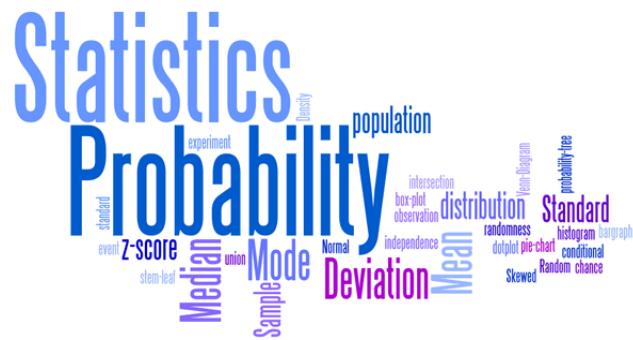


MSIT 431 Probability and Statistical Methods

Chapter 1 Looking at Data

Dongning Guo

Fall 2018



- Statistics is the science of data.
 - Probability provides the mathematical foundation.



Outline

1.1 Data

1.2 Displaying Distributions with Graphs

1.3 Describing Distributions with Numbers

1.4 Density Curves and Normal Distributions

3

1.1 Data

Key characteristics of a data set.



4

2

Data set examples

1. precipitation:

Los Angeles	Sacramento	San Francisco	Denver
14.0	17.2	20.7	13.0
Hartford	Wilmington	Washington	Jacksonville
43.4	40.2	38.9	54.5
Miami	Atlanta	Honolulu	Boise
59.8	48.3	22.9	11.5
Chicago	Peoria	Indianapolis	Des Moines
34.4	35.1	38.7	30.8

2. Call log:

Date	Number	Time	Duration	Cell Site	Name
1 Thursday, February 18, 1999	#4432539023	9:49:54 PM	0:00:03	BLTM2	Forward to Voicemail
2 Thursday, February 18, 1999	incoming	9:49:54 PM	0:00:03	TPA44	Incoming
3 Thursday, February 18, 1999	incoming	7:29:58 PM	0:00:30	L651C	Incoming
4 Wednesday, February 17, 1999	#4432539023	10:01:10 PM	0:00:01	BLTM2	Forward to Voicemail
5 Wednesday, February 17, 1999	incoming	10:01:10 PM	0:00:01	TPA44	Incoming
6 Wednesday, February 17, 1999	#4432539023	8:47:07 PM	0:00:05	BLTM2	Forward to Voicemail
7 Wednesday, February 17, 1999	incoming	8:47:07 PM	0:00:05	TPA44	Incoming

5

Terminology

- ✓ **Cases** are the objects described by a set of data. Cases may be customers, companies, experimental subjects, or other objects.
- ✓ A **variable** is a special characteristic of a case.
- ✓ A **label** is a special variable used in some data sets to distinguish between cases.
- ✓ Different cases can have different **values** of a variable.

Date	Number	Time	Duration	Cell Site	Name
1 Thursday, February 18, 1999	#4432539023	9:49:54 PM	0:00:03	BLTM2	Forward to Voicemail
2 Thursday, February 18, 1999	incoming	9:49:54 PM	0:00:03	TPA44	Incoming
3 Thursday, February 18, 1999	incoming	7:29:58 PM	0:00:30	L651C	Incoming
4 Wednesday, February 17, 1999	#4432539023	10:01:10 PM	0:00:01	BLTM2	Forward to Voicemail
5 Wednesday, February 17, 1999	incoming	10:01:10 PM	0:00:01	TPA44	Incoming
6 Wednesday, February 17, 1999	#4432539023	8:47:07 PM	0:00:05	BLTM2	Forward to Voicemail
7 Wednesday, February 17, 1999	incoming	8:47:07 PM	0:00:05	TPA44	Incoming

6

Types of variables

- A **categorical** variable places each case into one of several groups, or categories.
- A **quantitative** variable takes numerical values for which arithmetic operations such as adding and averaging make sense.
- The **distribution** of a variable tells us the values that a variable takes and how often it takes each value.

Date	Number	Time	Duration	Cell Site	Name
1 Thursday, February 18, 1999	#4432539023	9:49:54 PM	0:00:03	BLTM2	Forward to Voicemail
2 Thursday, February 18, 1999	incoming	9:49:54 PM	0:00:03	TPA44	Incoming
3 Thursday, February 18, 1999	incoming	7:29:58 PM	0:00:30	L651C	Incoming
4 Wednesday, February 17, 1999	#4432539023	10:01:10 PM	0:00:01	BLTM2	Forward to Voicemail
5 Wednesday, February 17, 1999	incoming	10:01:10 PM	0:00:01	TPA44	Incoming
6 Wednesday, February 17, 1999	#4432539023	8:47:07 PM	0:00:05	BLTM2	Forward to Voicemail
7 Wednesday, February 17, 1999	incoming	8:47:07 PM	0:00:05	TPA44	Incoming

7

Key Characteristics of a Data Set

Every data set is accompanied by important background information. In a statistical study, **always ask the following questions:**

- **Who?** What **cases?** How many cases?
- **What?** How many **variables?** How are they defined? What are the units of measurement?
- **Why?** What purpose? Do the data contain the information needed to answer the questions of interest?

Date	Number	Time	Duration	Cell Site	Name
1 Thursday, February 18, 1999	#4432539023	9:49:54 PM	0:00:03	BLTM2	Forward to Voicemail
2 Thursday, February 18, 1999	incoming	9:49:54 PM	0:00:03	TPA44	Incoming
3 Thursday, February 18, 1999	incoming	7:29:58 PM	0:00:30	L651C	Incoming
4 Wednesday, February 17, 1999	#4432539023	10:01:10 PM	0:00:01	BLTM2	Forward to Voicemail
5 Wednesday, February 17, 1999	incoming	10:01:10 PM	0:00:01	TPA44	Incoming
6 Wednesday, February 17, 1999	#4432539023	8:47:07 PM	0:00:05	BLTM2	Forward to Voicemail
7 Wednesday, February 17, 1999	incoming	8:47:07 PM	0:00:05	TPA44	Incoming

8

1.2 Displaying Distributions with Graphs



- Variables
- Examining distributions of variables
- Graphs for categorical variables
 - Bar graphs
 - Pie charts
- Graphs for quantitative variables
 - Histograms
 - Stemplots
 - Time plots

9

Variables

Exploring Data

- Begin by examining each variable by itself. Then move on to study the relationships among the variables.
- Begin with a graph or graphs. Then add numerical summaries of specific aspects of the data.

10

Distribution of a Variable

- To examine a single variable, we graphically display its **distribution**.

- The distribution of a variable tells us what values it takes and how often it takes these values.
- Distributions can be displayed using a variety of graphical tools. The proper choice of graph depends on the nature of the variable.

Categorical variable

Pie chart
Bar graph

Quantitative variable

Histogram
Stemplot

11

* Using



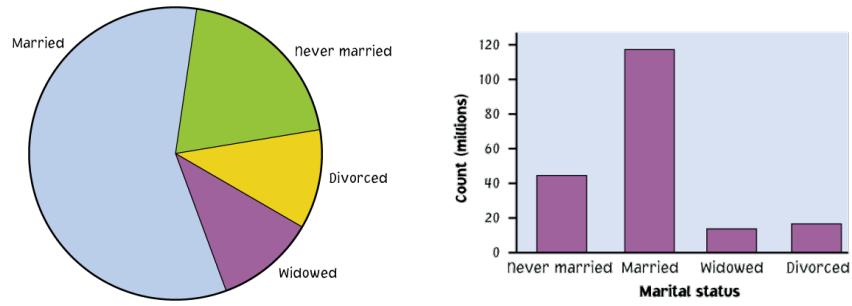
- An elegant open-source software for doing statistics.
- Widely used in industry.
- Numerous packages.
- Heavily used in MSIT 423 Data Mining.
- R will be used in class and in homework. Please install R to your computer. <http://www.R-project.org>
- To learn the basics of R, read “R: A self-learn tutorial” and do the exercises (Homework 1).
- A more detailed reference is:
<http://cran.r-project.org/web/packages/IPSUR/vignettes/IPSUR.pdf>
- Other reference: <http://www.statmethods.net/>

12

Categorical Variables

The **distribution of a categorical variable** lists the categories and gives the **count** or **percent** of individuals who fall into each category.

- **Pie charts** show the distribution of a categorical variable as a “pie” whose slices are sized by the counts or percents for the categories.
- **Bar graphs** represent categories as bars whose heights show the category counts or percents.



13

Quantitative Variables

The **distribution of a quantitative variable** tells us what values the variable takes on and how often it takes those values.

- **Histograms** show the distribution of a quantitative variable by using bars. The height of a bar represents the number of individuals whose values fall within the corresponding class.
- **Stemplots** separate each observation into a stem and a leaf that are then plotted to display the distribution while maintaining the original values of the variable.
- **Time plots** plot each observation against the time at which it was measured.

14

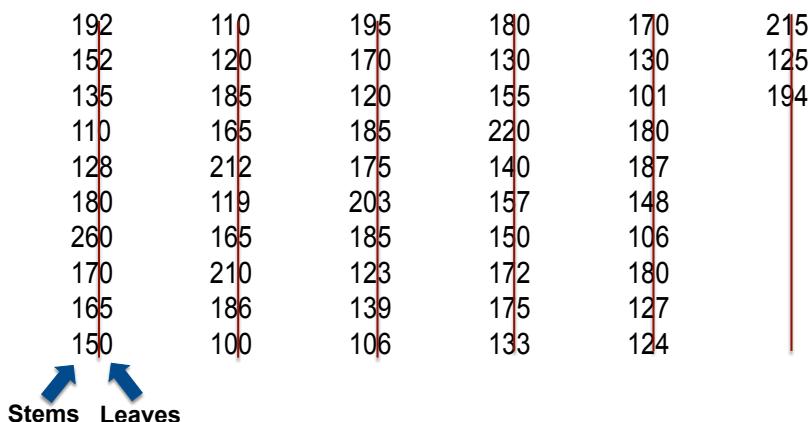
Stemplots

- Separate each observation into a **stem** (first part of the number) and a **leaf** (the remaining part of the number).
- Write the stems in a vertical column; draw a vertical line to the right of the stems.
- Write each leaf in the row to the right of its stem; order leaves if desired.

15

Stemplots

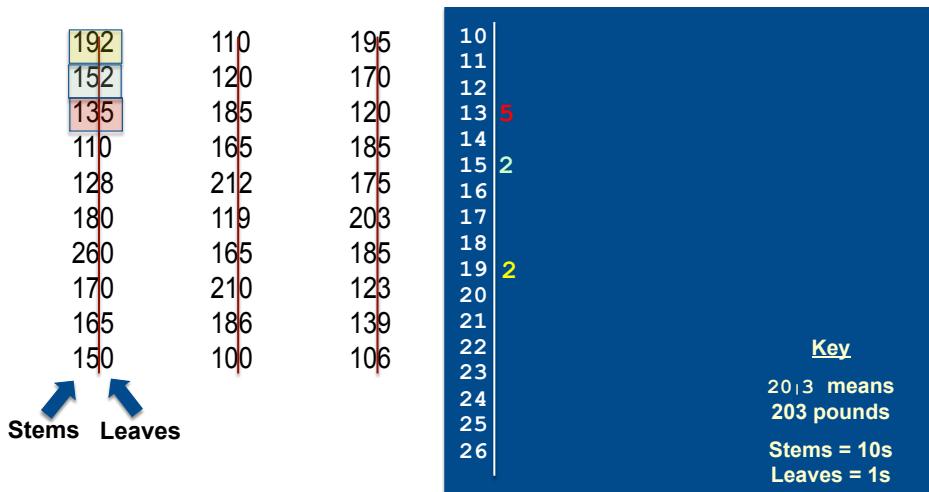
Example: Weight Data of a Class



16

Stemplots

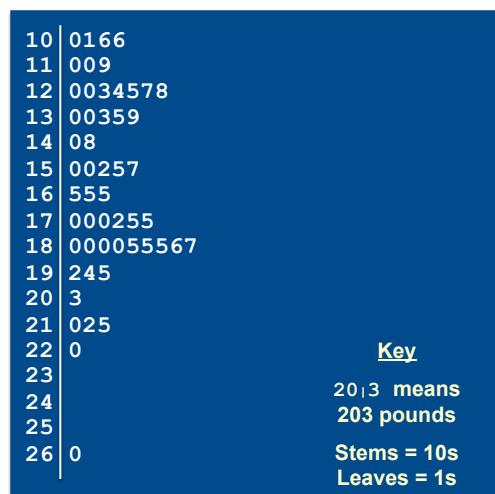
Example: Weight Data of a Class



17

Stemplots

Example: Weight Data of a Class



18

Stemplots

If there are very few stems (when the data cover only a very small range of values), then we may want to create more stems by **splitting** the original stems.

Example: If all of the data values are between 150 and 179, then we may choose to use the following stems:

15	
15	
16	Leaves 0–4 would go on each upper stem (first “15”), and leaves 5–9 would go on each lower stem (second “15”).
16	
17	
17	

19

Histograms

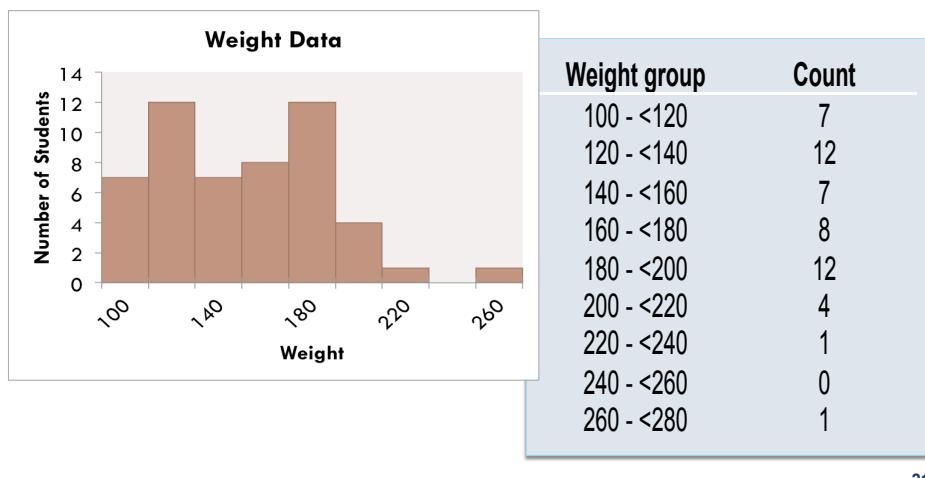
For large datasets and/or quantitative variables that take many values:

- Divide the possible values into **classes**, or **intervals** of equal widths.
- Count how many observations fall into each interval. Instead of counts, one may also use percents.
- Draw a picture representing the distribution—each bar height is equal to the number (percent) of observations in its interval.

20

Histograms

Example: Weight Data



21

Examining Distributions

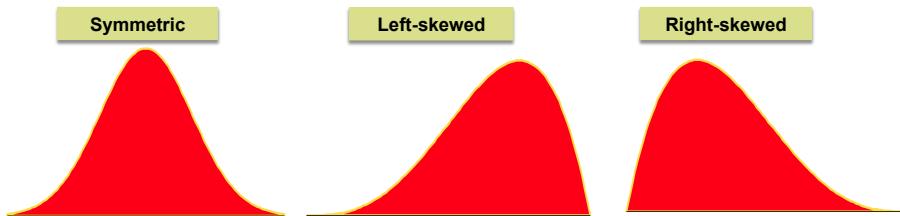
In any graph of data, look for the **overall pattern** and for striking **deviations** from that pattern.

- You can describe the overall pattern by its **shape, center, and spread**.
- An important kind of deviation is an **outlier**, an individual that falls outside the overall pattern.

22

Examining Distributions

- A distribution is **symmetric** if the right and left sides of the graph are approximately mirror images of each other.
- A distribution is **skewed to the left** (left-skewed) if the left side of the graph (containing the half of the observations with smaller values) is much longer than the right side.
- It is **skewed to the right** (right-skewed) if the right side of the graph is much longer than the left side.



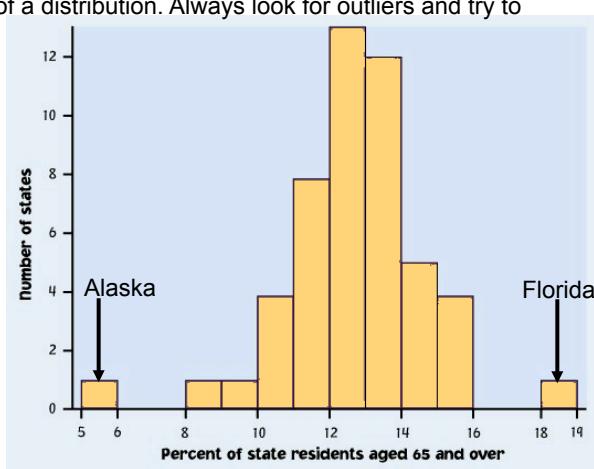
23

Outliers

An important kind of deviation is an **outlier**. Outliers are observations that lie outside the overall pattern of a distribution. Always look for outliers and try to explain them.

The overall pattern is fairly symmetrical except for two states that clearly do not belong to the main pattern. Alaska and Florida have unusually small and large percents, respectively, of elderly residents in their populations.

A large gap in the distribution is typically a sign of an outlier.



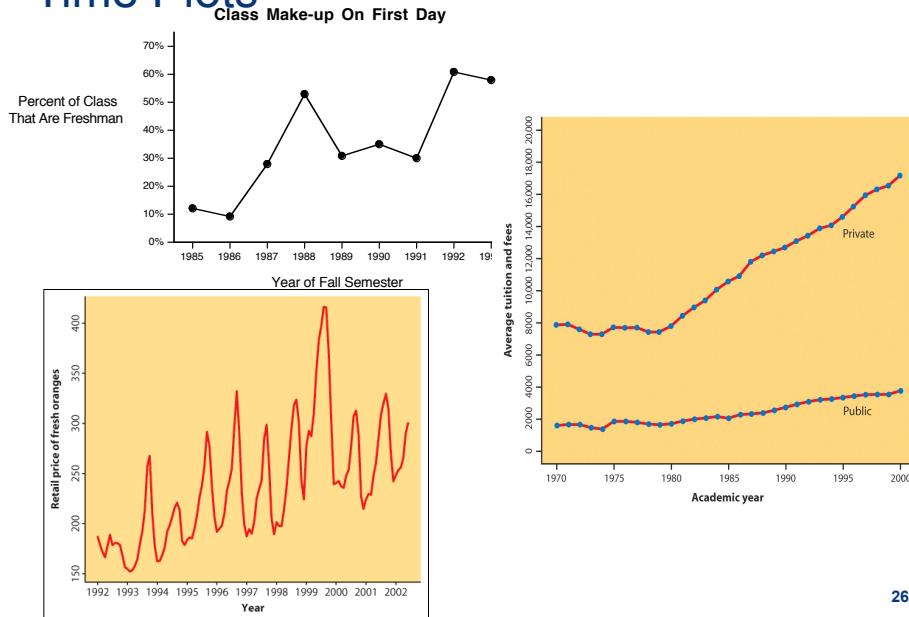
Time Plots

A **time plot** shows behavior over time.

- Time is always on the horizontal axis, and the variable being measured is on the vertical axis.
- Look for an overall pattern (trend) and deviations from this trend. Connecting the data points by lines may emphasize this trend.
- Look for patterns that repeat at known regular intervals (seasonal variations).

25

Time Plots



26.

1.3 Describing Distributions with Numbers



- Measures of center: mean, median
- Mean versus median
- Measures of spread: quartiles, standard deviation
- Five-number summary and boxplot
- Choosing among summary statistics
- Changing the unit of measurement

27

Measuring Center: The Mean

The most common measure of center is the arithmetic average, or **mean**.

To find the **mean** \bar{x} (pronounced “x-bar”) of a set of observations, add their values, and divide by the number of observations. If the n observations are $x_1, x_2, x_3, \dots, x_n$, their mean is:

$$\bar{x} = \frac{\text{sum of observations}}{n} = \frac{x_1 + x_2 + \dots + x_n}{n}$$

In more compact notation:

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

28

Measuring Center: The Median

Because the mean cannot resist the influence of extreme observations, it is not a **resistant measure** of center.

Another common measure of center is the **median**.

The **median M** is the midpoint of a distribution, the number such that half of the observations are smaller and the other half are larger.

To find the median of a distribution:

1. Arrange all observations from smallest to largest.
2. If the number of observations n is odd, the median M is the center observation in the ordered list.
3. If the number of observations n is even, the median M is the average of the two center observations in the ordered list.

29

Measuring Center: Example

Use the data below to calculate the mean and median of the commuting times (in minutes) of 20 randomly selected New York workers.

10	30	5	25	40	20	10	15	30	20	15	20	85	15	65	15	60	60	40	45
----	----	---	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----

$$\bar{x} = \frac{10 + 30 + 5 + 25 + \dots + 40 + 45}{20} = 31.25 \text{ minutes}$$

0	5
1	005 35
2	0005
3	00
4	005
5	
6	005
7	
8	5

Key: 4|5
represents a
New York
worker who
reported a 45-
minute travel
time to work.

$$M = \frac{20 + 25}{2} = 22.5 \text{ minutes}$$

30

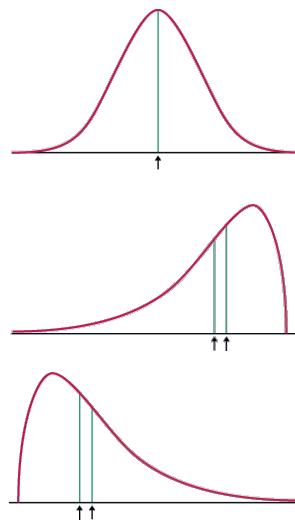
Comparing Mean and Median

The mean and median measure center in different ways, and both are useful.

The mean and median of a roughly **symmetric** distribution are close together.

If the distribution is exactly **symmetric**, the mean and median are exactly the same.

In a **skewed** distribution, the mean is usually farther out in the long tail than is the median.



31

Measuring Spread: The Quartiles

- A measure of center alone can be misleading.
- A useful numerical description of a distribution requires both a measure of center and a measure of spread.

How to Calculate the Quartiles and the Interquartile Range

To calculate the **quartiles**:

- Arrange the observations in increasing order and locate the **median M** .
- The **first quartile Q_1** is the median of the observations located to the left of the median in the ordered list.
- The **third quartile Q_3** is the median of the observations located to the right of the median in the ordered list.
- The **interquartile range (IQR)** is defined as: $IQR = Q_3 - Q_1$.

32

The Five-Number Summary

- The minimum and maximum values alone tell us little about the distribution as a whole. Likewise, the median and quartiles tell us little about the tails of a distribution.
- To get a quick summary of both center and spread, combine all five numbers.

The **five-number summary** of a distribution consists of the smallest observation, the first quartile, the median, the third quartile, and the largest observation, written in order from smallest to largest.

Minimum Q_1 M Q_3 Maximum

33

Suspected Outliers: $1.5 \times \text{IQR}$ Rule

In addition to serving as a measure of spread, the interquartile range (IQR) is used as part of a rule of thumb for identifying outliers.

The $1.5 \times \text{IQR}$ Rule for Outliers

Call an observation an outlier if it falls more than $1.5 \times \text{IQR}$ above the third quartile or below the first quartile.

In the New York travel time data, $Q_1 = 15$ minutes, $Q_3 = 42.5$ minutes, and so $\text{IQR} = 27.5$ minutes.

For these data, $1.5 \times \text{IQR} = 1.5(27.5) = 41.25$

$$Q_1 - 1.5 \times \text{IQR} = 15 - 41.25 = -26.25$$

$$Q_3 + 1.5 \times \text{IQR} = 42.5 + 41.25 = 83.75$$

Any travel time shorter than -26.25 minutes (of course no such thing) or longer than 83.75 minutes is considered an outlier.

0	5
1	005555
2	0005
3	00
4	005
5	
6	005
7	
8	5

34

Boxplots

The median and quartiles divide the distribution roughly into quarters. This leads to a new way to display quantitative data, the **boxplot**.

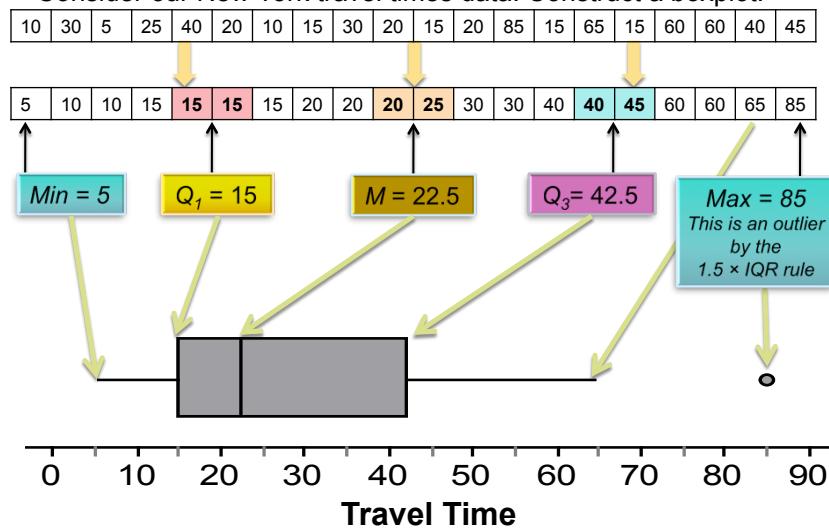
How to Make a Boxplot

- Draw and label a number line that includes the range of the distribution.
- Draw a central box from Q_1 to Q_3 .
- Note the median M inside the box.
- Extend lines (whiskers) from the box out to the minimum and maximum values that are not outliers.

35

Boxplots

- Consider our New York travel times data. Construct a boxplot.



37

Measuring Spread: Standard Deviation

The most common measure of spread looks at how far each observation is from the mean. This measure is called the **standard deviation**.

The **standard deviation** s_x measures the average distance of the observations from their mean. It is calculated by finding an average of the squared distances and then taking the square root. This average squared distance is called the **variance**.

$$\text{variance} = s_x^2 = \frac{(x_1 - \bar{x})^2 + (x_2 - \bar{x})^2 + \dots + (x_n - \bar{x})^2}{n-1} = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

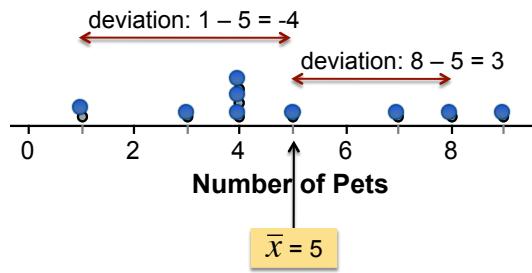
$$\text{standard deviation} = s_x = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2}$$

38

Calculating the Standard Deviation

Example: Consider the following data on the number of pets owned by a group of nine children.

1. Calculate the mean.
2. Calculate each *deviation*.
 $\text{deviation} = \text{observation} - \text{mean}$



39

Calculating the Standard Deviation

3. Square each deviation.
4. Find the “average” squared deviation. Calculate the sum of the squared deviations divided by $(n - 1)$. This is called the **variance**.
5. Calculate the square root of the variance. This is the **standard deviation**.

x_i	$(x_i\text{-mean})$	$(x_i\text{-mean})^2$
1	$1 - 5 = -4$	$(-4)^2 = 16$
3	$3 - 5 = -2$	$(-2)^2 = 4$
4	$4 - 5 = -1$	$(-1)^2 = 1$
4	$4 - 5 = -1$	$(-1)^2 = 1$
4	$4 - 5 = -1$	$(-1)^2 = 1$
5	$5 - 5 = 0$	$(0)^2 = 0$
7	$7 - 5 = 2$	$(2)^2 = 4$
8	$8 - 5 = 3$	$(3)^2 = 9$
9	$9 - 5 = 4$	$(4)^2 = 16$
Sum = ?		Sum = ?

“Average” squared deviation = $52/(9 - 1) = 6.5$. This is the **variance**.

Standard deviation = square root of variance = $\sqrt{6.5} = 2.55$

40

Properties of the Standard Deviation

- s measures spread about the mean and should be used only when the mean is the measure of center.
- $s = 0$ only when all observations have the same value and there is no spread. Otherwise, $s > 0$.
- s is not resistant to outliers.
- Variance and standard deviation both measure spread and are square/square root of one another. Standard deviation (s) has the same units of measurement as the original observations. Variance has different units (squared units).

41

Choosing Measures

We now have a choice between two descriptions for center and spread:

- ✓ Mean and standard deviation
- ✓ Median and interquartile range

Choosing Measures of Center and Spread

The median and *IQR* are usually better than the mean and standard deviation for describing a skewed distribution or a distribution with outliers.

Use mean and standard deviation only for reasonably symmetric distributions that don't have outliers.

NOTE: Numerical summaries do not fully describe the shape of a distribution. *ALWAYS PLOT YOUR DATA!*

42

Changing the Unit of Measurement

Variables can be recorded in different units of measurement. Most often (at least we wish), one measurement unit is a **linear transformation** of another measurement unit: $x_{\text{new}} = a + bx$.

Linear transformations do not change the basic shape of a distribution (skew, symmetry, multimodal). But they do change the measures of center and spread:

- Multiplying each observation by a positive number b multiplies both measures of center (mean, median) and spread (IQR, s) by b .
- Adding the same number a (positive or negative) to each observation adds a to measures of center and to quartiles, but it does not change measures of spread (IQR, s).

43

1.4 Density Curves and Normal Distributions



- Density curves
- Measuring center and spread for density curves
- Normal distributions
- The 68-95-99.7 rule
- Standardizing observations
- Using the standard Normal Table
- Inverse Normal calculations
- Normal quantile plots

44

Exploring Quantitative Data

We now have a kit of graphical and numerical tools for describing distributions. We also have a strategy for exploring data on a single quantitative variable. Now we'll add a fourth step to the strategy.

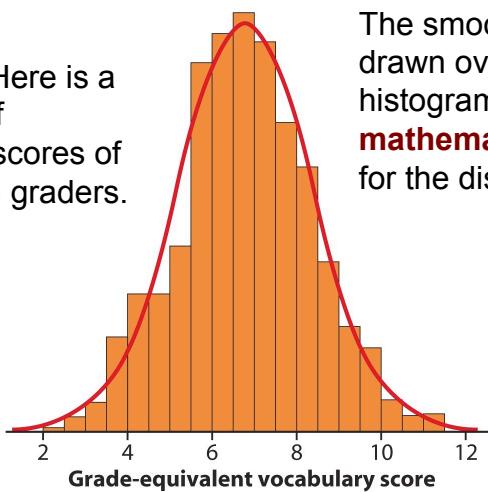
Exploring Quantitative Data

1. Always plot your data: make a graph.
2. Look for the overall pattern (shape, center, and spread) and for striking departures such as outliers.
3. Calculate a numerical summary to briefly describe center and spread.
4. Sometimes, the overall pattern of a large number of observations is so regular that we can describe it by a smooth curve.

45

Density Curves

Example: Here is a histogram of vocabulary scores of 947 seventh graders.

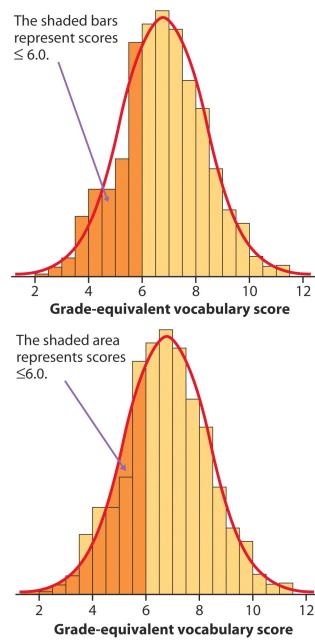


The smooth curve drawn over the histogram is a **mathematical model** for the distribution.

46

Density Curves

- The areas of the shaded bars in this histogram represent the proportion of scores in the observed data that are less than or equal to 6.0. This proportion is equal to 0.303.
- Now the area under the smooth curve to the left of 6.0 is shaded. If the scale is adjusted so the total area under the curve is exactly 1, then this curve is called a **density curve**. The proportion of the area to the left of 6.0 is now equal to 0.293.



47

Density Curves

- A **density curve** is a curve that:
 - is always on or above the horizontal axis
 - has an area of exactly 1 underneath it

- A density curve describes the overall pattern of a distribution. The area under the curve and above any range of values on the horizontal axis is the proportion of all observations that fall in that range.

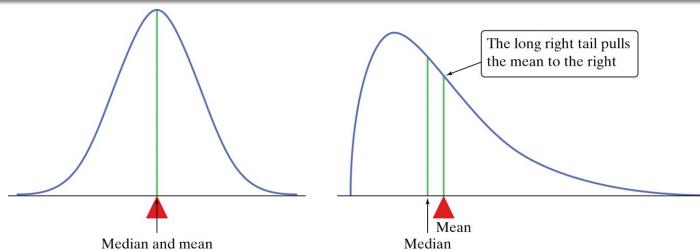
48

Density Curves

Our measures of center and spread apply to density curves as well as to actual sets of observations.

Distinguishing the Median and Mean of a Density Curve

- The **median** of a density curve is the equal-areas point—the point that divides the area under the curve in half.
- The **mean** of a density curve is the balance point, that is, the point at which the curve would balance if made of solid material.
- The median and the mean are the same for a symmetric density curve. They both lie at the center of the curve. The mean of a skewed curve is pulled away from the median in the direction of the long tail.



49

Density Curves

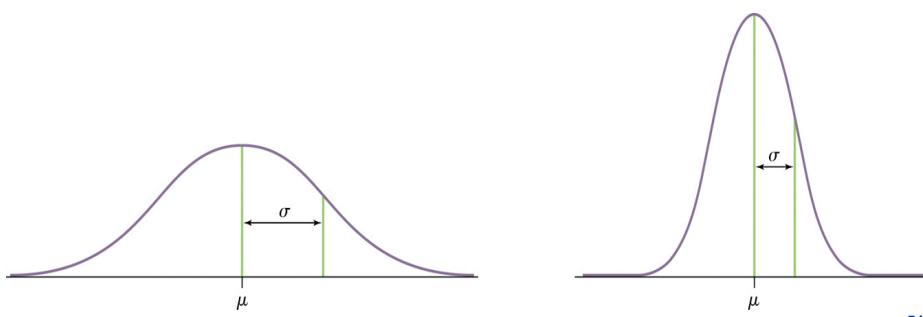
- The mean and standard deviation computed from actual observations (data) are denoted by \bar{x} and s , respectively.
- The mean and standard deviation of the actual distribution represented by the density curve are denoted by μ (“mu”) and σ (“sigma”), respectively.

50

Normal Distributions

One particularly important class of density curves is the class of Normal curves, which describe Normal distributions.

- All Normal curves are symmetric, single-peaked, and bell-shaped.
- A specific Normal curve is described by giving its mean μ and standard deviation σ .



51

Normal Distributions

A **Normal distribution** is described by a Normal density curve. Any particular Normal distribution is completely specified by two numbers: its mean μ and standard deviation σ .

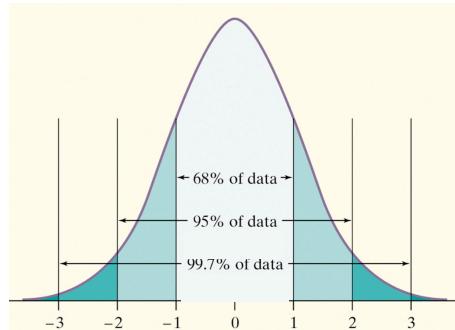
- The mean of a Normal distribution is the center of the symmetric **Normal curve**.
- The standard deviation is the distance from the center to the change-of-curvature points on either side.
- We abbreviate the Normal distribution with mean μ and standard deviation σ as $N(\mu, \sigma)$.

52

The 68-95-99.7 Rule

In the Normal distribution with mean μ and standard deviation σ :

- Approximately **68%** of the observations fall within σ of μ .
- Approximately **95%** of the observations fall within 2σ of μ .
- Approximately **99.7%** of the observations fall within 3σ of μ .

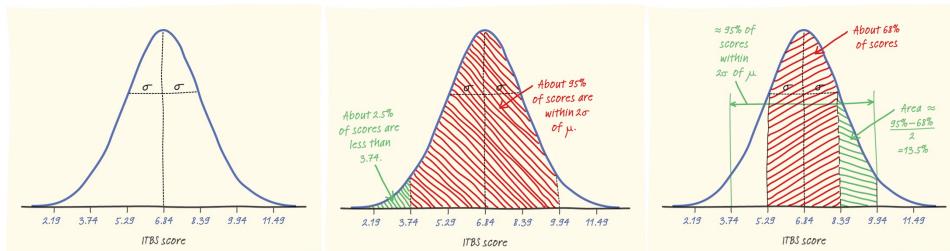


53

The 68-95-99.7 Rule

The distribution of Iowa Test of Basic Skills (ITBS) vocabulary scores for 7th-grade students in Gary, Indiana, is close to Normal. Suppose the distribution is $N(6.84, 1.55)$.

- ✓ Sketch the Normal density curve for this distribution.
- ✓ What percent of ITBS vocabulary scores are less than 3.74?
- ✓ What percent of the scores are between 5.29 and 9.94?



54

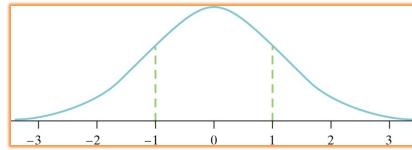
Standardizing Observations

If a variable x has a distribution with mean μ and standard deviation σ , then the **standardized value** of x , or its ***z-score***, is

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

All Normal distributions are the same if we measure in units of size σ from the mean μ as center.

The **standard Normal distribution** is the Normal distribution with mean 0 and standard deviation 1. That is, the standard Normal distribution is $N(0, 1)$.



55

The Standard Normal Table



TABLE A Standard normal probabilities										
<i>z</i>	.00	.01	.02	.03	.04	.05	.06	.07	.08	.09
-3.4	.0003	.0003	.0003	.0003	.0003	.0003	.0003	.0003	.0002	.0002
-3.3	.0005	.0005	.0005	.0004	.0004	.0004	.0004	.0004	.0003	.0002
-3.2	.0008	.0008	.0008	.0007	.0007	.0007	.0007	.0006	.0005	.0004
-3.1	.0010	.0009	.0009	.0008	.0008	.0008	.0008	.0007	.0006	.0005
-3.0	.0013	.0013	.0013	.0012	.0012	.0011	.0011	.0010	.0010	.0009
-2.9	.0019	.0018	.0018	.0017	.0016	.0016	.0015	.0015	.0014	.0014
-2.8	.0026	.0026	.0026	.0023	.0023	.0023	.0023	.0022	.0022	.0021
-2.7	.0035	.0034	.0033	.0032	.0031	.0030	.0029	.0028	.0027	.0026
-2.6	.0047	.0045	.0044	.0043	.0041	.0040	.0039	.0038	.0037	.0036
-2.5	.0062	.0060	.0059	.0057	.0056	.0055	.0053	.0052	.0051	.0050
-2.4	.0082	.0080	.0078	.0075	.0073	.0071	.0069	.0068	.0066	.0064
-2.3	.0107	.0104	.0102	.0099	.0096	.0094	.0091	.0089	.0087	.0084
-2.2	.0139	.0136	.0132	.0129	.0125	.0122	.0119	.0116	.0113	.0110
-2.1	.0179	.0176	.0172	.0168	.0164	.0161	.0158	.0154	.0150	.0146
-2.0	.0228	.0222	.0217	.0212	.0207	.0202	.0197	.0192	.0188	.0183
-1.9	.0287	.0281	.0274	.0268	.0262	.0256	.0250	.0244	.0239	.0233
-1.8	.0359	.0351	.0344	.0336	.0329	.0322	.0314	.0307	.0301	.0294
-1.7	.0438	.0430	.0424	.0418	.0410	.0403	.0396	.0387	.0375	.0365
-1.6	.0528	.0537	.0526	.0516	.0505	.0495	.0485	.0475	.0465	.0455
-1.5	.0628	.0655	.0643	.0630	.0618	.0606	.0594	.0582	.0571	.0559
-1.4	.0740	.0767	.0794	.0815	.0835	.0853	.0870	.0889	.0901	.0913
-1.3	.0868	.0905	.0934	.0918	.0901	.0885	.0869	.0853	.0838	.0823
-1.2	.1151	.1131	.1112	.1093	.1075	.1056	.1038	.1020	.1003	.0985
-1.1	.1357	.1335	.1314	.1292	.1271	.1251	.1230	.1210	.1190	.1170
-1.0	.1573	.1570	.1565	.1555	.1542	.1530	.1518	.1504	.1491	.1479
-0.9	.1841	.1814	.1788	.1762	.1736	.1711	.1685	.1660	.1635	.1611
-0.8	.2119	.2090	.2061	.2033	.2005	.1977	.1949	.1922	.1894	.1867
-0.7	.2402	.2380	.2357	.2327	.2296	.2266	.2236	.2206	.2177	.2148
-0.6	.2743	.2709	.2676	.2643	.2610	.2576	.2541	.2507	.2473	.2438
-0.5	.3085	.3050	.3015	.2981	.2946	.2912	.2877	.2843	.2810	.2776
-0.4	.3446	.3409	.3372	.3336	.3300	.3264	.3228	.3192	.3156	.3121
-0.3	.3821	.3787	.3753	.3717	.3681	.3645	.3609	.3573	.3537	.3501
-0.2	.4207	.4168	.4129	.4090	.4052	.4013	.3974	.3936	.3897	.3859
-0.1	.4602	.4562	.4522	.4483	.4443	.4404	.4364	.4325	.4286	.4247
0.0	.5000	.4960	.4920	.4880	.4840	.4801	.4761	.4721	.4681	.4641

TABLE A Standard normal probabilities (continued)										
<i>z</i>	.00	.01	.02	.03	.04	.05	.06	.07	.08	.09
0.0	.5000	.5040	.5080	.5120	.5160	.5199	.5239	.5279	.5319	.5359
0.1	.5398	.5438	.5478	.5517	.5557	.5596	.5636	.5675	.5714	.5753
0.2	.5797	.5837	.5876	.5915	.5954	.5993	.6032	.6071	.6110	.6149
0.3	.6179	.6217	.6255	.6293	.6331	.6368	.6406	.6443	.6480	.6517
0.4	.6554	.6591	.6628	.6664	.6700	.6736	.6772	.6808	.6844	.6879
0.5	.6915	.6950	.6985	.7019	.7054	.7088	.7123	.7157	.7190	.7224
0.6	.7257	.7291	.7324	.7357	.7382	.7412	.7442	.7471	.7501	.7549
0.7	.7580	.7611	.7642	.7673	.7704	.7734	.7764	.7794	.7823	.7852
0.8	.7881	.7910	.7939	.7967	.7994	.8023	.8051	.8079	.8106	.8133
0.9	.8159	.8186	.8216	.8242	.8269	.8293	.8319	.8343	.8367	.8389
1.0	.8413	.8448	.8481	.8515	.8544	.8577	.8609	.8641	.8673	.8701
1.1	.8643	.8665	.8688	.8708	.8729	.8749	.8770	.8790	.8810	.8830
1.2	.8849	.8869	.8888	.8907	.8925	.8944	.8962	.8980	.8997	.9015
1.3	.9032	.9050	.9068	.9085	.9103	.9121	.9139	.9157	.9174	.9192
1.4	.9192	.9207	.9222	.9236	.9251	.9265	.9279	.9292	.9306	.9319
1.5	.9332	.9345	.9357	.9370	.9382	.9396	.9418	.9439	.9459	.9471
1.6	.9452	.9463	.9474	.9484	.9495	.9505	.9515	.9525	.9535	.9545
1.7	.9554	.9564	.9573	.9583	.9592	.9600	.9609	.9618	.9627	.9633
1.8	.9641	.9649	.9656	.9664	.9671	.9678	.9685	.9693	.9700	.9706
1.9	.9713	.9719	.9726	.9732	.9738	.9744	.9750	.9756	.9761	.9767
2.0	.9772	.9778	.9784	.9790	.9796	.9802	.9808	.9814	.9819	.9824
2.1	.9821	.9826	.9830	.9834	.9838	.9842	.9846	.9850	.9854	.9857
2.2	.9861	.9864	.9868	.9871	.9875	.9878	.9881	.9884	.9887	.9890
2.3	.9893	.9896	.9898	.9901	.9904	.9906	.9909	.9911	.9913	.9916
2.4	.9919	.9923	.9926	.9929	.9932	.9935	.9937	.9940	.9942	.9946
2.5	.9938	.9940	.9941	.9943	.9945	.9946	.9948	.9949	.9951	.9952
2.6	.9953	.9955	.9956	.9957	.9959	.9960	.9961	.9962	.9963	.9964
2.7	.9965	.9966	.9967	.9968	.9969	.9970	.9971	.9972	.9973	.9974
2.8	.9974	.9975	.9976	.9977	.9978	.9979	.9980	.9981	.9982	.9983
2.9	.9981	.9982	.9982	.9983	.9984	.9984	.9985	.9985	.9986	.9986
3.0	.9987	.9987	.9987	.9988	.9988	.9989	.9989	.9989	.9990	.9990
3.1	.9990	.9991	.9991	.9992	.9992	.9993	.9993	.9994	.9994	.9994
3.2	.9993	.9994	.9994	.9994	.9994	.9994	.9994	.9995	.9995	.9995
3.3	.9995	.9995	.9995	.9996	.9996	.9996	.9996	.9996	.9996	.9997
3.4	.9997	.9997	.9997	.9997	.9997	.9997	.9997	.9997	.9997	.9997

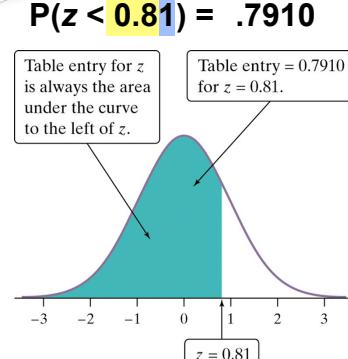
56

The Standard Normal Table

- Suppose we want to find the proportion of observations from the standard Normal distribution that are less than 0.81.

- We can use Table A:

$$P(z < 0.81) = .7910$$

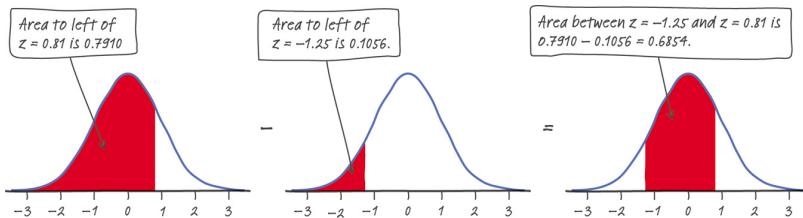


57

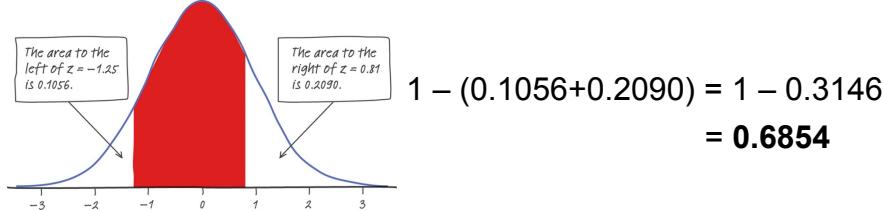
Z	.00	.01	.02
0.7	.7580	.7611	.7642
0.8	.7881	.7910	.7939
0.9	.8159	.8186	.8212

Normal Calculations

Find the proportion of observations from the standard Normal distribution that are between -1.25 and 0.81 .



Can you find the same proportion using a different approach?



58

Normal Calculations

How to Solve Problems Involving Normal Distributions

Express the problem in terms of the observed variable x .

Draw a picture of the distribution and shade the area of interest under the curve.

Perform calculations.

- **Standardize** x to restate the problem in terms of a standard Normal variable z .
- **Use Table A** and the fact that the total area under the curve is 1 to find the required area under the standard Normal curve.

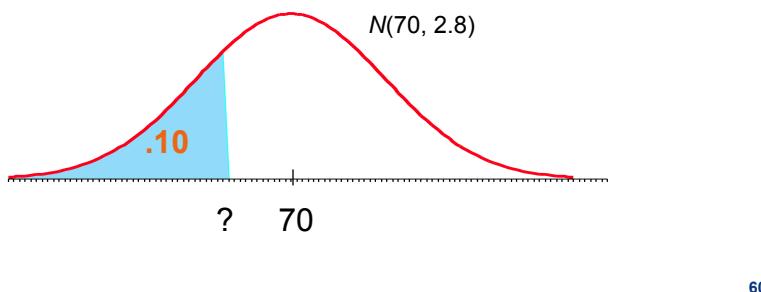
Write your conclusion in the context of the problem.

59

Normal Calculations

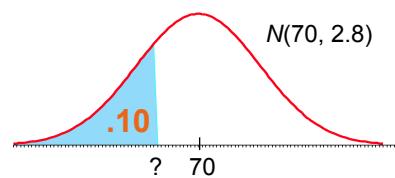
According to the Health and Nutrition Examination Study of 1976–1980, the heights (in inches) of adult men aged 18–24 are $N(70, 2.8)$.

If exactly 10% of men aged 18–24 are shorter than John Doe, how tall is he?



Normal Calculations

How tall is a man who is taller than exactly 10% of men aged 18–24?



Look up the probability closest to 0.10 in the table.

Find the corresponding **standardized score**.

The value you seek is that many standard deviations from the mean.

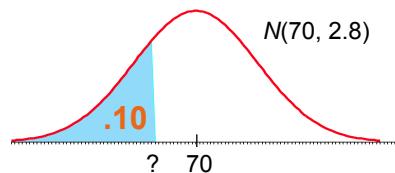
z	.07	.08	.09
-1.3	.0853	.0853	.0823
-1.2	.0985	.1003	.1003
-1.1	.1170	.1190	.1190

$$Z = -1.28$$

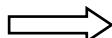
Normal calculations

How tall is a man who is taller than exactly 10% of men aged 18–24?

$$Z = -1.28$$



We need to “unstandardize” the z-score to find the observed value (x):



$$\begin{aligned} x &= 70 + z(2.8) \\ &= 70 + [(-1.28) \times (2.8)] \\ &= 70 + (-3.58) = \underline{\underline{66.42}} \end{aligned}$$

A man would have to be approximately 66.42 inches tall or less to be in the lower 10% of all men in the population.

62

Normal Quantile Plots

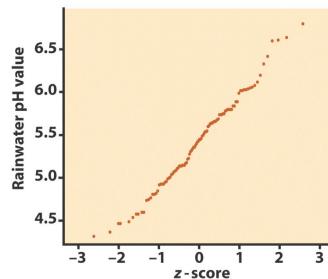
One way to assess if a distribution is indeed approximately Normal is to plot the data on a **Normal quantile plot**.

The data points are ranked and the percentile ranks are converted to z-scores with Table A. The z-scores are then used for the x-axis against which the data are plotted on the y-axis of the Normal quantile plot.

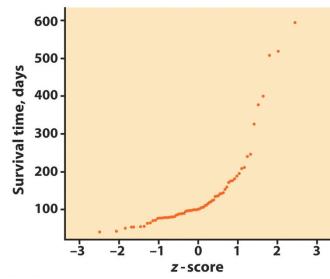
- If the distribution is indeed Normal, the plot will show a straight line, indicating a good match between the data and a Normal distribution.
- Systematic deviations from a straight line indicate a non-Normal distribution. Outliers appear as points that are far away from the overall pattern of the plot.

63

Normal Quantile Plots



Good fit to a straight line: The distribution of rainwater pH values is close to Normal.



Curved pattern: The data are not Normally distributed. Instead, the data are right skewed: A few individuals have particularly long survival times.

Normal quantile plots are complex to do by hand, but they are standard features in most statistical software.

64

Using R: examples you can try

- `barplot(table(state.region))`
- `barplot(prop.table(table(state.region)))`
- `stem(trees$Girth)`
- `hist(precip)`
- `hist(precip, breaks=500)`
- `plot(LakeHuron)`
- `plot(LakeHuron, type = "h")`
- `plot(LakeHuron, type = "p")`
- `boxplot(precip)`
- `summary(precip)`
- `var(precip)`
- `sd(precip)`
- `pnorm(q, mean=0, sd=1)`
- `qnorm(p, mean=0, sd=1)`
- `qqnorm(trees$Height)`

65