$\sigma^2 = \text{Sum}$

**Quarterback A:**  $\sigma^2 = \underline{\hspace{2cm}}$      $\sigma = \underline{\hspace{2cm}}$

**Quarterback B:**

Observations	Deviation	Squared Deviation	Sq. Dev. $\times$ Rel. Freq.
$O_i$	$O_i - 28$	$(O_i - 28)^2$	$(O_i - 28)^2 \frac{f_i}{16}$
12			
17			
19			
25			
30			
32			
33			
35			
39			
40			
			$\sigma^2 = \text{Sum}$

**Quarterback B:**  $\sigma^2 = \underline{\hspace{2cm}}$      $\sigma = \underline{\hspace{2cm}}$

**Sample Variance and Standard Deviation :** If we calculate the variance according to the formula given above, for a sample from a particular population, it is not accurate (biased) as an estimate for the population variance. So for a sample from a given population, we use the **sample variance** as an unbiased estimator of the population variance.

Given a sample,  $\{x_1, x_2, \dots, x_n\}$ , of size  $n$  from a population, where the sample mean is given by  $\bar{x}$ , the **sample variance** is given by

$$s^2 = \frac{(x_1 - \bar{x})^2 + (x_2 - \bar{x})^2 + \dots + (x_n - \bar{x})^2}{n - 1}.$$

The **sample standard deviation** is given by

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

If the data is given in a frequency distribution, we can shorten the calculations. If the outcomes in the sample are given by  $\{O_1, O_2, \dots, O_m\}$  with respective frequencies given by  $\{f_1, f_2, \dots, f_m\}$ , then

$$s^2 = \frac{(O_1 - \bar{x})^2 f_1 + (O_2 - \bar{x})^2 f_2 + \dots + (O_m - \bar{x})^2 f_m}{n - 1}.$$

**Example 3** A random sample of size twenty of a golfer's scores for nine-hole rounds of golf over the past year are as follows:

39, 40, 40, 41, 39, 40, 44, 43, 40, 41, 40, 41, 41,  
42, 43, 40, 41, 41, 41, 43.

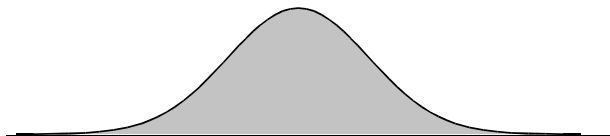
Compute the mean, sample variance and the sample standard deviation for the sample of the golfer's scores. You can view the sample variance as an estimate of the overall variance of the golfer's scores.

## Interpreting The Standard Deviation

When presented with raw scores for performance, it is difficult to interpret their meaning without some measure of center and variability for the population from which they come. In any set of data, whether it is population data or a sample, observations that are more than 3 standard deviations from the mean are rare and exceptional. One such rule demonstrating this is the empirical rule for mound shaped data shown below. We will explore this rule in more detail when we study the normal distribution.

### The Empirical Rule for Mound Shaped Data

The empirical rule given below applies to data sets with frequency distributions that are mound shaped and symmetric, like the one shown below.



**Empirical Rule** If the data has a frequency distribution which is mound shaped and symmetric, we have the following empirical rule:

- Approximately 68% of the measurements will fall within 1 standard deviation of the mean i.e. within the interval  $(\bar{x} - s, \bar{x} + s)$  for a sample and  $(\mu - \sigma, \mu + \sigma)$  for a population.
- Approximately 95% of the measurements will fall within 2 standard deviations of the mean, i.e. within the interval  $(\bar{x} - 2s, \bar{x} + 2s)$  for samples and  $(\mu - 2\sigma, \mu + 2\sigma)$  for a population.
- Approximately 99.7% of the measurements(essentially all) will fall within 3 standard deviations of the mean, i.e. within the interval  $(\bar{x} - 3s, \bar{x} + 3s)$  for samples and  $(\mu - 3\sigma, \mu + 3\sigma)$  for a population.

Mound shaped distributions are very important because they frequently occur as population distributions. Even more importantly, the central limit theorem says that if we take all samples of a given size from a population and calculate all of the means, then the distribution of the means is mound shaped (Normal). We will study Normal distributions in more detail later.

### Numerical Measures of Relative Standing

Quite often when interpreting a data observation, such as a baby's height and weight, we are interested in how it compares to the rest of the relevant population. Measures of relative standing describe the location of a particular measurement relative to the rest of the data. We explore some of the standard measures of relative standing below.

**Z-Scores** The z-score for a particular measurement in a set of data, measures how many standard deviations that measurement lies away from the mean.

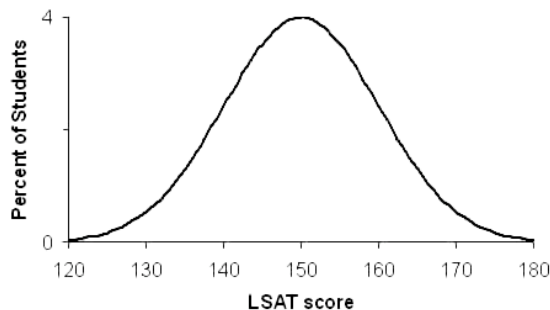
**Definition** The **sample z-score** for a measurement  $x$  in a set of data is

$$z = \frac{x - \bar{x}}{s}.$$

The **population z-score** for a data measurement,  $x$ , is

$$z = \frac{x - \mu}{\sigma}.$$

**Example** The scores on the LSAT for a particular year have a mound shaped distribution. The mean is  $\mu = 150$  and the standard deviation is  $\sigma = 10$ . The distribution is shown below.



(a) Use the empirical rule to determine what percentage of prospective law students have z scores between -2 and 2.

(b) If you scored 175 on the exam, what would your z-score be?

**Example** In 2013 Mary was among the college bound seniors who took the SAT and ACT exams. Her composite score on the SAT was 2500 and her composite score on the ACT was 34. The national average for the composite score on the SAT among college bound seniors for that year was 1499 and the standard deviation was 319. The national average for the ACT among college bound seniors for that year was 22.5 and the standard deviation was 4.9.

Find Mary's z-score for both exams and use the z-scores to compare Mary's performance on both exams.

**Rank** We can also use rank to measure relative standing, by ranking the data as 1st, 2nd, 3rd, according to the size of the data measurement. This is commonly used in racing, where a lower time leads to a higher position, also in many competitions a higher number of points or wins leads to a higher rank. When two data measurements are the same (a tie) we can give both the same rank and skip a rank. A closely related measure of relative standing is given by the percentile:

### Percentiles

Recall that the **median** of a set of data is number for which 50% of the measurements lie at or below

the median and 50% lie at or above it. This is the 50th percentile of the distribution.

For any set of  $n$  measurements, (arranged in ascending or descending order), the  $p$ th percentile is a number such that  $p\%$  of the measurements fall at or below that number and  $(100 - p)\%$  of the measurements fall above it. The calculation of percentiles is not well defined and there are a few conventions for one might adopt for choosing a value for a percentile. We will adopt a relatively simple convention which agrees with our calculation of the median from before.

For a set of data of size  $N$ , **to calculate the  $P$ -th percentile** we :

- order our data from smallest to largest,
- and calculate  $\frac{P}{100} \times N$ .
- If  $\frac{P}{100} \times N$  is not a whole number, we round it up to the nearest integer  $n$  and the  $n$  th data point on the list is the  $P$  th percentile.
- If  $\frac{P}{100} \times N$  is a whole number  $n$ , we calculate the average of the  $n$  th data point and the  $n + 1$  th data point to get the  $P$  th percentile.

**Example 1** Find the 10th Percentile, the 25th percentile, the 50th percentile, the 75th percentile and the 90th percentile of the following set of 20 exam scores:

60, 71, 85, 99, 100, 76, 98, 61, 75, 82, 95, 72, 88, 61, 72, 80, 100, 90, 60, 70.

**Example** Attached you will find list of the top forty players from the NBA with the number of rebounds for the 2013-2014 regular season for each given in the highlighted column. The players are ranked 1-40 according to the number of rebounds.

What is the 95-th percentile for the number of rebounds among all NBA players for the 2013-2014 regular season.

**Rebounds Leaders - All Players**

RK	PLAYER	TEAM	GP	MPG	OFF	ORPG	DEF	DRPG	REB	RPG	RP48
1	DeAndre Jordan, C	LAC	82	35.0	331	4.0	783	9.5	<b>1114</b>	13.6	18.6
2	Andre Drummond, C	DET	81	32.3	440	5.4	631	7.8	<b>1071</b>	13.2	19.6
3	Kevin Love, PF	MIN	77	36.3	224	2.9	739	9.6	<b>963</b>	12.5	16.5
4	Joakim Noah, C	CHI	80	35.3	282	3.5	618	7.7	<b>900</b>	11.3	15.3
5	Dwight Howard, C	HOU	71	33.7	231	3.3	635	8.9	<b>866</b>	12.2	17.3
6	DeMarcus Cousins, C	SAC	71	32.4	218	3.1	613	8.6	<b>831</b>	11.7	17.4
7	Zach Randolph, PF	MEM	79	34.2	265	3.4	530	6.7	<b>795</b>	10.1	14.1
8	Al Jefferson, C	CHA	73	35.0	156	2.1	636	8.7	<b>792</b>	10.8	14.9
9	Marcin Gortat, C	WSH	81	32.8	202	2.5	565	7.0	<b>767</b>	9.5	13.9
10	LaMarcus Aldridge, PF	POR	69	36.2	166	2.4	600	8.7	<b>766</b>	11.1	14.7
RK	PLAYER	TEAM	GP	MPG	OFF	ORPG	DEF	DRPG	REB	RPG	RP48
11	Greg Monroe, PF	DET	82	32.8	256	3.1	504	6.1	<b>760</b>	9.3	13.6
12	Blake Griffin, PF	LAC	80	35.8	192	2.4	565	7.1	<b>757</b>	9.5	12.7
13	Tristan Thompson, PF	CLE	82	31.6	269	3.3	485	5.9	<b>754</b>	9.2	14.0
14	Tim Duncan, PF	SA	74	29.2	158	2.1	563	7.6	<b>721</b>	9.7	16.0
15	Jonas Valanciunas, C	TOR	81	28.2	226	2.8	488	6.0	<b>714</b>	8.8	15.0
16	Serge Ibaka, PF	OKC	81	32.9	224	2.8	485	6.0	<b>709</b>	8.8	12.8
17	Robin Lopez, C	POR	82	31.8	326	4.0	373	4.5	<b>699</b>	8.5	12.9
18	Kenneth Faried, PF	DEN	80	27.2	238	3.0	446	5.6	<b>684</b>	8.6	15.1
19	Anthony Davis, PF	NO	67	35.2	207	3.1	466	7.0	<b>673</b>	10.0	13.7
20	Andrew Bogut, C	GS	67	26.4	182	2.7	489	7.3	<b>671</b>	10.0	18.2
RK	PLAYER	TEAM	GP	MPG	OFF	ORPG	DEF	DRPG	REB	RPG	RP48
21	Spencer Hawes, PF	CLE/PHI	80	30.9	131	1.6	529	6.6	<b>660</b>	8.3	12.8
22	David Lee, PF	GS	69	33.2	182	2.6	461	6.7	<b>643</b>	9.3	13.5
23	Derrick Favors, C	UTAH	73	30.2	199	2.7	438	6.0	<b>637</b>	8.7	13.9
24	J.J. Hickson, C	DEN	69	26.9	206	3.0	426	6.2	<b>632</b>	9.2	16.3
	Carlos Boozer, PF	CHI	76	28.2	137	1.8	495	6.5	<b>632</b>	8.3	14.2
26	Anderson Varejao, C	CLE	65	27.7	187	2.9	442	6.8	<b>629</b>	9.7	16.8
27	Paul Millsap, PF	ATL	74	33.5	154	2.1	473	6.4	<b>627</b>	8.5	12.1
28	Nikola Vucevic, C	ORL	57	31.8	185	3.2	441	7.7	<b>626</b>	11.0	16.6
	Miles Plumlee, C	PHX	80	24.6	198	2.5	428	5.4	<b>626</b>	7.8	15.3
30	Carmelo Anthony, SF	NY	77	38.7	145	1.9	477	6.2	<b>622</b>	8.1	10.0
RK	PLAYER	TEAM	GP	MPG	OFF	ORPG	DEF	DRPG	REB	RPG	RP48
31	Nicolas Batum, SF	POR	82	36.0	116	1.4	495	6.0	<b>611</b>	7.5	9.9
32	Jared Sullinger, C	BOS	74	27.6	241	3.3	360	4.9	<b>601</b>	8.1	14.1
33	Enes Kanter, C	UTAH	80	26.7	222	2.8	376	4.7	<b>598</b>	7.5	13.4
	Kevin Durant, SF	OKC	81	38.5	58	0.7	540	6.7	<b>598</b>	7.4	9.2
35	Pau Gasol, PF	LAL	60	31.4	124	2.1	456	7.6	<b>580</b>	9.7	14.8
36	Lance Stephenson, SG	IND	78	35.3	95	1.2	463	5.9	<b>558</b>	7.2	9.7
	Taj Gibson, PF	CHI	82	28.7	200	2.4	358	4.4	<b>558</b>	6.8	11.4
38	David West, PF	IND	80	30.9	120	1.5	422	5.3	<b>542</b>	6.8	10.5
	Paul George, SF	IND	80	36.2	64	0.8	478	6.0	<b>542</b>	6.8	9.0
40	Samuel Dalembert, C	DAL	80	20.2	200	2.5	341	4.3	<b>541</b>	6.8	16.1