

Data and Data Sets

 $\succ \underline{\bf Data}$ are the facts and figures collected, analyzed, and summarized for presentation and interpretation.

>All the data collected in a particular study are referred to as the $\underline{\text{data}}$ $\underline{\text{set}}$ for the study.

Elements, Variables, and Observations

▶ Elements are the entities on which data are collected.

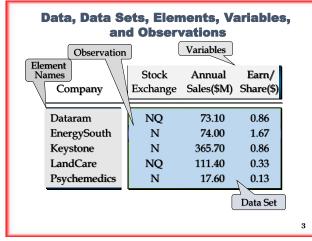
A variable is a characteristic of interest for the elements

>The set of measurements obtained for a particular element is called an observation.

A data set with n elements contains n observations.

>The total number of data values in a complete data set is the is the number of elements multiplied by the number of variables.

2



Scales of Measurement

The scale determines the amount of information contained in the data and the statistical analyses that are most appropriate.

■Nominal

Data are <u>labels or names</u> used to identify an attribute of the element.

A nonnumeric label or numeric code may be used.

Example:

2

Students of a university are classified by the school in which they are enrolled using a nonnumeric label such as Business, Humanities, Education, and so on.

Alternatively, a numeric code could be used for the school variable (e.g. 1 denotes Business, 2 denotes Humanities, 3 denotes Education, and so on).

4

3

Scales of Measurement Ordinal The data have the properties of nominal data and the order or rank of the data is meaningful. A nonnumeric label or numeric code may be used. Example: Students of a university are classified by their class standing using a nonnumeric label such as Freshman, Sophomore, Junior, or Senior. Alternatively, a numeric code could be used for the class standing variable (e.g. 1 denotes Freshman, 2 denotes Sophomore, and so on).

Scales of Measurement

Interval

The data have the properties of ordinal data, and the interval between observations is expressed in terms of a fixed unit of measure.

Interval data are always numeric.

Example:

Melissa has an SAT score of 1885, while Kevin has an SAT score of 1780. Melissa scored 105 points more than Kevin.

5

Data can be further classified as being categorical or quantitative.

Categorical Data

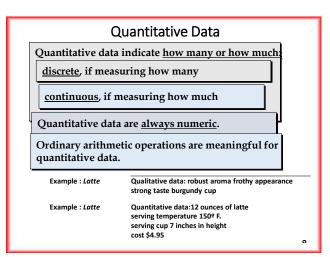
Labels or names used to identify an attribute of each element

Often referred to as qualitative data

Use either the nominal or ordinal scale of measurement

Can be either numeric or nonnumeric

Appropriate statistical analyses are rather limited



7

Summarizing Data

- Frequency Distribution
- Relative Frequency and Percent Frequency Distributions
- **□**Bar Chart
- ☐ Pie Chart

9

10

Categorical data use labels or names to identify categories of like items.

Quantitative data are numerical values that indicate how much or how many.

Frequency Distribution

A frequency distribution is a tabular summary of data showing the frequency (or number) of items in each of several non-overlapping classes.

The objective is to provide insights about the data that cannot be quickly obtained by looking only at the original data.

Frequency Distribution

Example: Marada Inn

Guests staying at Marada Inn were asked to rate the quality of their accommodations as being excellent, above average, average, below average, or poor. The ratings provided by a sample of 20 guests are:

Below Average	Average	Above Average	
Above Average	Above Average	Above Average	
Above Average	Below Average	Below Average	
Average	Poor	Poor	
Above Average	Excellent	Above Average	
Average	Above Average	Average	
Above Average	Average		

■ Excel Formula Worksheet

Note: Rows 9-21 are not shown.

A B C

1 Quality Rating Quality Rating Frequence

	adduncy reading		addition reading	
2	Above Average		Poor	=COUNTIF(\$A\$2:\$A\$21,C2)
3	Below Average		Below Average	=COUNTIF(\$A\$2:\$A\$21,C3)
4	Above Average		Average	=COUNTIF(\$A\$2:\$A\$21,C4)
5	Average		Above Average	=COUNTIF(\$A\$2:\$A\$21,C5)
6	Average		Excellent	=COUNTIF(\$A\$2:\$A\$21,C6)
7	Above Average		Total	=SUM(D2:D6)
8	Above Average			
	Α	В	С	D
1	Quality Rating		Quality Rating	Frequency
2	Above Average		Poor	2
3	Below Average		Below Average	3
4	Above Average		Average	5
5	Average		Above Average	9
6	Average		Excellent	1
7	rtvorago			
- /	Above Average		Total	20

12

11 12

Frequency Distribution Example: Marada Inn Rating Frequency 2 Poor 3 Below Average 5 Average 9 Above Average Excellent 1 Total 20 13

Relative Frequency Distribution

The <u>relative frequency</u> of a class is the fraction or proportion of the total number of data items belonging to the class.

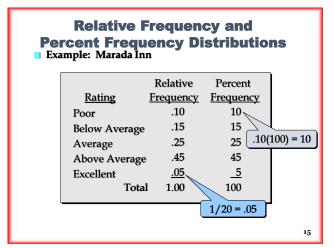
A <u>relative frequency distribution</u> is a tabular summary of a set of data showing the relative frequency for each class.

Percent Frequency Distribution

The <u>percent frequency</u> of a class is the relative frequency multiplied by 100.

14

13 14



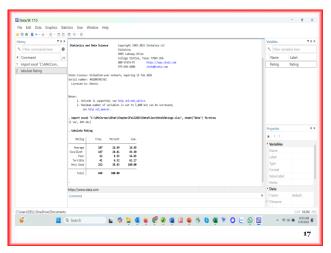
But what if the sample size is much larger?

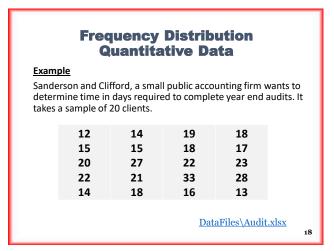
Chapter2DataFiles\HotelRatings.xlsx

Looking at the Soft Drink Market

 $\underline{Chapter2DataFiles \backslash SoftDrink.xlsx}$

16





17 18

Frequency Distribution

Example: Sanderson and Clifford

- If we choose five classes:
- Approximate Class Width = (33 12)/5 = 4.2 ≅ (4)

Time in days	Frequency		
10-14	4		
15-19	8		
20-24	5		
25-29	2		
30-34	1		
Total	20		

19

Using Excel's PivotTable to Construct a Pivot Table

> Step 1 – Select any cell in the data set

➤ Step 2 – Click Insert on the Ribbon

> Step 3 – In the Tables group click Pivot Table

> Step 4 - When the Create PivotTable dialog box appears: Click

(a PivotTable and PivotTable Fields dialog box will appear in a new worksheet)

20

19 20

Using Excel's *PivotTable* to Construct a Pivot Table

Step 5 - In the PivotTable Fields dialog box: Drag Audit time to the Rows area Drag Audit time to the Values area

- > Step 6 Click on Sum of Audit time in the Values area
- > Step 7 Click Value Field Settings from the list of options
- > Step 8 When the Value Field Settings dialog box appears: Under Summarize value field by, choose Count Click OK

21

Using Excel's Pivot table to construct a frequency distribution 3 Row Lables Count of audit time PivotTable Fields 5 15-19 Choose fields to add to report: 45 -6 20-24 7 25-29 8 30-34 9 Grand total Drag fields between areas below: 11 12 13 14 **▼** Filters III Columns 15 16

21 22

Using Excel's PivotTable to Construct a Frequency Distribution

To construct the frequency distribution, we must group the rows containing audit time.

- Step 1 Right click any cell in the PivotTable report containing a an audit time.
- >Step 2 Choose Group from the list of options that appears
- > Step 3 When the Grouping dialog box appears:
 - ✓ Enter 10 in the **Starting at** box ✓ Enter 34 in the **Ending at** box
 - ✓ Enter 5 in the **By** box
 - ✓ Click **OK**

23

Relative Frequency and Percent Frequency Distributions

Example: Sanderson and Clifford

Audit time (in days)	Relative Frequency	Percent Frequency	
10 – 14	.20 (4/20)	20 (0.2 * 100)	
15 – 19	.40	40	
20 – 25	.25	25	
25 – 29	.10	10	
30 – 34	.05	5	
	Total 1.00	100	

24

23 24

Relative Frequency and Percent Frequency Distributions

Example: Sanderson and Clifford

Insights obtained from the Percent Frequency Distribution:

- ✓ 40% of the audits required from 15 to 19 days.
- ✓ Another 25% of the audits required 20 to 25 days.
- ✓ Only 5% of the audits required more than 30 days.

25

Example— Frequency Distribution

Work Status in the General Social Survey

GSS asked: "Last week were you working full time, part time, going to school, keeping house, or what?"

The responses were as follows:

1. Working full time
2. Working part time
3. Temporally not working
4. Unemployed, laid off
5. Relited
6. School
7. Keeping house
8. Other
DATA: GSS2018

Variable: WRKSTAT (Column X)

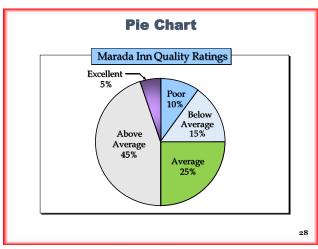
Variable: WRKSTAT (Column X)

Chapter 2 Data Files \ GSS2018.xlsx

25 26

Pie Chart

- The <u>pie chart</u> is a commonly used graphical device for presenting relative frequency and percent frequency distributions for categorical data.
- First draw a <u>circle</u>; then use the relative frequencies to subdivide the circle into sectors that correspond to the relative frequency for each class.
- Since there are 360 degrees in a circle, a class with a relative frequency of .25 would consume .25(360) = 90 degrees of the circle.



27 28

7

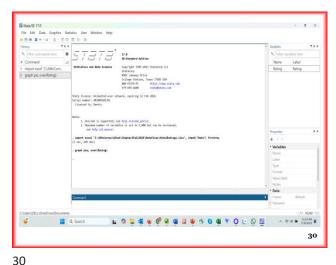
But what if the sample size is much larger?

 $\underline{DataFiles}\backslash \underline{HotelRatings.xlsx}$

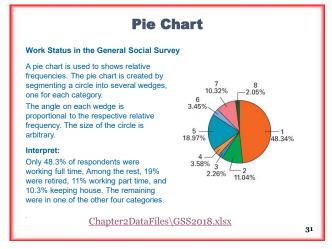
Looking at the Soft Drink Market

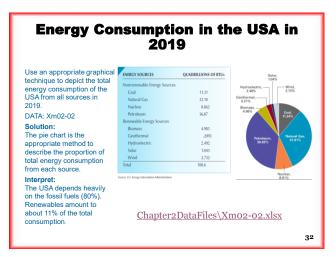
DataFiles\SoftDrink.xlsx

29



29





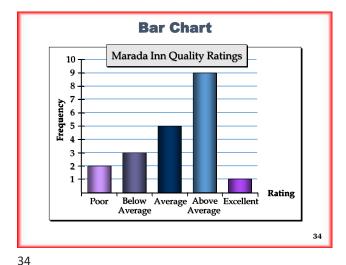
Bar Chart

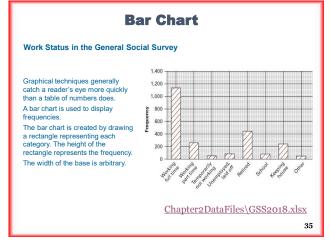
- A <u>bar chart</u> is a graphical device for depicting qualitative data.
- On one axis (usually the horizontal axis), we specify the labels that are used for each of the classes.
- A frequency, relative frequency, or percent frequency scale can be used for the other axis (usually the vertical axis).
- Using a <u>bar of fixed width</u> drawn above each class label, we extend the height appropriately.

33

The <u>bars are separated</u> to emphasize the fact that each class is a separate category.

33





Frequency Distribution and Bar Chart

Example 2.1 – Stata Output and Instructions
Chapter 2Stata instructions.pdf

| WRKSTAT | Freqs. | Percent | Cum. | 1 | 1,134 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 | 48,34 |

35 36

Cumulative Distributions

Cumulative frequency distribution - shows the number of items with values less than or equal to the upper limit of each class...

<u>Cumulative relative frequency distribution</u> – shows the *proportion* of items with values less than or equal to the upper limit of each class.

Cumulative percent frequency distribution - shows the percentage of items with values less than or equal to the upper limit of each class.

37

Cumulative Distributions

Example: Sanderson and Cliffords

Audit time (Days)	Cumulative Frequency	Cumulative Relative Frequency	Cumulative Percent Frequency
≤ 14	4	.20	20
≤ 19	12	.60	60
≤ 24	17	.85	85
≤ 29	19	.95	95
≤ 34	20	1.00	100

37 38

Crosstabulation

- A crosstabulation is a tabular summary of data for two variables.
- Crosstabulation can be used when:
 - one variable is qualitative and the other is quantitative,
 - both variables are qualitative, or
 - both variables are quantitative.
- The left and top margin labels define the classes for the two variables.

Crosstabulation

Example: Finger Lakes Homes

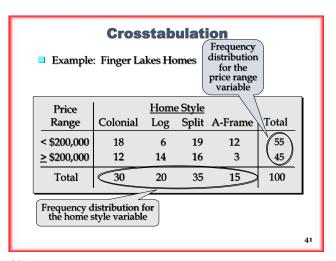
The number of Finger Lakes homes sold for each style and price for the past two years is shown below. quantitative

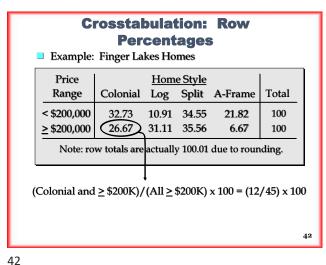
	variable			varia	able
Price	Home Style				1
Range	Colonial	Log	Split	A-Frame	Total
< \$200,000	18	6	19	12	55
<u>></u> \$200,000	12	14	16	3	45
Total	30	20	35	15	100

catogorical

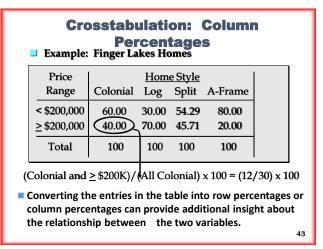
39 40

10



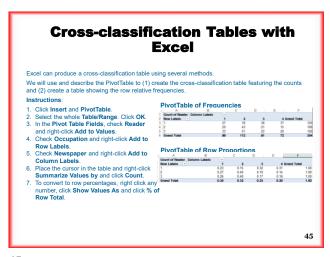


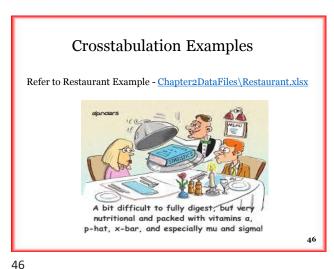
41





43 44





45

