

GE2262 Business Statistics

Topic 1

Introduction to Statistics

Reference

Levine, D.M., Krehbiel, T.C. and Berenson, M.L., *Business Statistics: A First Course*, Pearson Education Ltd, Chapter 1 & 2 & 3

Liu, K. I., To K. M., *Speaking of Statistics*, Pearson Education Ltd, Chapter 1

Outline

■ Introduction

- ❑ What is/are Statistics?
- ❑ Why Study Statistics?
- ❑ Types of Variables
- ❑ Organizing and Visualizing Data
- ❑ Use and Misuse of Statistics

■ Descriptive Statistics

- ❑ Measures of Central Tendency
- ❑ Measures of Variation
- ❑ Distribution Shape
- ❑ Use of Excel in Organizing Data
- ❑ Use of Excel in Descriptive Statistics

What is/are Statistics?

Statistics

The branch of mathematics that transforms data into useful information for decision makers.

Of what?
What types?

Descriptive Statistics

How to summarize data
with tables & charts?

Collecting, summarizing, and
describing data

What measures of
central tendency &
variation to use?
What is the shape of
the distribution?

Difference?

Inferential Statistics

Drawing conclusions and/or
making decisions concerning a
population based only on
sample data

What is the difference between
sample & population? What is the
difference between sample statistics
and population parameters?



《誇世代》首播吸164萬觀眾

無綫50周年台慶劇《誇世代》劇情搞笑，劇中演員歐陽震華、陳豪、田蕊妮等大門演技，劇情已說到歐陽震華與吳業坤已交換靈魂，而之後陸續亦都會不少藝人客串出場，相信觀眾都相當期待。而此劇首播的收視相當不俗，首播24小時跨平台總收視有25.2點，有164萬觀眾人次收看，成績理想。

How to estimate the number of people watching a TV programme? *Cont'd*

- Nielsen Media Research uses a “**People Meter**” device to record the viewing behavior of members of a sample of households (minute level viewing information)
- In USA, **5,000 homes** (13,000 individuals) are selected; in Hong Kong, **815 households** (2,300 individuals) are selected



How to estimate the number of people watching a TV programme? *Cont'd*

People Meter Data

Panel_ID	Member Number	Eff_Date	Channel Code	Start Time	End Time
1124339	2	2016-4-19	99	776	806
1124339	2	2016-4-19	1	1108	1262
1124339	2	2016-4-19	99	776	806
....

It's a **STATISTICAL STUDY!**

Start time = **776** means **12:56 pm** and End time = **806** means **1:26 pm** etc.
Channel Code = 1 means watching TVB Jade; Channel Code = 99 means watching Viu TV

Hence, on “19 April 2016”, the individual “2” of household “1124339” watched channel “99” (Viu TV) from 12:56 pm (“776”) to 1:26 pm (“806”)

Basic Steps in a Statistical Study

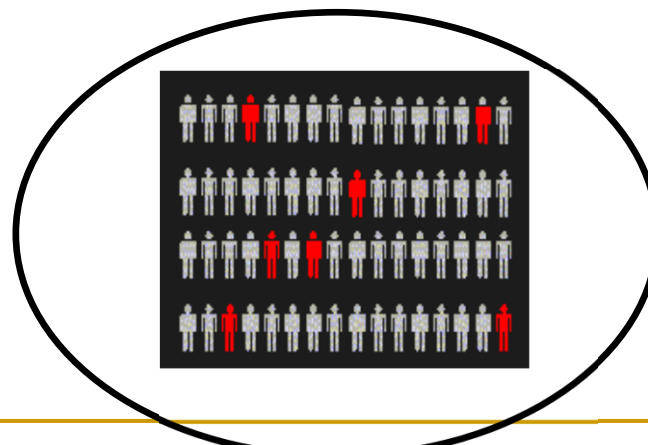
- Step 1: State the goal of your study precisely; that is, determine the **population** you want to study and exactly what you'd like to learn about it (population parameters).
- Step 2: Choose a **sample** from the population. (Be sure to use an appropriate sampling technique)
- Step 3: **Collect** raw data from the sample and **summarize** these data by finding sample statistics of interest.
- Step 4: Use the sample statistics to **make inference** about the population.
- Step 5: Draw **conclusions**; determine what you learned and whether you achieved your goal.

Population



Measures used to describe the population are called **parameters**

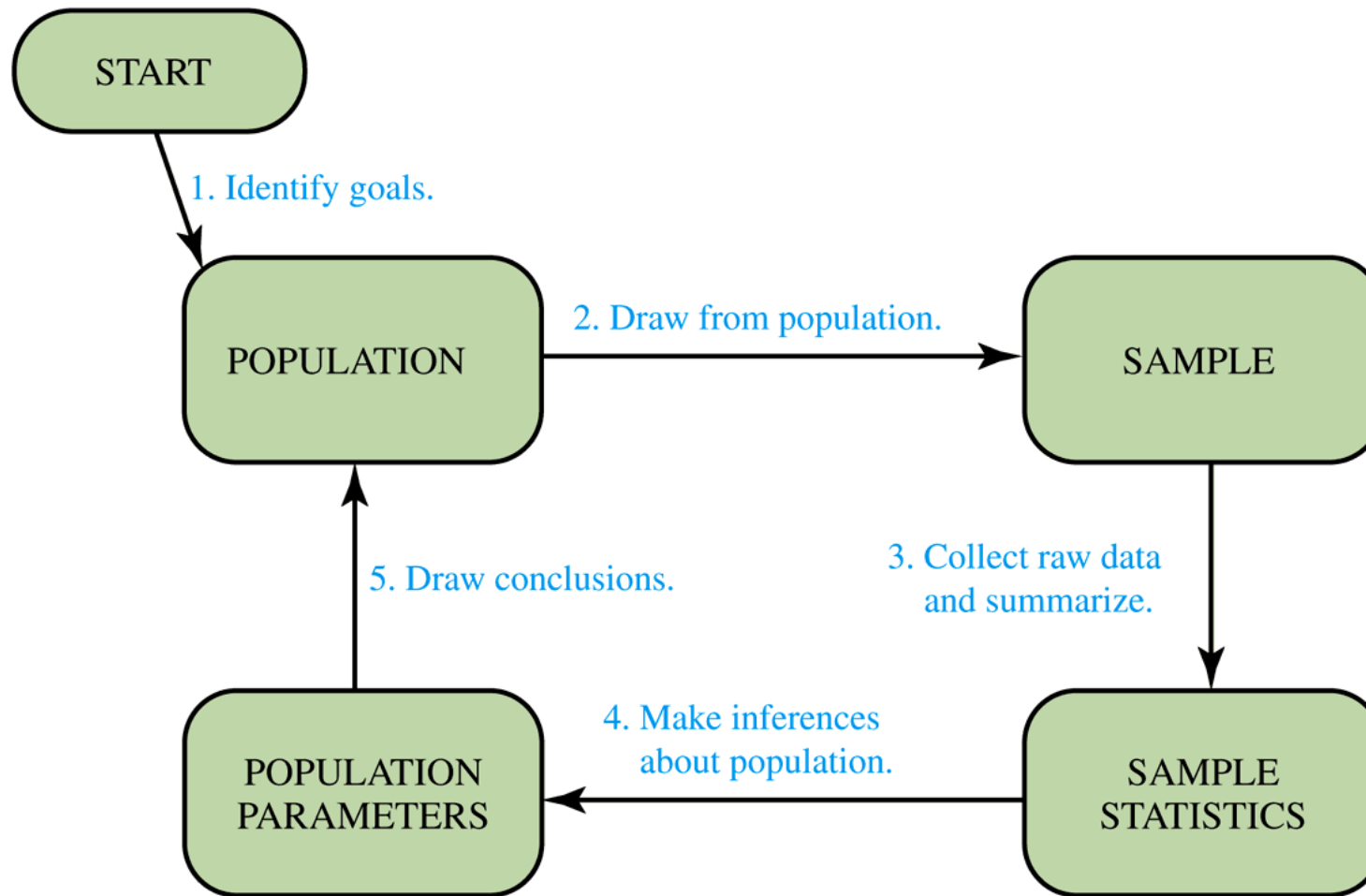
Sample



Measures computed from sample data are called **statistics**

Process of a Statistical Study

Cont'd



Why use Nielsen Media Research?

- Nielsen Media Research earns money by charging television stations and networks for its services. For example, TVB pays Nielsen to provide ratings for its television shows. Why doesn't TVB simply do its own ratings, instead of paying a company like Nielsen to do them?
- Cost of Advertising in TVB Jade

Rate Card No. 47
(Effective January 1, 2013)
Basic Spot Announcement Rates (In Hong Kong Dollars)

J7	1855-2255 (Mon-Fri)																					(HK\$)	
	RB	FB	F1	F2	F3	F3A	F4	F4A	F5	F5A	F6	F6A	F7	F7A	F8	F8A	F9	F10	F11	F12	F13		F14
30-Second	112620	177360	203940	234540	269700	289920	310140	333420	356640	383400	410160	440940	471660	507060	542400	583080	623760	670560	720840	774900	833040	895500	962640
25-Second	93850	147800	169950	195450	224750	241600	258450	277850	297200	319500	341800	367450	393050	422550	452000	485900	519800	558800	600700	645750	694200	746250	802200
20-Second	75080	118240	135960	156360	179800	193280	206760	222280	237760	255600	273440	293960	314440	338040	361600	388720	415840	447040	480560	516600	555360	597000	641760
15-Second	56310	88680	101970	117270	134850	144960	155070	166710	178320	191700	205080	220470	235830	253530	271200	291540	311880	335280	360420	387450	416520	447750	481320
10-Second	37540	59120	67980	78180	89900	96640	103380	111140	118880	127800	136720	146980	157220	169020	180800	194360	207920	223520	240280	258300	277680	298500	320880
5-Second	18770	29560	33990	39090	44950	48320	51690	55570	59440	63900	68360	73490	78610	84510	90400	97180	103960	111760	120140	129150	138840	149250	160440

Statistical Study – Exercise

Cont'd

- Describe how you would apply the five basic steps in a statistical study to estimate the average time that local CityU students use to travel from home to campus.
 - ❑ Step 1: Goal of the study: to estimate the average time that local CityU students use to travel from home to campus; Target population: all local CityU students
 - ❑ Step 2: To select a sample from local CityU students
 - ❑ Step 3: Collect the time in minutes that the student at the sample traveled from home to campus. Calculate the sample means of travelling time for the students
 - ❑ Step 4: Use statistical techniques to infer the likely results for the entire population of local CityU students
 - ❑ Step 5: Based on the likely population results, draw conclusions about the average time that local CityU students use to travel from home to campus

How to Make Money Nowadays?



Any relationship?

- Walmart put all its checkout-counter data into a giant digital warehouse and set the disk drives spinning.
- Out popped a most unexpected **correlation**: **diapers** and **beer** at the same cart usually on **Fridays**.

How to Make Money Nowadays?

Cont'd

- With statistical analysis, Walmart found that young mothers always ask fathers to purchase diapers for babies after work
- Evidently, young fathers would make a late-night run to the store to pick up **Huggies** and get some **Blue Light** while they were there
- Capitalizing on the discovery, the store
 - placed the disparate items together
 - placed high-price diapers beside beer (as males don't concern the price)
- **Sales zoomed!!!**



Why Study Statistics?

- Knowing how to do statistics can lead to become a well recognized and respected profession: Statistician
- “I don’t like numbers. Statistics are not for me!!!”
 - Unfortunately, statistics are there for everybody, like it or not!
 - “Smaller average pay rise of 3.8% likely to hit Hong Kong workers next year (2016).” The Institute of Human Resources Management
 - “The fresh graduate’s average expected monthly salary in 2014 is HK\$17,314.” APAC at Universum
 - “TVB’s ratings of weekday programmes from 8:00pm – 10:30pm losing some 10,000 viewers compared with a week ago.” SCMP, 16 April 2016
 - “Chief Executive Leung Chun-Ying’s popularity rating is 37.5 in January 2016.” The University of Hong Kong

Why Study Statistics?

Cont'd

Accounting

Information
Technology

Marketing

Economics

Finance

Why Study Statistics?

Accounting

Scope of this ISA

1. This International Standard on Auditing (ISA) applies when the auditor has decided to use audit sampling in performing audit procedures. It deals with the auditor's use of statistical and non-statistical sampling when designing and selecting the audit sample, performing tests of controls and tests of details, and evaluating the results from the sample.
2. This ISA complements ISA 500,¹ which deals with the auditor's responsibility to design and perform audit procedures to obtain sufficient appropriate audit evidence to be able to draw reasonable conclusions on which to base the auditor's opinion. ISA 500 provides guidance on the means available to the auditor for selecting items for testing, of which audit sampling is one means.

Objective

4. The objective of the auditor, when using audit sampling, is to provide a reasonable basis for the auditor to draw conclusions about the population from which the sample is selected.



■ Audit Sampling

- The application of audit procedures to **less than 100%** of items within a population of audit relevance such that all sampling units have a chance of selection in order to provide the auditor with a reasonable basis on which to **draw conclusions about the entire population**
- The auditor shall determine a **sample size** sufficient to **reduce sampling risk to an acceptably low level**

Source: <http://www.ifac.org/sites/default/files/publications/files/A028%20201>

How to obtain a
sufficient
SAMPLE SIZE?

pdf

Why Study Statistics?

Information Technology

■ IT Auditing

- Almost all computer-assisted audit tools (CAATs) have a command for Benford's Law
- Benford's Law holds true for a data set that grows **exponentially**, but also appears to hold true for many cases in which an exponential growth pattern is not obvious
- Beneficial tool for **fraud detection**, e.g. credit card transactions, purchase orders, loan data, customer refunds, ...

□ Example

- If a bank's policy is to refer loans at or above US \$50,000 to a loan committee, looking just below that approval threshold gives a loan officer the potential to discover loan frauds.
- Note that 4 is aberrantly high in occurrence, and 5 is too low, indicating the possible manipulation of the natural occurrence of loans beginning with 5 (US \$50,000 loans) possibly being switched to just under the cutoff or indicating that the suspect could be issuing a lot of \$49,999.99 loans fictitiously to embezzle funds.

How to obtain a
**PROBABILITY
DISTRIBUTION?**

Figure 1—Benford's Law Distribution Leading Digit

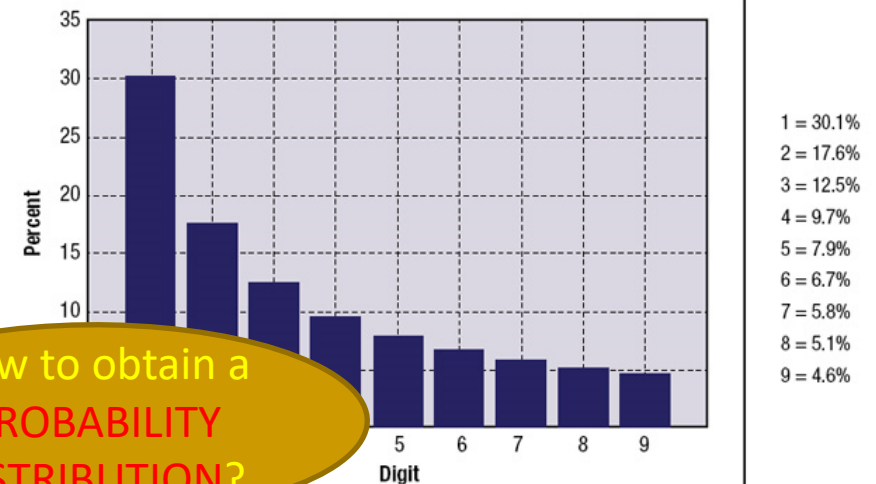
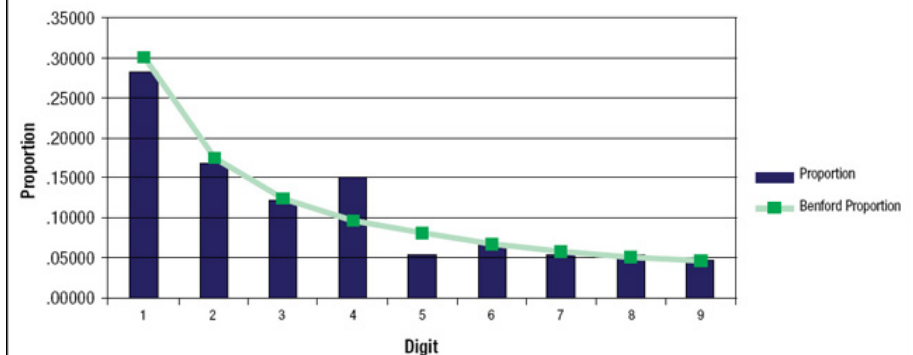


Figure 2—Benford's Law Test/Comparison



Source: <http://www.isaca.org/Journal/Past-Issues/2011/Volume-3/Pages/Understanding-and-Appling-Benford's-Law.aspx>

Why Study Statistics?

Marketing

How to draw conclusions on the needs of a **POPULATION** based only on **SAMPLE** data?

■ Marketing Research

- ❑ Construction of questionnaires and scales
- ❑ **Understand the needs** of individuals in marketplace, to create **marketing strategies and plans**

Interesting facts on mobile website users expectations

*The following stats are compiled from a study from Google (conducted by Sterling Research and SmithGeiger, independent market research firms). The report surveyed 1,088 US adult smartphone Internet users in July 2012.

Friendly = More likely to buy

Unfriendly = More likely to leave

67%

"A mobile-friendly site makes me more likely to buy a product or use a service."



61%

"If I don't see what I'm looking for right away on a mobile site, I'll quickly move on to another site."



Turning visitors into customers*

If your site offers a great mobile experience your chance of getting new customers that visit your site increases dramatically.

- 74% of people say they are more likely to return to a website if it is mobile friendly.

Hurting your business and helping your competition*

If you have a poor mobile experience and your competitors have built a great experience for mobile users chances are they will benefit and you will be hurt.

- 61% of users said if they don't find what they are after right away on a mobile site, they will

Non mobile sites can damage your reputation*

If a site isn't designed for mobile users it can leave users feeling frustrated.

- 48% of users say they feel frustrated and annoyed when they get to a site that is not mobile friendly.
- 52% of users said that a bad mobile experience made them less likely to engage with a company.
- 48% said that if a site didn't work well on their smartphones, it made them feel like the company didn't care about their business.

While you may agree or disagree with some of the responses the thing to take away from this is that it's imperative to your business to ensure you provide a great mobile experience.

Source: <http://www.onlinemarketing.fuelgroup.com.au/mobile-websites>

Why Study Statistics?

Economics

■ Economics Indicators

- Allow **analysis of economic performance** and **predictions of future performance**

(July 1997 = 100)

Centa-City Leading Index CCL

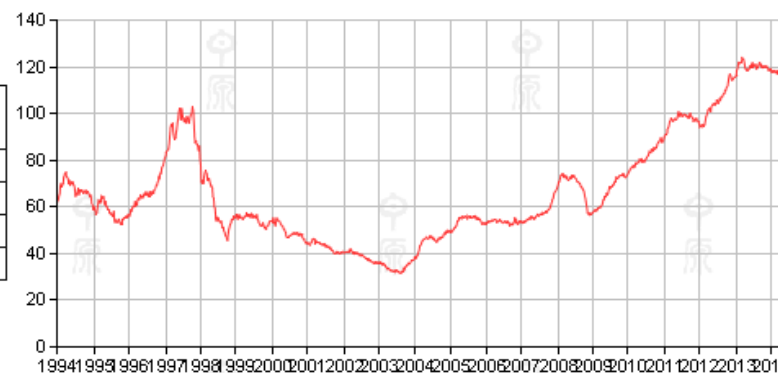
Announced every Friday, latest on 2014/08/15; reflecting seasonally adjusted annual rate of change in the index of leading indicators for the Hong Kong economy from 2014/08/04 to 2014/08/10 (based on scheduled formal sale & purchase of new housing units within 14 days after preliminary S&P)

How to construct
ECONOMICS INDEX?

	This Week	Previous Week	Previous Month
[Centa-City Leading Index]	125.66	↑ 1.34 %	↑ 1.86 %
[Centa-City (large units) Leading Index]	131.41	↑ 2.03 %	↑ 3.53 %
[Centa-City (small/medium units) Leading Index]	124.01	↑ 1.21 %	↑ 1.55 %
[Mass Centa-City Leading Index]	125.64	↑ 1.14 %	↑ 1.43 %

[Centa-City Leading Sub-index]

	This Week	Previous Week	Previous Month
HK	135.79	↑ 1.23 %	↑ 1.7 %
KLN	125	↑ 1.86 %	↑ 1.7 %
NT (East)	126.35	↑ 0.76 %	↑ 0.42 %
NT (West)	107.26	↑ 0.13 %	↑ 1.31 %



Source: <http://hk.centadata.com/ccl/ccie.htm>

Why Study Statistics?

Finance

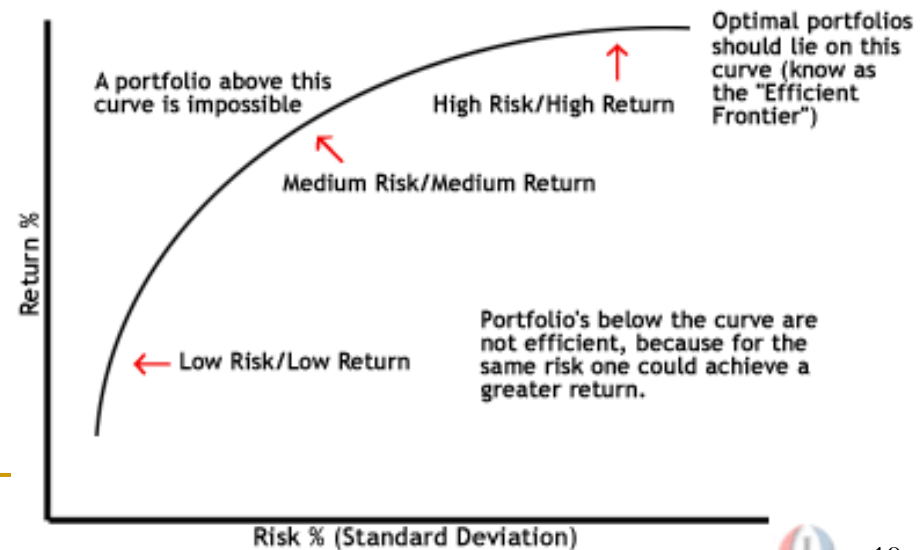
■ Risk and Portfolio Management

- ❑ Use statistical models to analyze the market
- ❑ Efficient frontier of portfolio (a basket of stocks)

How to measure the portfolio **EXPECTED RETURN** and **RISK**?

Modern portfolio theory (MPT) is a theory of **finance** that attempts to maximize portfolio expected return for a given amount of portfolio risk, or equivalently minimize risk for a given level of expected return, by carefully choosing the proportions of various **assets**. Although MPT is widely used in practice in the financial industry and several of its creators won a **Nobel memorial prize** for the theory,^[1] in recent years the basic assumptions of MPT have been widely challenged by fields such as **behavioral economics**.

More technically, MPT models an asset's return as a normally distributed function (or more generally as an **elliptically distributed random variable**), defines **risk** as the **standard deviation** of return, and models a portfolio as a weighted combination of assets, so that the return of a portfolio is the weighted combination of the assets' returns. By combining different assets whose returns are not perfectly positively **correlated**, MPT seeks to reduce the total **variance** of the portfolio return. MPT also assumes that investors are **rational** and markets are **efficient**.



Source:

http://en.wikipedia.org/wiki/Modern_portfolio_theory#The_efficient_frontier_with_no_risk-free_asset

Why Study Statistics?

Cont'd

Big DATA

- In May 2011, McKinsey published “**Big data: The next frontier for innovation, competition, and productivity**”. By 2018, the United States alone could face a **shortage** of 140,000 to 190,000 **people with deep analytical skills** as well as 1.5 million **managers and analysts with the know-how to use the analysis of big data** to make effective decisions
- In March 2012, the Obama administration announced the **Big Data Research and Development Initiative**, which explored how big data could be used to address important problems faced by the government
- The White House announced a national "Big Data Initiative" that consisted of six Federal departments and agencies committing **more than \$200 million** to big data research projects

Why Study Statistics?

Cont'd

Statisticians – Dream job of the next

“I keep saying that the **sexy job** in the next 10 years will be **statisticians**.”



Google's Chief Economist, Hal Varian, interviewed by McKinsey Quarterly in January 2009

More Than Just Numbers

8.32, 7.91, 9.64, 9.18, 10.33, 7.46

- As just numbers, this list is uninteresting, but what can you say if this list represents:
 - ❑ Weight of a newborn puppy?
 - ❑ Minutes to run a mile?

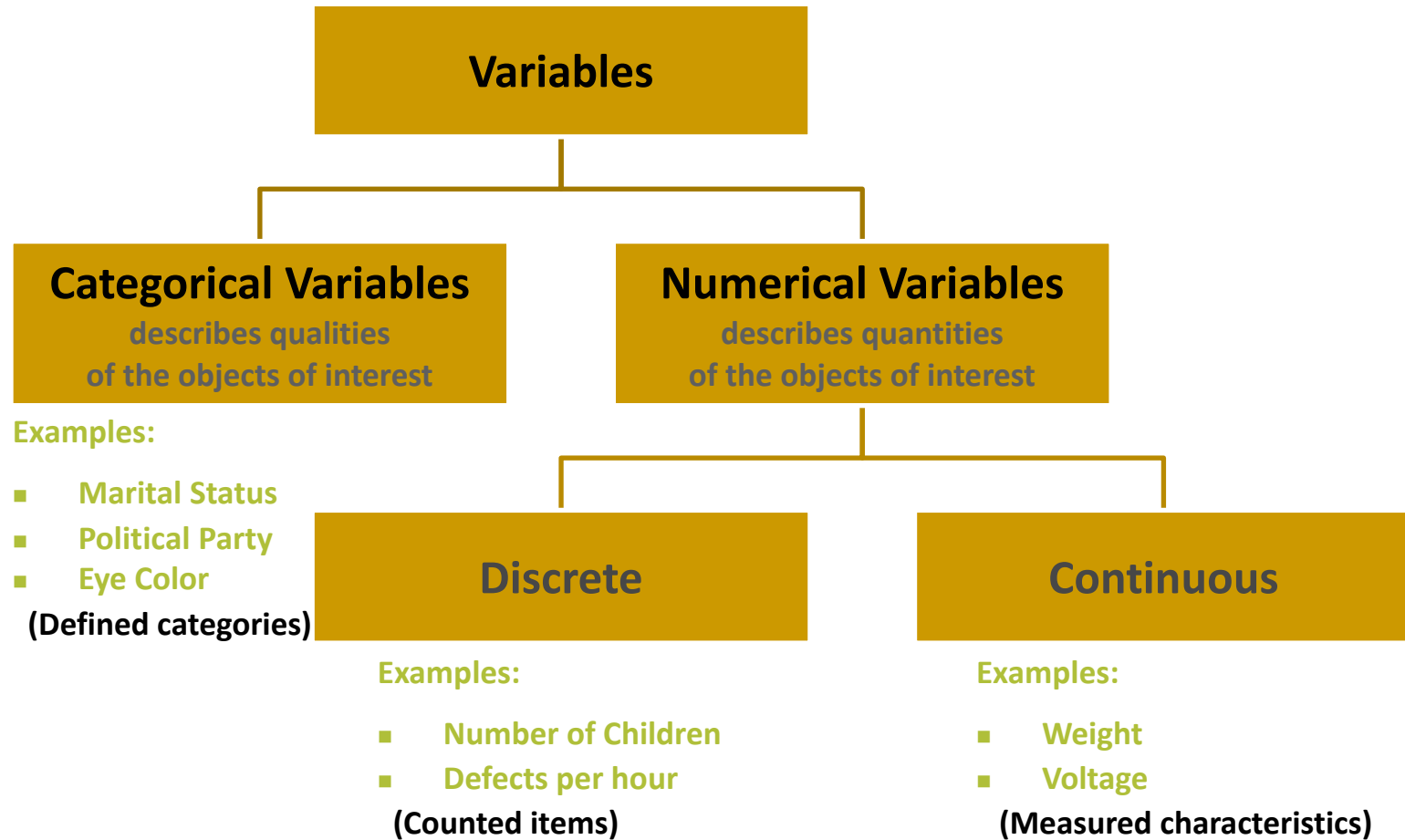
Variables

- A **variable** is any characteristic, number, or quantity that can be measured or counted
- E.g.
 - ❑ A person's gender
 - ❑ The weight of a newborn puppy
 - ❑ The concentration of CO₂ in the atmosphere
 - ❑ People's income
 - ❑ Examination grade
 - ❑ Vehicle type

Data

- **Data** are the values measured or observed for each variable of each object
 - Data can be numeric or categorical
 - Numeric data takes numeric values and work well for statistics
 - Satisfaction rating ranging from 1 to 10
 - Number of students attending the lecture
 - People's income
 - Categorical data
 - People's gender
 - Examination grade
-
- Vehicle type

Types of Variables



Types of Variables – Exercise

Cont'd

Age	Gender	Major	Credits	District	GPA
18	Male	Management Sciences	16	Hong Kong Island	3.6
21	Male	Accountancy	18	New Territories	3.1
20	Female	Marketing Information Mgt	16	Kowloon	2.8

■ Numerical

■ Categorical

Numerical or Categorical?

Why are you in college? Answer:

1. Personal Growth 2. Career Opportunities
3. Parental Pressure 4. Personal Networking

Results: 1, 4, 3, 2, 2, 1, 2, 3, 3, 1, 4, 2

- **Coding categorical data** with numbers:
Although the above data values are numbers,
the variable is still categorical
- **Reason for coding:** Easier to input into a
computer

Coding Yes/No Questions

- Use 0 for “No” and 1 for “Yes”
- Useful for data with only two possible values
 - ❑ True or False
 - ❑ Black or White
 - ❑ Success or Failure
 - ❑ Dead or Alive

Organizing and Visualizing Data

Variables

Categorical Variables

describes qualities
of the objects of interest

- Summary Table
- Bar Chart
- Pie Chart

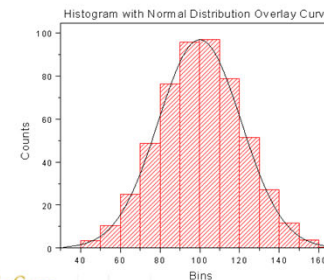
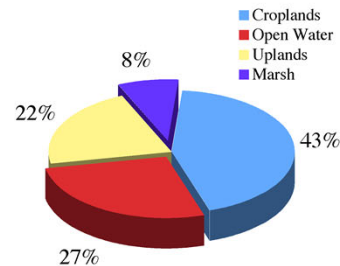
Numerical Variables

describes quantities
of the objects of interest

- Frequency Distribution
- Histogram

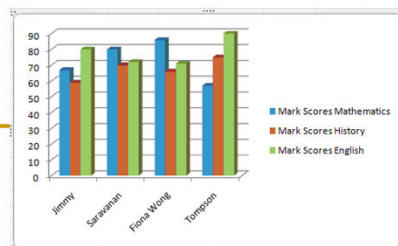
SUMMARY TABLE

Effective tax rates on bonus remuneration and dividends				
	Bonus	Dividends		
		Small Co	Medium Co	Large Co
2012/13				
Higher-rate taxpayer	49.04%	40.00%	43.75%	43.00%
Top-rate taxpayer	57.82%	48.89%	52.08%	51.44%
2013/14				
Higher-rate taxpayer	49.04%	40.00%	43.75%	43.00%
Top-rate taxpayer	53.43%	44.44%	47.91%	47.22%



CV_HGC_0607 2006				
S7513600	Frequency	Percent	Cumulative Frequency	Cumulative Percent
5	2	0.03	2	0.03
6	19	0.25	21	0.28
7	41	0.55	62	0.83
8	260	3.48	322	4.31
9	394	5.27	716	9.58
10	491	6.57	1207	16.15
11	580	7.76	1787	23.91
12	5681	76.02	7468	99.93
95	5	0.07	7473	100.00

Frequency Missing = 1511



Organizing and Visualizing Data

Cont'd

■ Organizing Categorical Data

- Suppose you asked 60 customers to pick which of the three colours, say green, red, or blue they like best for a product

- The data might look like this: green, red, green, green, red, red, blue, blue, green, red, green, blue, red, blue, green, green, blue, green, green, blue, green, blue, green, red, blue, green, green, green, green, red, red, red, blue, green, green, green, green, blue, red, red, green, green, red, blue, green, red, green, green, blue, red, green, red, green, blue, blue, blue, green, green, green, green, green

Organizing and Visualizing Data

Cont'd

- A natural way to describe the data is counting how many of each colour you have got

- A summary table:

Colour	Number of Customers
blue	15
green	30
red	15
Total	60

- It is accustomed to list the values of the variable in alphabetical order of the category, or in descending (or ascending) order of the count

- In statistical context, the proper name for count is called **frequency**

Organizing and Visualizing Data

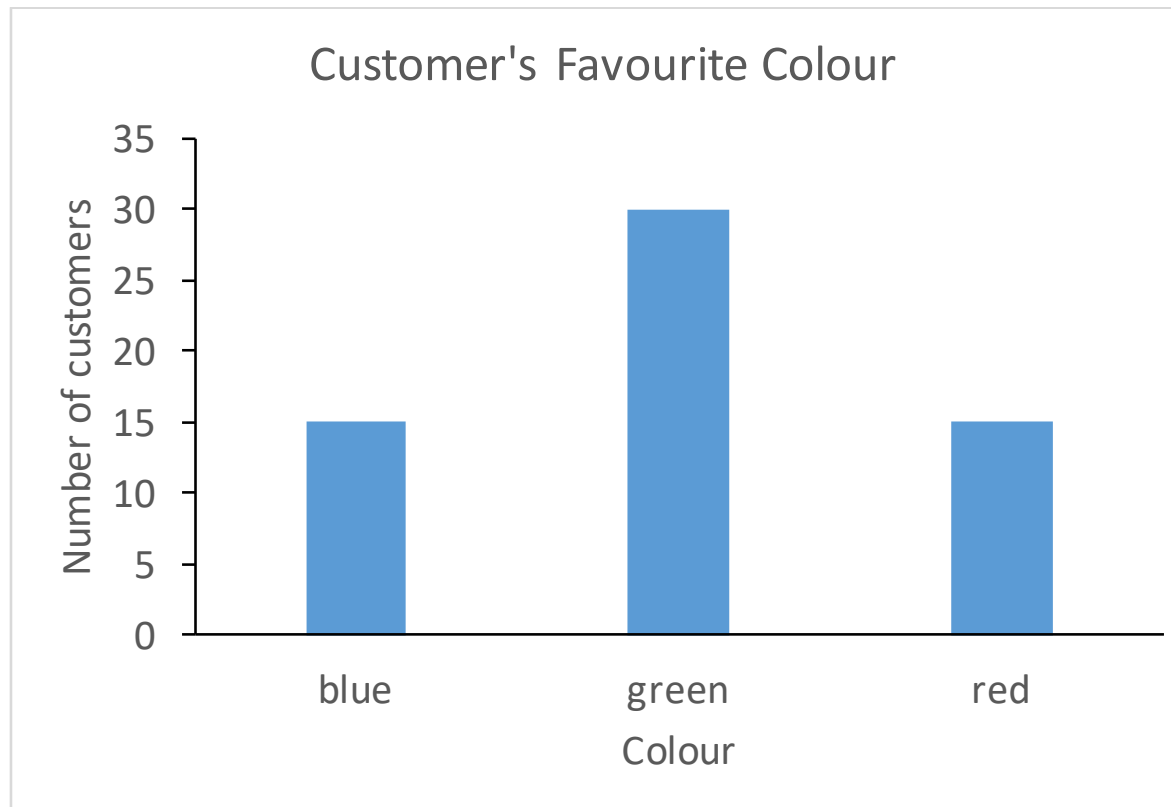
Cont'd

- You cannot tell from the table which customer picked what colour. This information is often unimportant in reporting
- The table tells us
 - ❑ 15 customers picked blue, 30 customers picked green, and 15 customers picked red
 - ❑ More customers picked green than the other two colours
 - ❑ About the same number of customers picked the two other colours

Organizing and Visualizing Data

Cont'd

■ (Frequency) Bar Chart



Organizing and Visualizing Data

Cont'd

■ Features of a Bar Chart

- ❑ It is accustomed to arrange the bars in the alphabetical order of the categories of the variable, or in descending (or in ascending) order of the count
- ❑ It is up to you to decide the gap between two bars, as long as the gaps are the same
- ❑ It is up to you to decide the width of each bar, as long as they all have the same width
 - Keeping the widths of the bars equal ensuring the area of each bar proportional to the number of individuals in that category
- ❑ The height of each bar is proportional to the number of individuals in that category

Organizing and Visualizing Data

Cont'd

- The proportion of each category can also be included in the summary table and bar chart

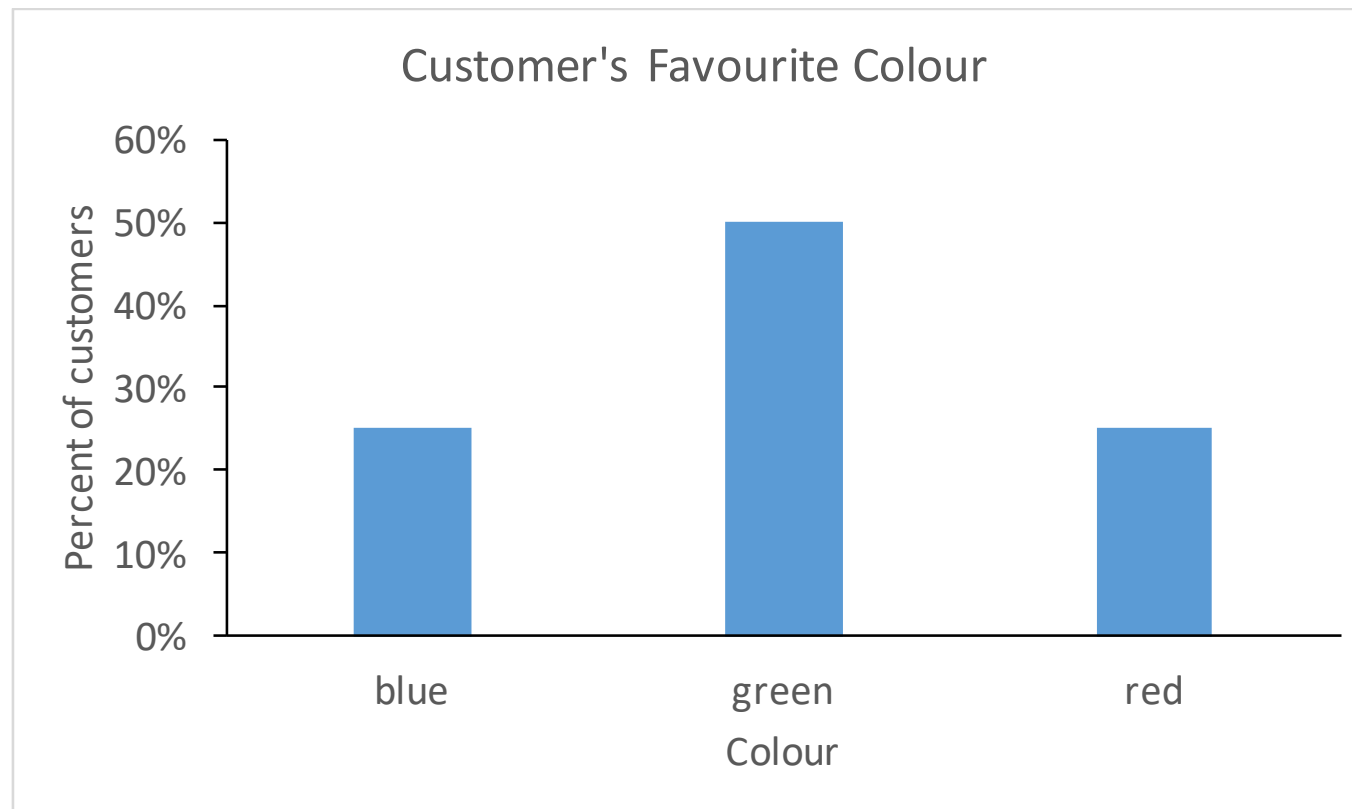
Colour	Number of Customers	Percent of Customers
blue	15	25%
green	30	50%
red	15	25%
Total	60	100%

- In statistical context, percent is called **relative frequency**

Organizing and Visualizing Data

Cont'd

■ (Relative Frequency) Bar Chart



Organizing and Visualizing Data

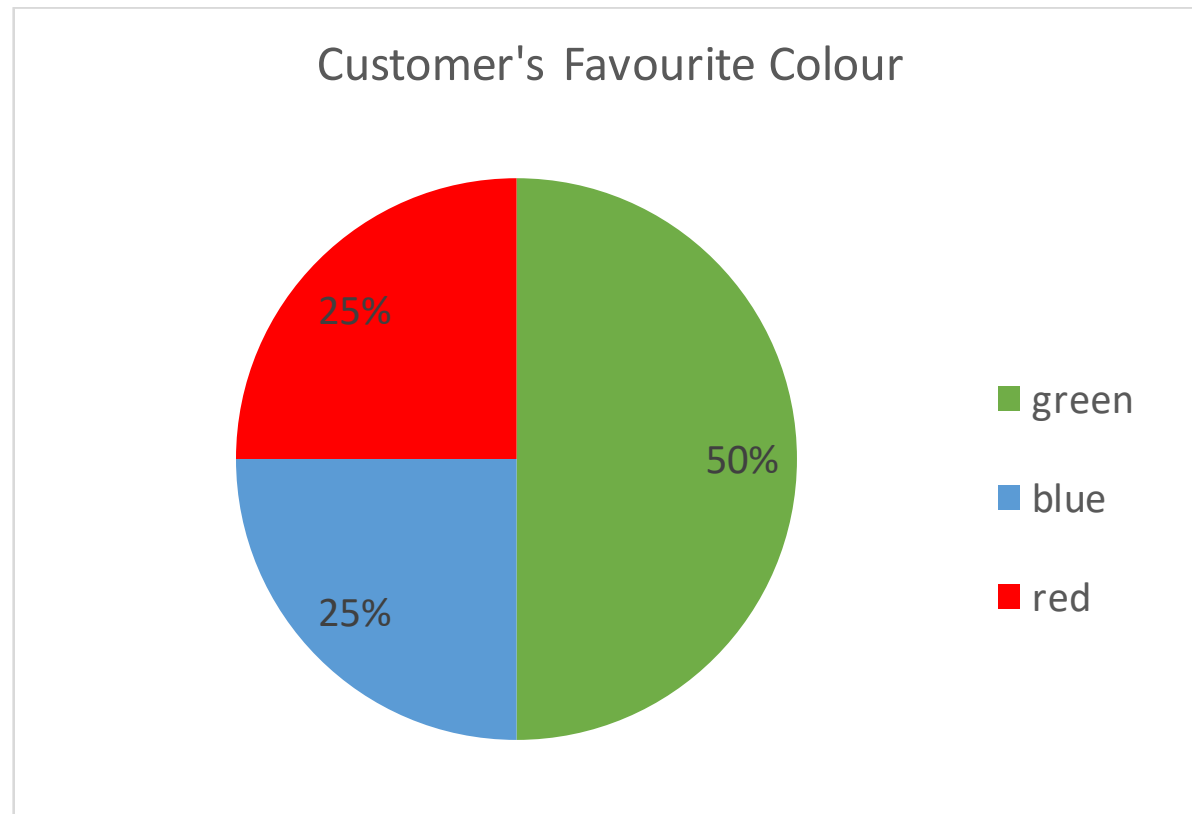
Cont'd

- The two bar charts, frequency and relative frequency, looks the same if the vertical scales are removed
- In fact, the sizes of the bars depend on percent, not on exactly how many
 - The bar chart would look the same if you had started out with 150 blue, 300 green, and 150 red

Organizing and Visualizing Data

Cont'd

■ Pie Chart



Organizing and Visualizing Data

Cont'd

■ Features of a Pie Chart

- ❑ It shows the size relationship between the categories of the variable and the variable itself
 - ❑ The slices are mutually exclusive. The sum of all slices equal to 100 percent
 - ❑ It is accustomed to arrange the slices in the alphabetical order of the categories of the variable, or in descending (or in ascending) order of the count
 - ❑ Slices of very low percent may need to be combined with others
-
- ❑ Numbers (percent or count) should be shown as it is difficult to compare slices of similar size

Organizing and Visualizing Data

Cont'd

■ Organizing Numerical Data

- Suppose you asked 100 people about the amount they spent in their last visit to supermarket

- The data might look like this: 44.8, 230.5, 303.6, 70.8, 534.4, 166.2, 466, 85.1, 63, 47.8 36.5, 35.7, 12.7, 11.9, 297.5, 74.1, 77.1, 251.2, 127.1, 118.6, 211.2, 221.9, 49.1, 349.1, 556.6, 768, 231.7, 247.2, 87.4, 304.3, 311.3, 825.8, 15.9, 526, 5.2, 156.7, 65.2, 143.3, 138.5, 478.4, 124.2, 205.1, 90.8, 3.1, 334.8, 7.4, 113.8, 79.2, 128.8, 26.6, 15.2, 554.4, 2.9, 70.2, 540.7, 36.4, 588.9, 151.5, 14.2, 235.7, 13.7, 187.4, 817.8, 140.3, 114.9, 219.5, 31.4, 99.4, 47.3, 111.8, 230.2, 478.2, 4.6, 783.5, 483.5, 99.3, 92.8, 464.2, 172.9, 380.1, 234.5, 120.2, 100.3, 109.8, 276.1, 157.7, 192.9, 13.1, 62.2, 44.2, 35.9, 239.9, 193.8, 591.9, 249.1, 17.9, 89.3, 369.1, 38.2, 154.3

Organizing and Visualizing Data

Cont'd

- Similar to categorical data, numerical data can be presented in the form of table. It is called **frequency distribution**

- The frequency distribution is a summary table in which the data are arranged into numerically ordered classes

Amount Spent (\$)	Frequency
0 - < 100	40
100 - < 200	22
200 - < 300	15
300 - < 400	7
400 - < 500	5
500 - < 600	7
600 - < 700	0
700 - < 800	2
800 - < 900	2
900 - < 1000	0
Total	100

Organizing and Visualizing Data

Cont'd

- Steps to construct a frequency distribution
 1. Sort data in ascending order: 2.9, 3.1, 4.6, 5.2, 7.4, ...
 2. Find the **range**: $825.8 - 2.9 = 822.9$
 3. Select the number of classes: 10
 4. Compute the **class interval** (width): $822.9 / 10 = 82.29$
 - Round up to a convenient number, say 100
 5. Determine **class boundaries** (limits):
 - Class 1: 0 but less than 100
 - Class 2: 100 but less than 200
 - ...
 6. Assign the observation to each class and count the number of observations

Organizing and Visualizing Data

Cont'd

- Features of frequency distribution
 - ❑ Exact value of each observation is lost
 - ❑ The width of each interval is identical
 - Width can be unequal. However, it should be done so only under very special circumstances, such as the data is sparsely distributed, or have a very long tail at one or both ends
 - ❑ The lower value of the first class interval is often the smallest value in the data, or a smaller value which is selected for the reason of convenience, such as 0
 - ❑ Class boundaries include the left endpoint, but not the right
 - Other endpoint policy can be adopted, but need to be consistent

Organizing and Visualizing Data

Cont'd

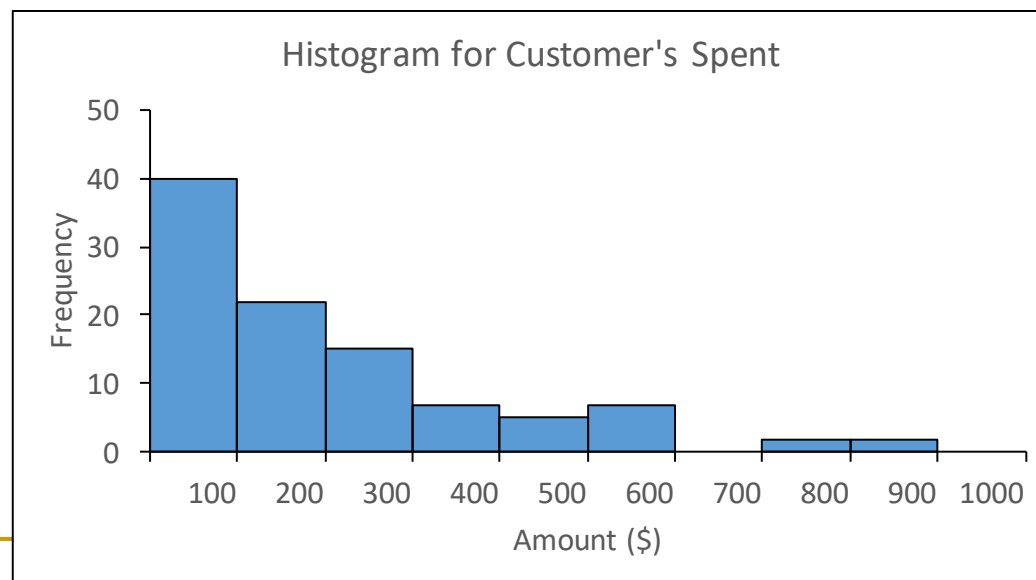
- Features of frequency distribution
 - The number of classes depends on the number of values in the data
 - Large data range and high number of observations allow a larger number of class. In general, 5 to 15 classes will be sufficient
 - The width of a class depends on the number of classes adopted and the data range
 - To determine the width of a class, you divide the range (highest value – lowest value) of the data by the number of classes desired

Organizing and Visualizing Data

Cont'd

■ Histogram

- A histogram is a bar chart for grouped numerical data in which the frequencies of each group of numerical data are represented as individual vertical bars
- For example:



Organizing and Visualizing Data

Cont'd

■ Features of a histogram

- ❑ The chart is made from the constructed frequency distribution
- ❑ The height of the bars is in proportion to the frequency of intervals
- ❑ There is no gap between bars
 - If an interval has 0 frequency, the height of the bar in the histogram is 0
- ❑ The width of the bars must be identical because the width of the intervals are identical
- ❑ The bar must be drawn in the same sequence as of the intervals in the frequency distribution

Organizing and Visualizing Data

Cont'd

- The proportion of each class can also be displayed in frequency distribution (or histogram)

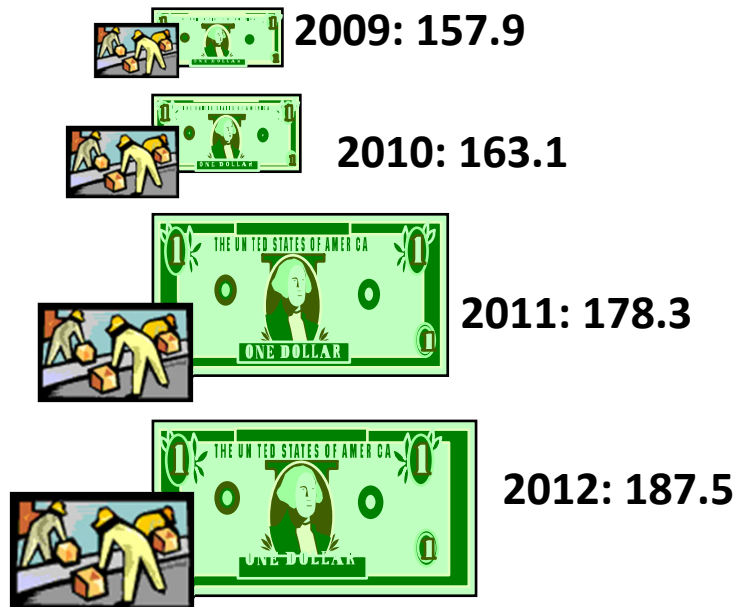
Amount Spent (\$)	Frequency	Relative Frequency
0 - < 100	40	0.40
100 - < 200	22	0.22
200 - < 300	15	0.15
300 - < 400	7	0.07
400 - < 500	5	0.05
500 - < 600	7	0.07
600 - < 700	0	0.00
700 - < 800	2	0.02
800 - < 900	2	0.02
900 - < 1000	0	0.00
Total	100	1.00

Faulty Graphs: Chart Junk



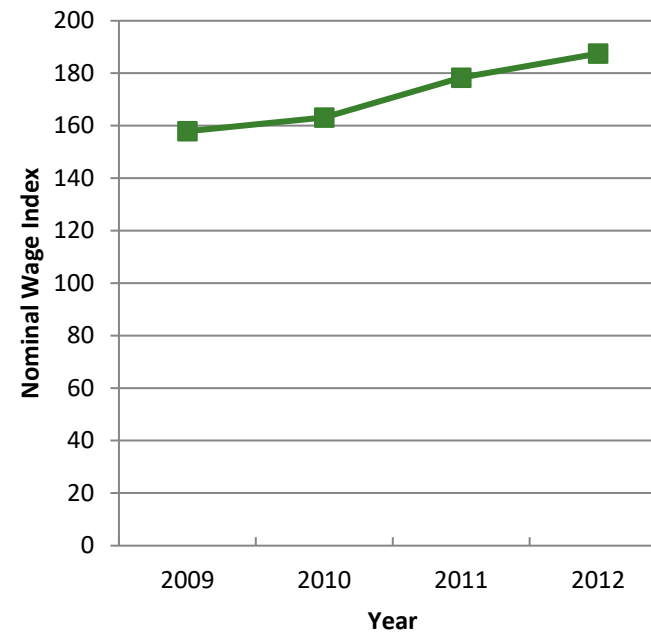
Bad Presentation

Nominal Wage Index in Hong Kong
(Sept 1992=100)



Good Presentation

Nominal Wage Index in Hong Kong
(Sept 1992=100)



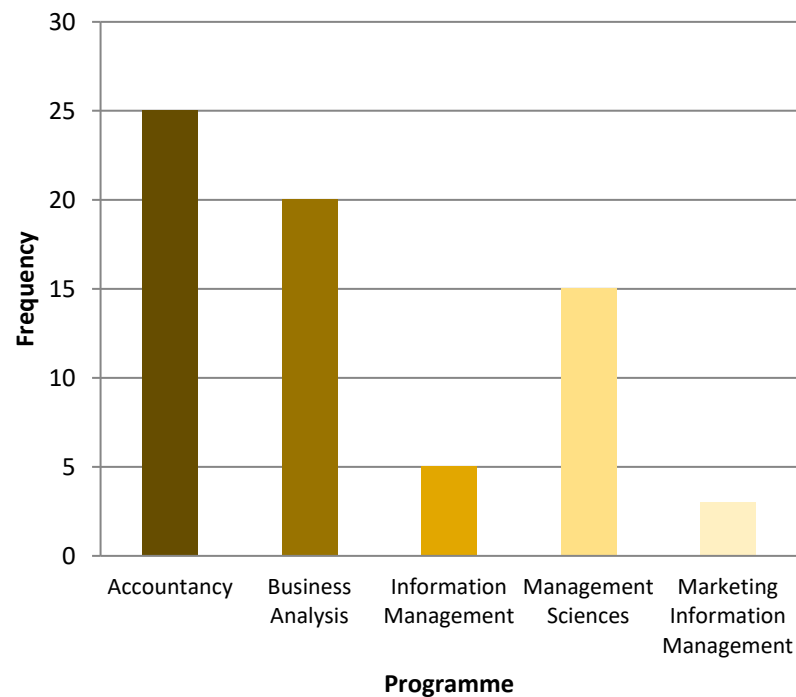
Faulty Graphs: No Relative Basis

Cont'd



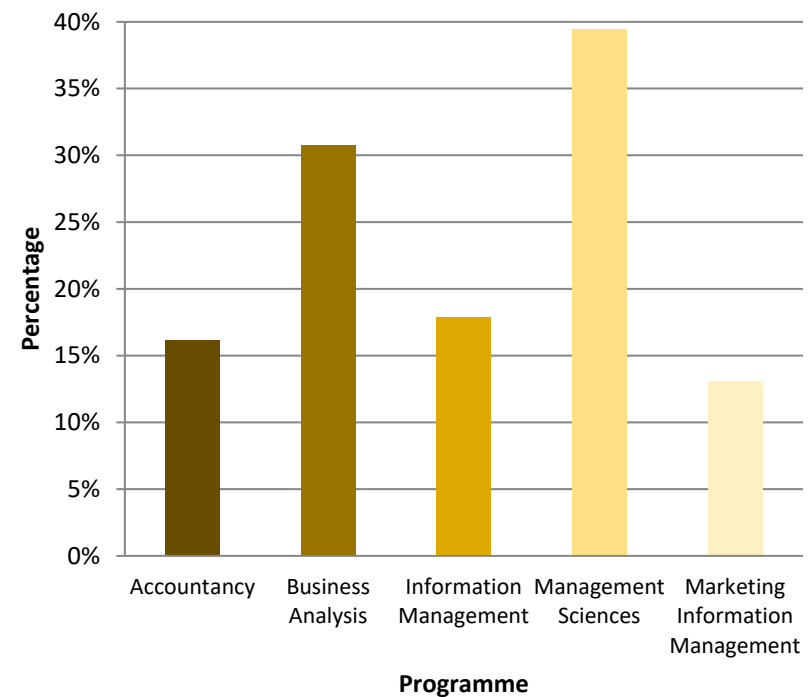
Bad Presentation

A's Obtained by Students in MS2200



Good Presentation

A's Obtained by Students in MS2200

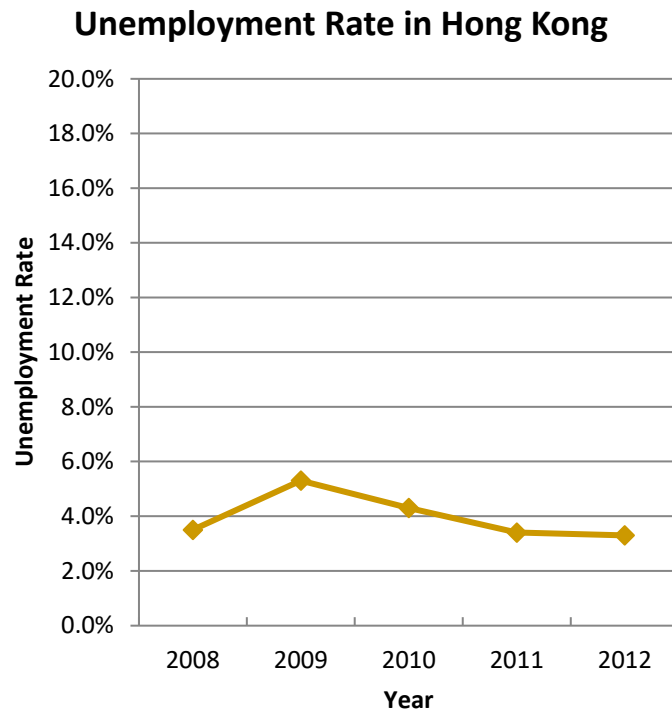


Faulty Graphs: Compressing the Vertical Axis

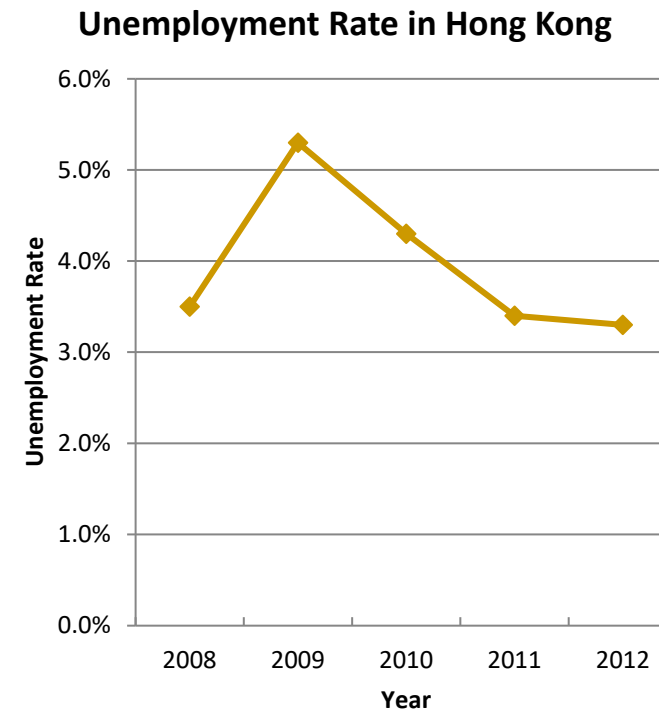
Cont'd



Bad Presentation



Good Presentation



Faulty Graphs: No Zero Point on the Vertical Axis

Cont'd



Bad Presentation



5/11/07重點新聞：恒指收市報28942點，跌1526點



Good Presentation



5/11/07重點新聞：恒指收市報28942點，跌1526點

Principles of Excellent Graphs

- The graph should **not distort** the data
- The graph should **not contain unnecessary adornments** (sometimes referred to as chart junk)
- The scale on the vertical axis should **begin at zero**
- All **axes** should be properly **labeled**
- The graph should contain a **title**
- The **simplest** possible graph should be used for a given set of data

Use and Misuse of Statistics

HP to Lay Off 9,000 in Enterprise Services Revamp

by Douglas McIntyre  Jun 1st 2010 9:35AM

Updated Jun 1st 2010 9:43AM

Hewlett-Packard (HPQ) announced Tuesday that it would put \$1 billion into its enterprise services division -- and lay off 9,000 workers in the process. In the last quarter for which it has reported results, which ended on Jan. 31, HP had revenue of \$21 billion and net income of \$2.3 billion. The enterprise unit provides consulting, outsourcing and technology services.



Paul Sakuma, AP

HP's 10-Q shows that the revenue from the division was \$8.7 billion, down slightly from the same period a year ago. Operating income was \$1.3 billion. So, the business is critical to HP's success, but it isn't doing terribly well.

The enterprise operation is part of HP's plan to diversify beyond hardware and become more competitive with rival IBM (IBM). HP bought information-technology consulting firm EDS in May 2008 for \$13.1 billion. In March 2009, HP said it would eliminate 24,600 jobs during the integration of EDS -- but it's not entirely clear whether the new layoffs are a subset of those or in addition to them.

Use and Misuse of Statistics: Faulty Percentages

Cont'd




$$\frac{1200 - 400}{400} \times 100\% = 200\%$$

$$\frac{1200 - 800}{800} \times 100\% = 50\%$$

Summary Definitions

- The **central tendency** is the extent to which all the data values group around a typical or central value
- The **variation** is the amount of dispersion or scattering of values
- The **shape** is the pattern of the distribution of values from the lowest value to the highest value

Measures of Central Tendency

 **太陽報** 太陽報 – 2015年4月1日星期三上午5:50

Cont'd

沙田住64萬人全港稱冠

【本報訊】政府統計處公布去年人口及住戶統計資料，全港人口七百一十五萬二千，沙田最多人居住，人口達六十四萬多，其次是觀塘。全港最富貴是中西區，住戶每月入息中位數是三萬五千元，收入最低則是深水埗，約為中西區一半。黃大仙陰盛陽衰情況最嚴重，而大埔的男女比例則最平均。

統計處昨發表《2014年按區議會分區劃分的人口及住戶統計資料》報告書，刊載人口特徵如年齡、性別、婚姻狀況、住戶入息等。全港最富貴是中西區，住戶每月入息中位數達三萬五千元，其次是灣仔區三萬四千元，拋離西貢的三萬零八百元。

收入最低是深水埗，為一萬八千元，觀塘區則稍高，有一萬九千元，葵青區一萬九千六百元，這些地區的住戶每月入息中位數均低過全港中位數二萬三千五百元。自置居所比例最高是西貢，達百分之六十四點七。

全港人口最多的沙田區，居民達六十四萬八千二百，其次是觀塘六十三萬九千九百人，元朗五十九萬五千一百人，東區五十七萬九千四百人。最年輕是元朗和離島，年齡中位數為三十九歲，最多長者在東區和黃大仙，年齡中位數是四十五歲，而全港中位數則是四十二歲。

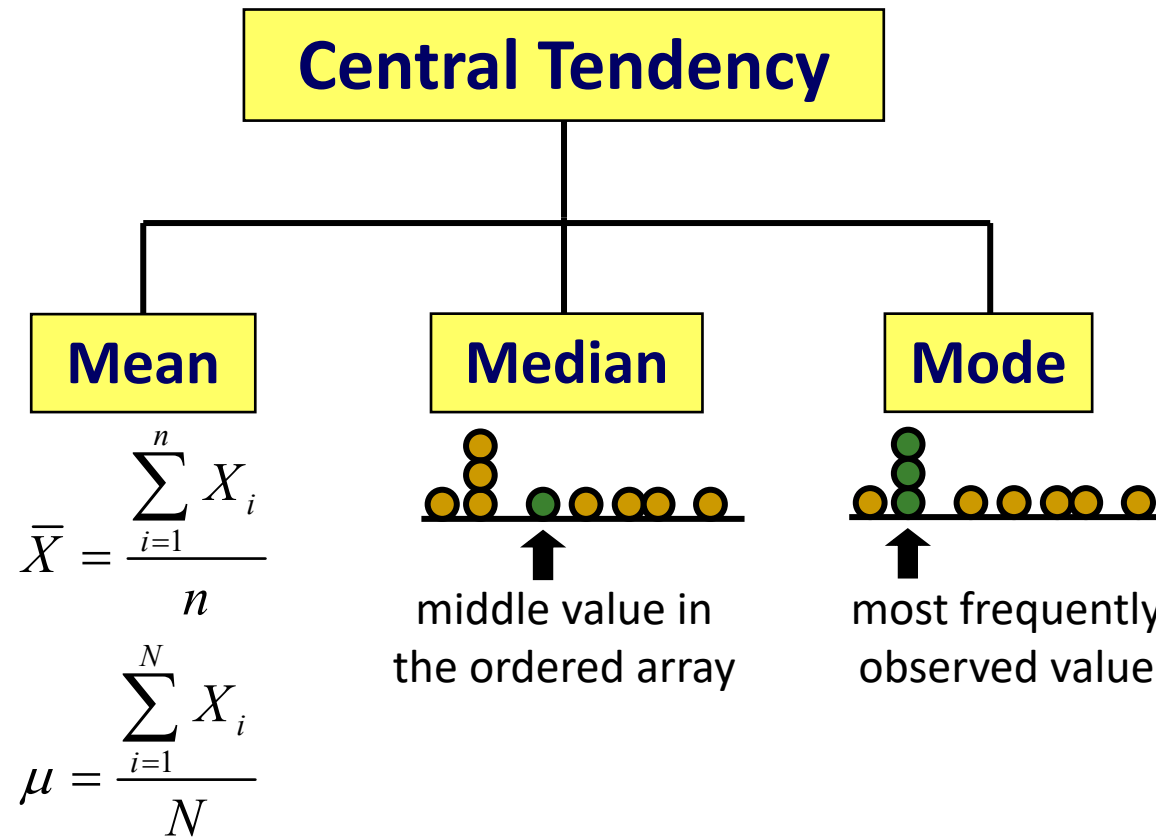
黃大仙陰盛陽衰

男性已婚比率是百分之六十一點四，高於女性的百分之五十六點一。女多男少最嚴重是黃大仙，每千名女性只有八百九十六名男性，大埔男女比例最平均，每千名女性有九百五十四名男性。

What is
Median? Why
it is used?

Measures of Central Tendency

Cont'd



Measures of Central Tendency:

The Mean

Cont'd

■ Sample mean

pronounced x-bar

$$\bar{X} = \frac{\sum_{i=1}^n X_i}{n} = \frac{X_1 + X_2 + \cdots + X_n}{n}$$

← Sample Size

■ Population mean

pronounced mu

$$\mu = \frac{\sum_{i=1}^N X_i}{N} = \frac{X_1 + X_2 + \cdots + X_N}{N}$$

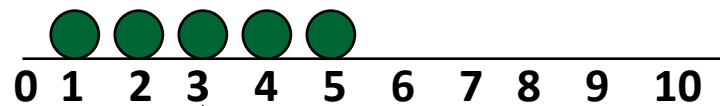
← Population Size

Measures of Central Tendency:

The Mean

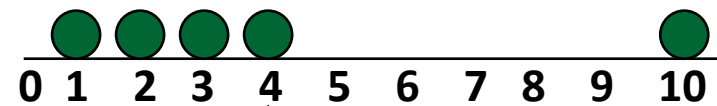
Cont'd

- The most common measure of central tendency
- Affected by extreme values (outliers)



Mean = 3

$$\frac{1+2+3+4+5}{5} = \frac{15}{5} = 3$$



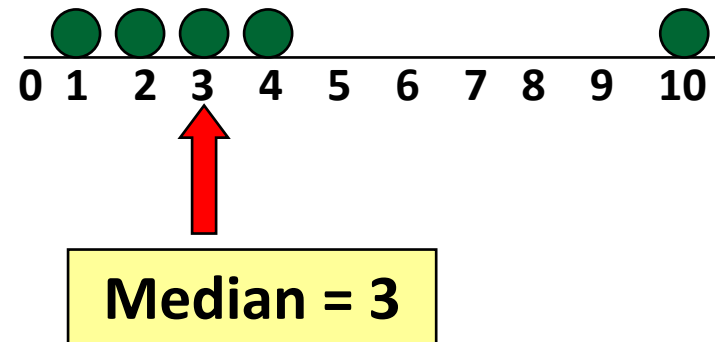
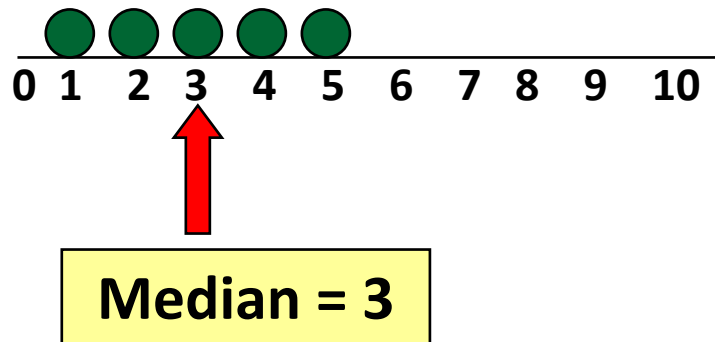
Mean = 4

$$\frac{1+2+3+4+10}{5} = \frac{20}{5} = 4$$

Measures of Central Tendency: The Median

Cont'd

- Robust measure of central tendency
- In an **ordered array**, the median is the “middle” number (50% above, 50% below)
 - If n or N is odd, the median is the middle number
 - If n or N is even, the median is the average of the 2 middle numbers
- Not affected by extreme values (outliers)

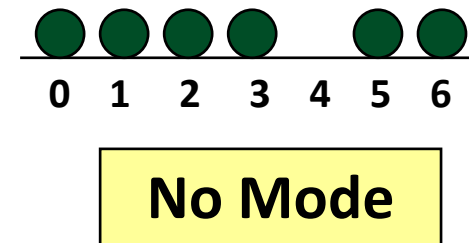
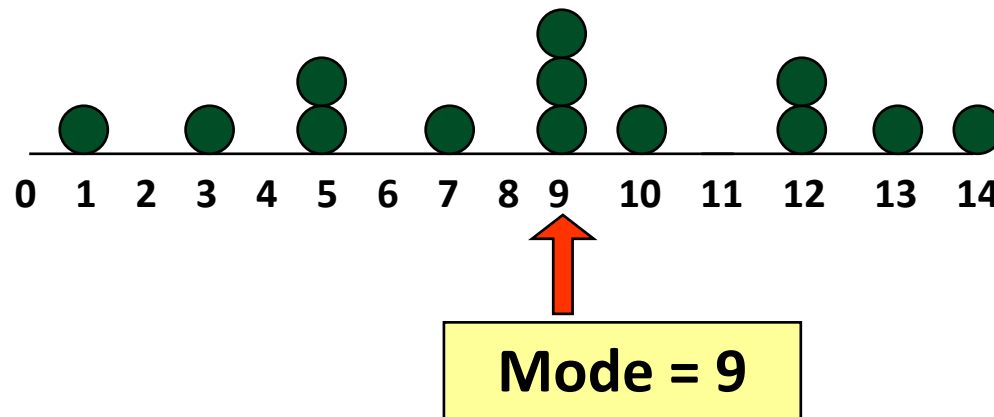


Measures of Central Tendency:

The Mode

Cont'd

- Value that occurs most often
- Not affected by extreme values (outliers)
- Used for both numerical and categorical data
- There may be no mode
- There may be several modes



Effects of Outliers

Imagine that the five graduating seniors on a college basketball team receive the following first-year contract offers to play in the National Basketball Association (zero indicates that the player did not receive a contract offer):

0 0 0 0 \$3,500,000

The mean contract offer is:

$$\text{mean} = \frac{0 + 0 + 0 + 0 + \$3,500,000}{5} = \$700,000$$

Is it therefore fair to say that the average senior on this basketball team received a \$700,000 contract offer?

Why Census & Statistics Department reported the median household income?

When the Hong Kong Housing Authority revises the rent of public housing, they need a reference for the average income of tenants. Should they consider the mean or the median?

Comparison of Mean, Median & Mode

Measure	Definition	How common?	Existence	Takes every value into account?	Affected by outliers?	Advantages
Mean	$\frac{\text{sum of all values}}{\text{total number of values}}$	most familiar "average"	always exists	yes	yes	commonly understood; works well with many statistical methods
Median	middle value	common	always exists	no (aside from counting the total number of values)	no	when there are outliers, may be more representative of an "average" than the mean
Mode	most frequent value	sometimes used	may be no mode, one mode, or more than one mode	no	no	most appropriate for qualitative data (see Section 2.1)

Measures of Variation

2016 JUPAS (HKDSE) Admission Scores

JS1006 - Department of Management Sciences (Bachelor of Business Administration)											
Admission Score Calculation											
Weighted Admission Score		Sample Cases	HKDSE Results							Score Calculation with Weighting Applied (4 core + 2 elective subjects with the highest weighted scores)	
			English Language	Chinese Language	Mathematics (Compulsory Part)	Liberal Studies	Elective Subject: M1/M2	Other Elective Subjects			
								Subject 1	Subject 2		Subject 3
			Weighting: 2	Weighting: 1	Weighting: 1.5	Weighting: 1	Weighting: 1.5	Weighting: 1			
Median	32.5	Student A	3	3	5*	3	5	4	3		3x2+3x1+6x1.5+3x1+5x1.5+4x1=32.5
		Student B	4	3	5	4	4	4	3	2	4x2+3x1+5x1.5+4x1+4x1.5+4x1=32.5
Lower Quartile	32	Student C	4	4	4	5	4	3	3		4x2+4x1+4x1.5+5x1+4x1.5+3x1=32
		Student D	3	4	5	3	5	4	4		3x2+4x1+5x1.5+3x1+5x1.5+4x1=32

Scoring Scale:

Core and Category A Elective Subjects							
Level	5**	5*	5	4	3	2	1
Score	7	6	5	4	3	2	1

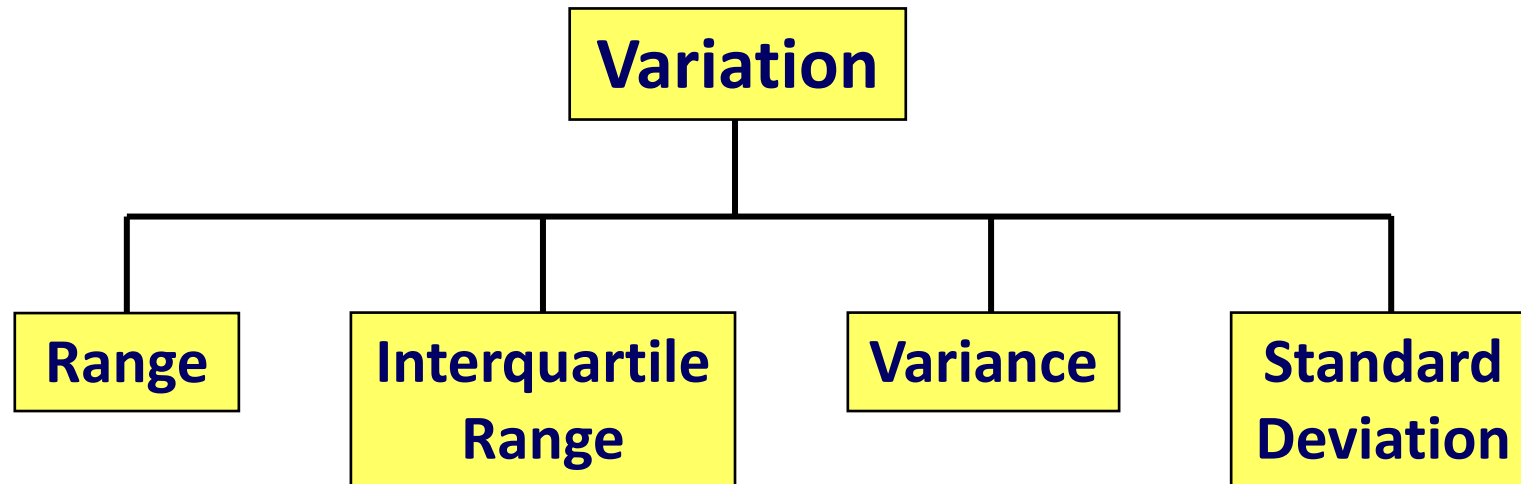
Category C Elective Subjects (Other Language Subjects)					
Grade	A	B	C	D	E
Score	5	4	3	2	1

Category B Applied Learning subjects are not considered in the above score calculations.

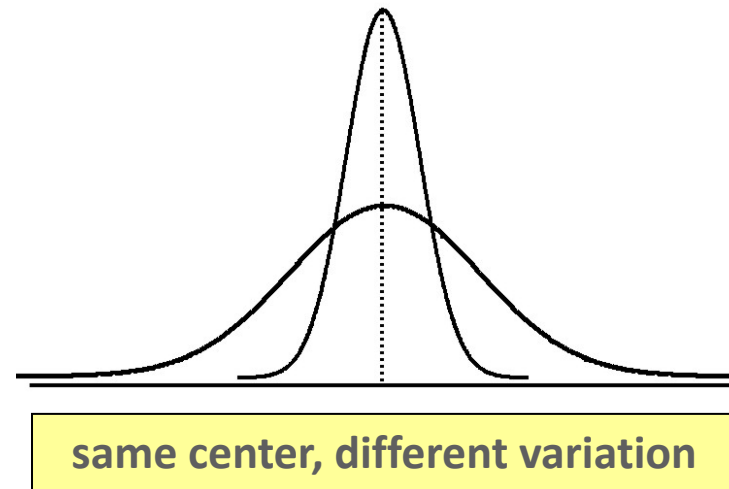
What is Lower Quartile? Why we need it?

Measures of Variation

Cont'd



Measures of variation give information on the **spread** or **variability** or **dispersion** of the data values



Measures of Variation:

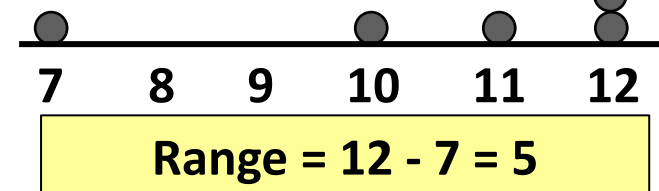
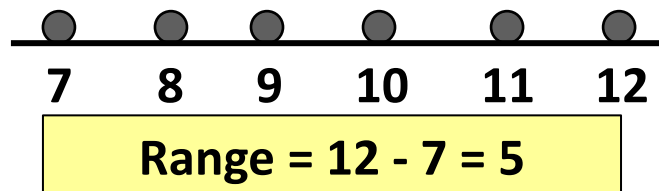
The Range

Cont'd

- Simplest measure of variation
- Difference between the largest and the smallest values

$$\text{Range} = X_{\text{Largest}} - X_{\text{Smallest}}$$

- Ignores the way in which data are distributed



- Sensitive to outliers

1,1,1,1,1,1,1,1,1,1,1,2,2,2,2,2,2,2,2,2,3,3,3,3,4,5

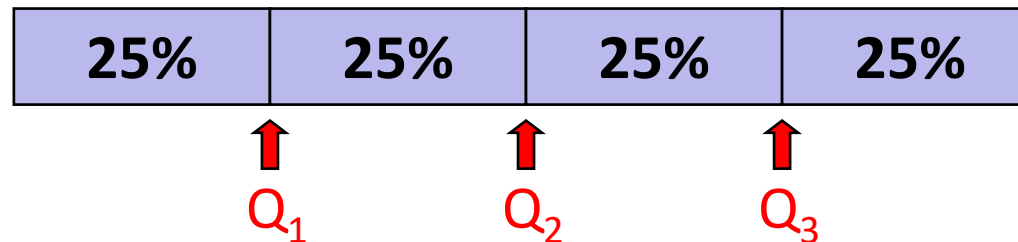
$$\text{Range} = 5 - 1 = 4$$

1,1,1,1,1,1,1,1,1,1,1,2,2,2,2,2,2,2,2,2,3,3,3,3,4,120

$$\text{Range} = 120 - 1 = 119$$

Quartiles

- Quartiles split the ranked data into 4 segments with an equal number of values per segment



- The **first quartile**, Q_1 , is the value for which 25% of the observations are smaller and 75% are larger
- Q_2 is the same as the **median** (50% of the observations are smaller and 50% are larger)
- Only 25% of the observations are greater than the **third quartile**, Q_3

Quartiles

Cont'd

$$Q_1 \text{ position: } \frac{n+1}{4}$$

$$Q_2 \text{ position: } \frac{2(n+1)}{4}$$

$$Q_3 \text{ position: } \frac{3(n+1)}{4}$$

where n is the number of observed values

When calculating the ranked position, use the following rules:

- If the result is a **whole number**, it is the ranked position to use
- If the result is a **fractional half** (e.g. 2.5, 8.5, ...), average the two corresponding data values
- If the result is **not a whole number or a fractional half**, round the result to the nearest integer to find the ranked position

Quartiles

Cont'd

Sample Data in Ordered Array: 11 12 13 16 16 17 18 21 22

Q_1 is in the $(9+1)/4 = 2.5$ position of the ranked data,
so $Q_1 = (12+13)/2 = 12.5$

Q_2 is in the $(9+1)/2 = 5^{\text{th}}$ position of the ranked data,
so $Q_2 = \text{median} = 16$

Q_3 is in the $3(9+1)/4 = 7.5$ position of the ranked data,
so $Q_3 = (18+21)/2 = 19.5$

Quartiles – Exercise

Cont'd

Data in ordered array:

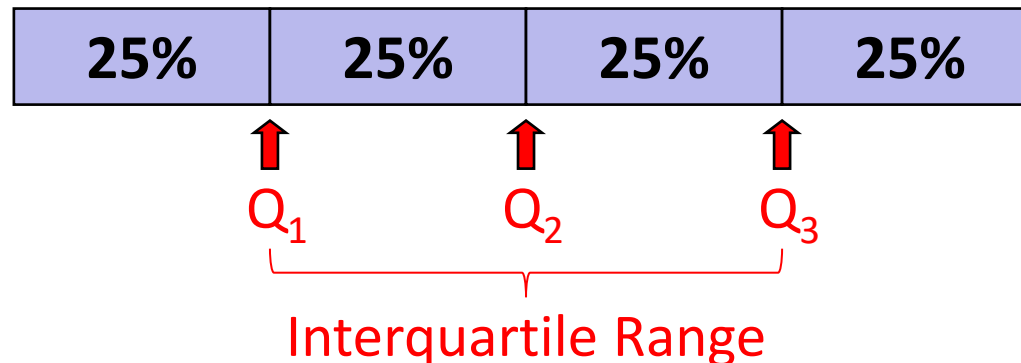
3 6 7 7 9 12

Measures of Variation:

Interquartile Range

Cont'd

- Quartiles split the ranked data into 4 segments with an equal number of values per segment



- Interquartile range is $Q_3 - Q_1$ and measures the spread in the middle 50% of the data
- Interquartile range is also called the **midspread** because it covers the middle 50% of the data
- Not influenced by outliers or extreme values

Measures of Variation:

Interquartile Range

Cont'd

Sample Data in Ordered Array: 11 12 13 16 16 17 18 21 22

Q_1 is in the $(9+1)/4 = 2.5$ position of the ranked data,
so $Q_1 = (12+13)/2 = 12.5$

Q_2 is in the $(9+1)/2 = 5^{\text{th}}$ position of the ranked data,
so $Q_2 = \text{median} = 16$

Q_3 is in the $3(9+1)/4 = 7.5$ position of the ranked data,
so $Q_3 = (18+21)/2 = 19.5$

Interquartile range = $19.5 - 12.5 = 7$

Measures of Variation:

Variance

Cont'd

- Important measure of variation
- Shows variation about the mean

- **Sample Variance**

$$S^2 = \frac{\sum_{i=1}^n (X_i - \bar{X})^2}{n-1}$$

- **Population Variance**

$$\sigma^2 = \frac{\sum_{i=1}^N (X_i - \mu)^2}{N}$$

pronounced
sigma squared

Measures of Variation:

Standard Deviation

Cont'd

- Important measure of variation
- Shows variation about the mean
- Square-root of variance
- Has the same units as the original data

□ Sample Standard Deviation

$$S = \sqrt{\frac{\sum_{i=1}^n (X_i - \bar{X})^2}{n-1}}$$

□ Population Standard Deviation

$$\sigma = \sqrt{\frac{\sum_{i=1}^N (X_i - \mu)^2}{N}}$$

pronounced
sigma

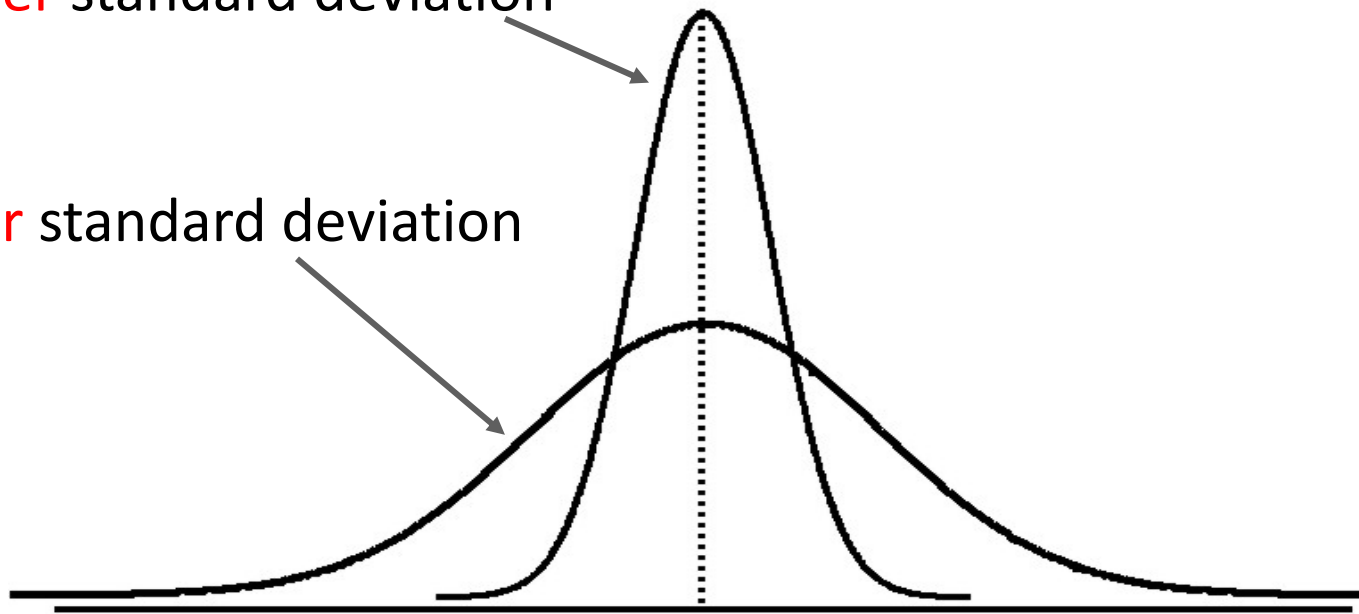
→

Measures of Variation: Standard Deviation

Cont'd

Smaller standard deviation

Larger standard deviation

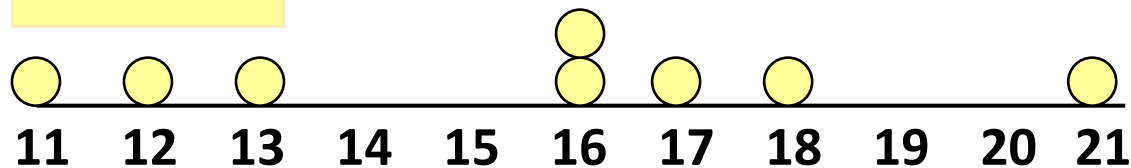


Measures of Variation: Standard Deviation

Cont'd

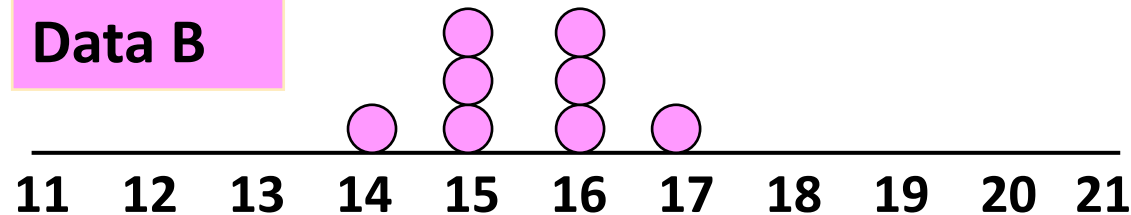
Note: All the data set are random samples from the population

Data A



$$\bar{X} = 15.5$$
$$s = 3.338$$

Data B



$$\bar{X} = 15.5$$
$$s = 0.926$$

Data C



$$\bar{X} = 15.5$$
$$s = 4.570$$

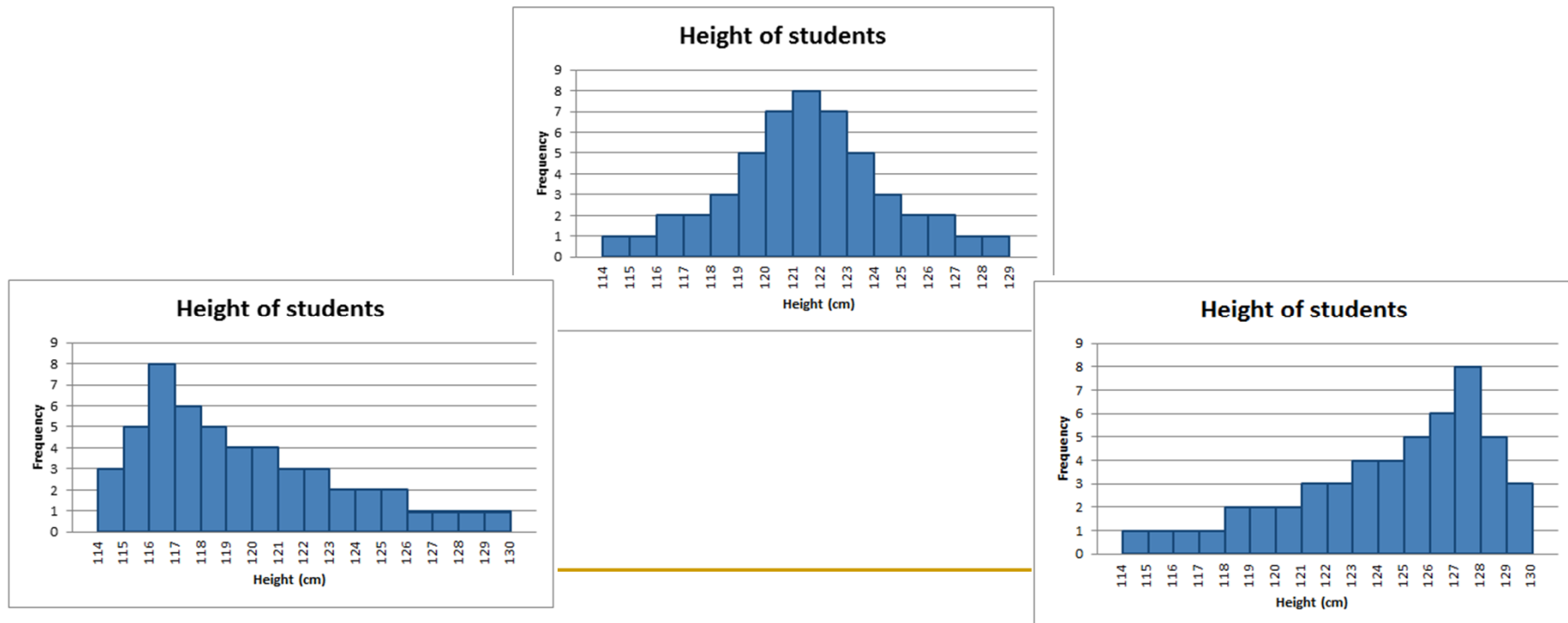
Measures of Variation

Cont'd

- The more the data are **spread out**, the **greater** the range, variance, and standard deviation
- The more the data are **concentrated**, the **smaller** the range, variance, and standard deviation
- If the values are all the same (no variation), all these measures will be **zero**
- None of these measures are **ever negative**

Distribution Shape

- Data sets may have similar central tendency measures, similar standard deviations, but different in shape



Distribution Shape

Cont'd

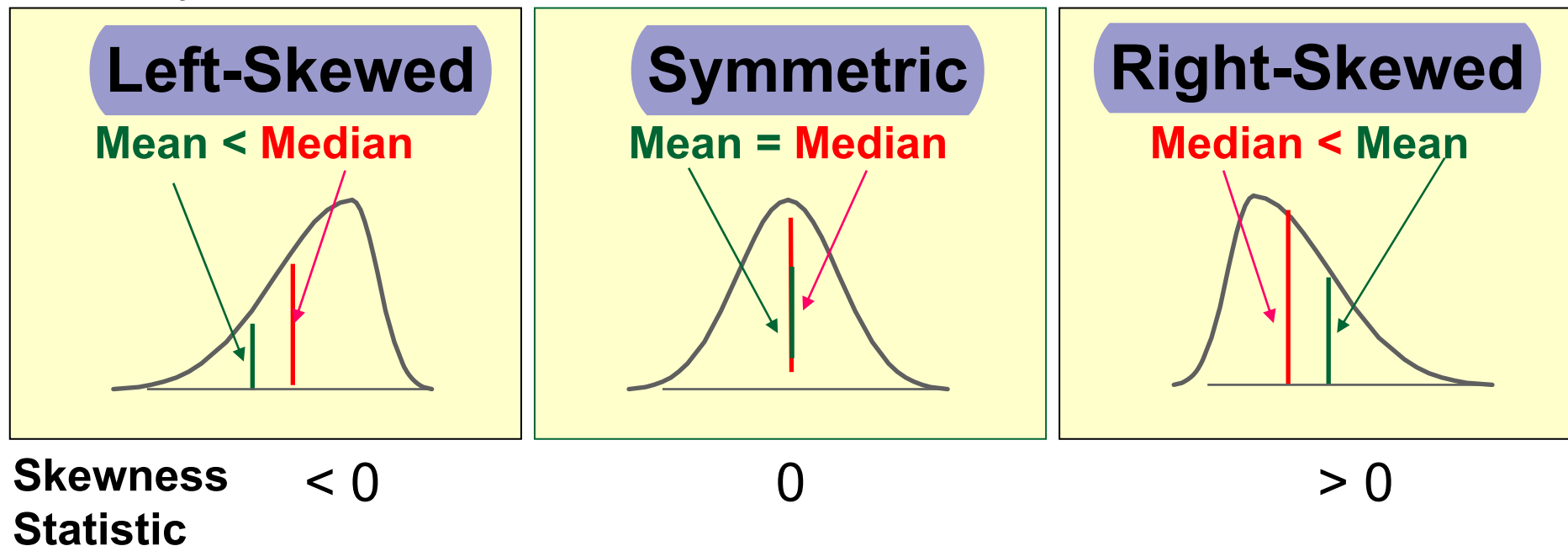
■ The Skewness

- ❑ The skewness measures the extent to which data values are not symmetrical
- ❑ The skewness equals to 0 if the distribution of the variable is symmetrical
- ❑ The skewness lesser than 0 if the distribution is left skewed, larger than 0 if the distribution is right skewed
- ❑ The skewness can help us to decide which type of central tendency is appropriate to use

Distribution Shape

Cont'd

■ Symmetric or skewed



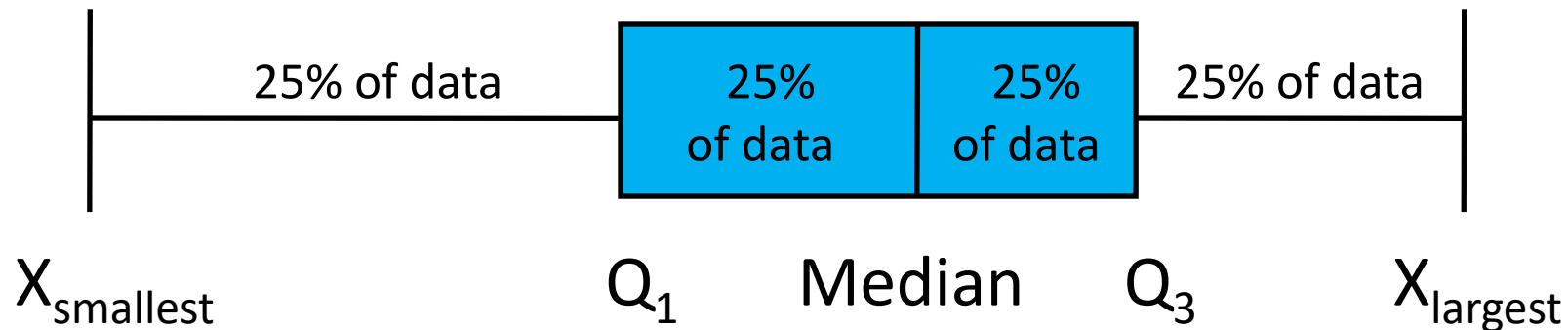
- If data are skewed, the median may be a more appropriate measure of central tendency

The Five Number Summary and Boxplot

- The five numbers that help describe the center, spread and shape of data are

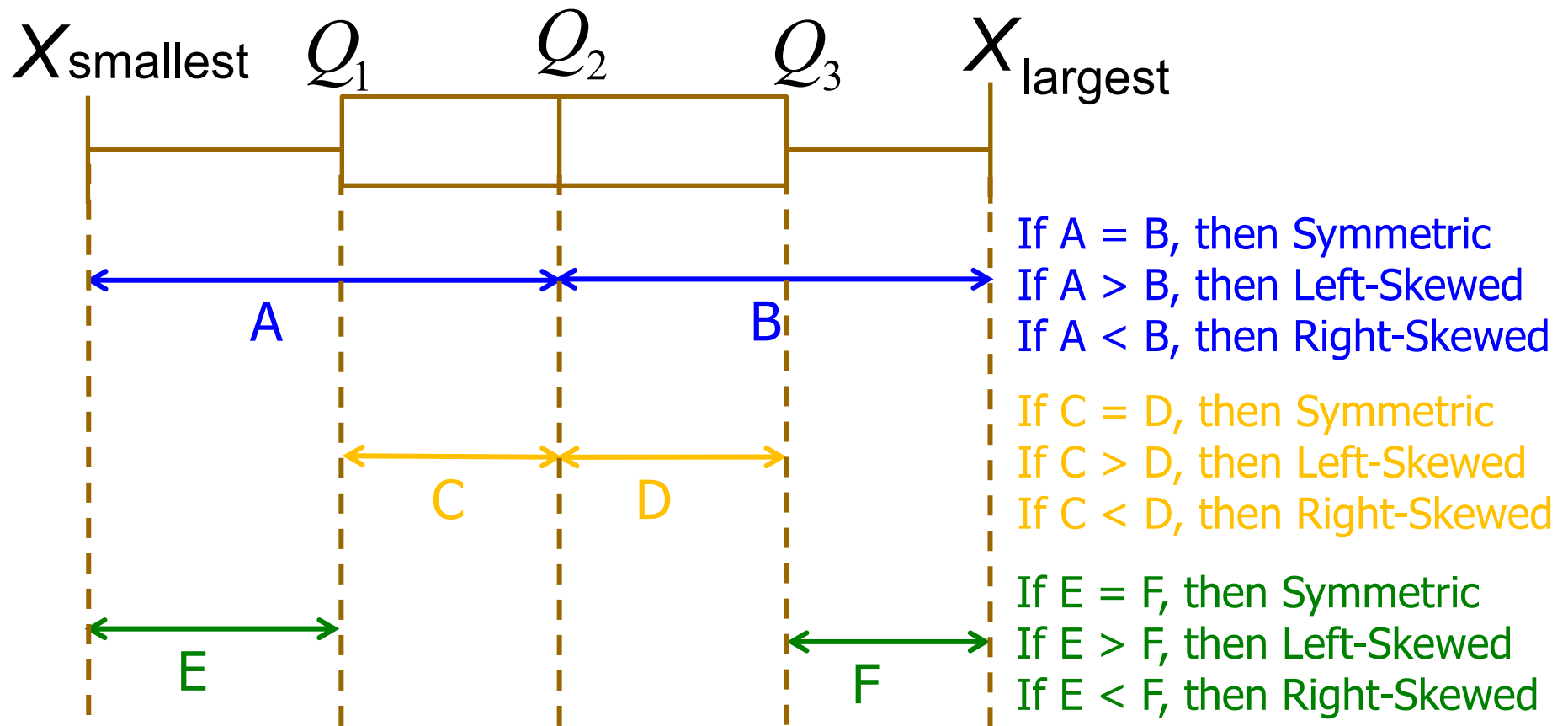
X_{smallest} -- Q_1 -- Median -- Q_3 -- X_{largest}

- Boxplot



Distribution Shape and Boxplot

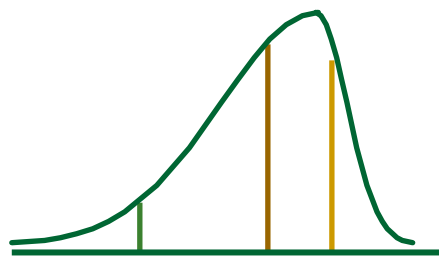
Cont'd



Distribution Shape and Boxplot

Cont'd

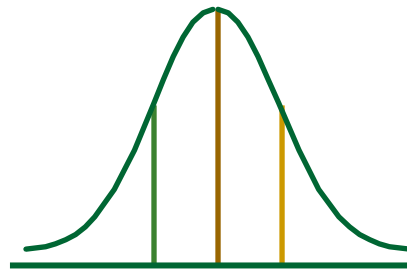
Left-Skewed



Q_1 Q_2 Q_3



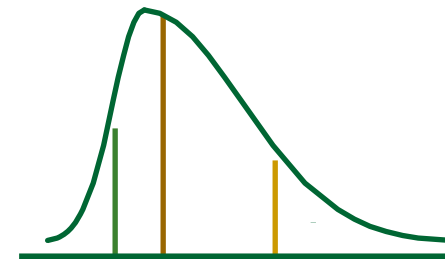
Symmetric



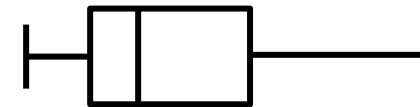
Q_1 Q_2 Q_3



Right-Skewed

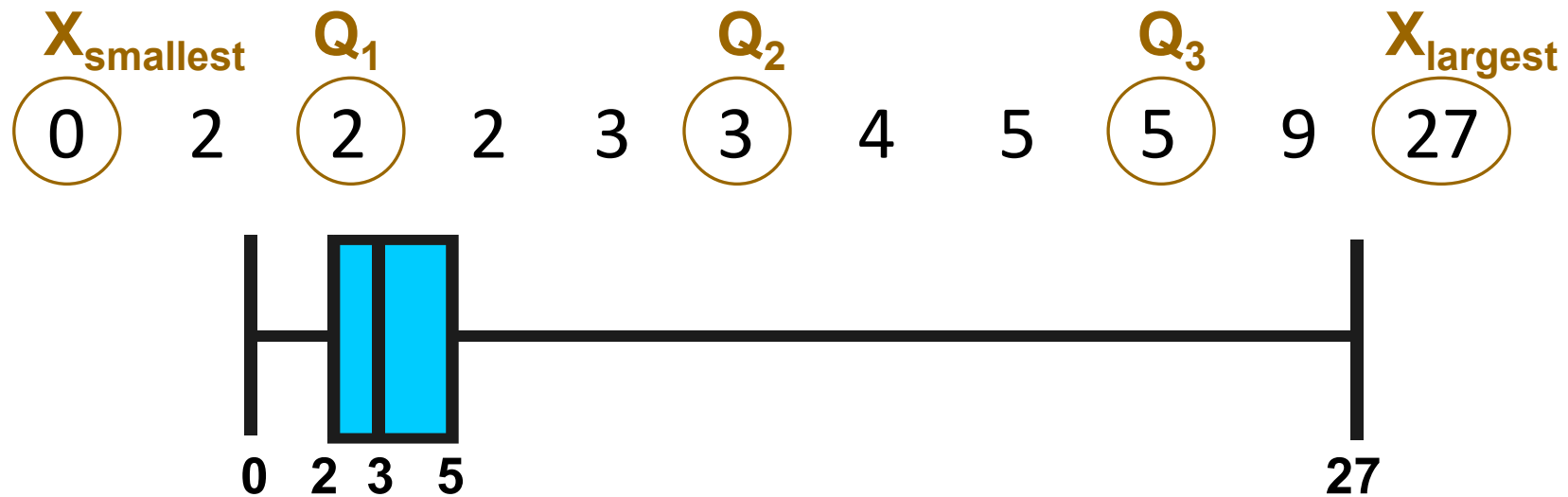


Q_1 Q_2 Q_3



Boxplot Example

Cont'd



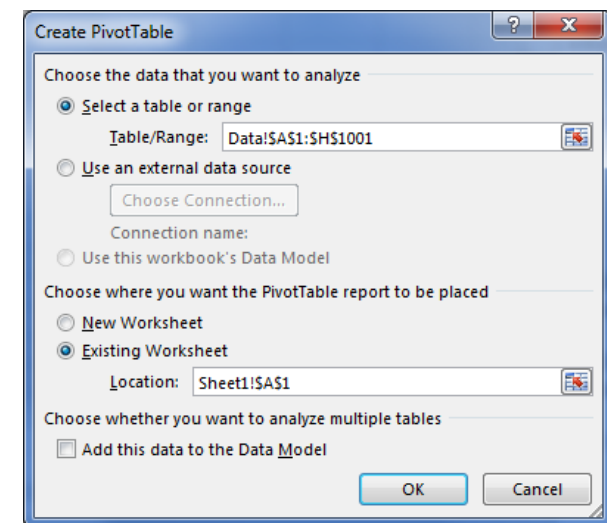
The data are right skewed, as the plot depicts

Use of Excel in Organizing Data

■ PivotTable

- ❑ PivotTable can be used to create summary table for categorical variables
- ❑ Steps to create a pivot table manually

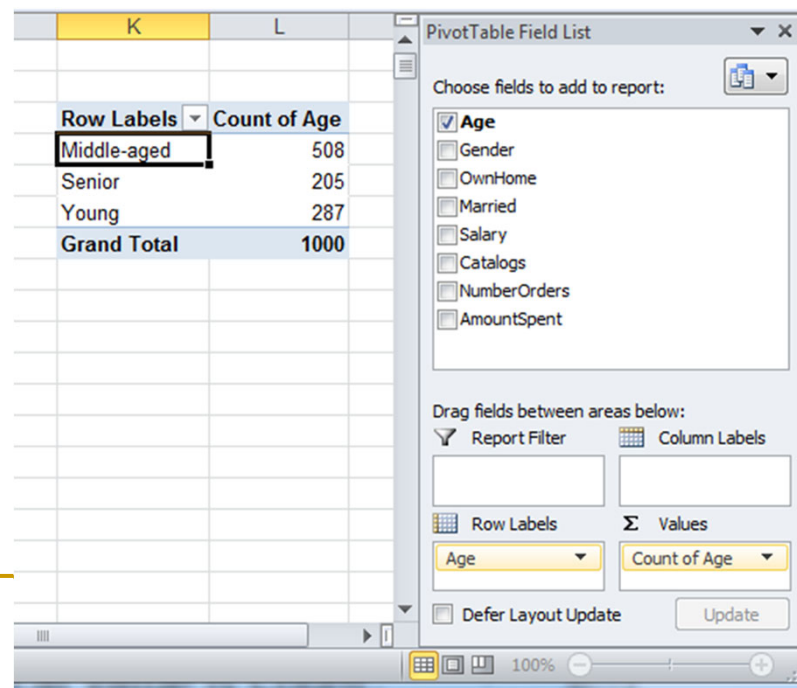
1. Click Pivot Table in Insert ribbon. In Create PivotTable dialog box, enter or confirm the address in the Table Range as the data that you want to analyze. Choose Existing worksheet and specify a location (A1, say) for the PivotTable report to be placed



Use of Excel in Organizing Data

Cont'd

2. Drag Age to the Row Labels area, and Age to Values area to create the frequency table
 - ❑ The default reported value for categorical value is Count. This can be changed from the dropdown list under Values area
 - ❑ Grand Total is included by default



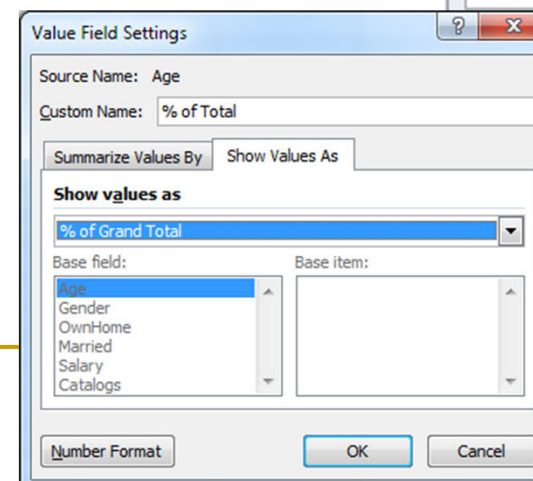
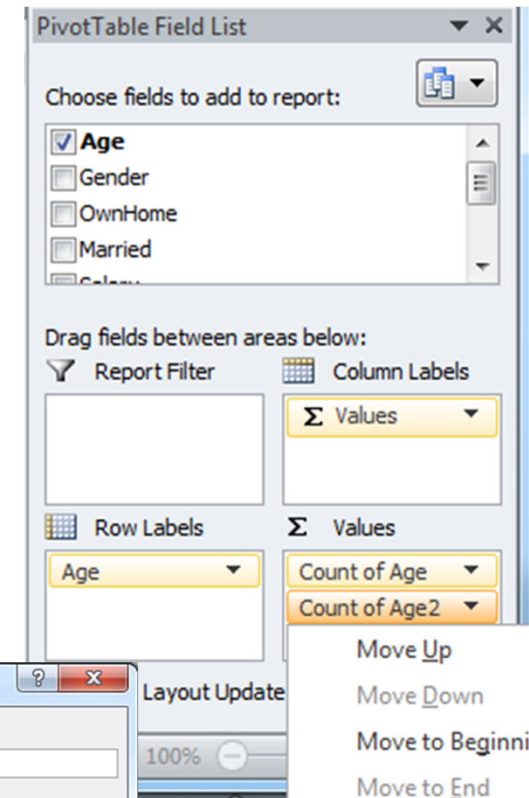
Use of Excel in Organizing Data

Cont'd

3. If the relative frequency is also wanted, drag Age into the Values field again. From the dropdown list of Count of Age 2, select Value Field Setting, then set Show Value As % of Grand Total. Enter "% of Total" into the Custom Name box

- If you do not see the Field List, click FieldList in Show group under the Analyze ribbon of PivotTable Tools

	K	L	M
Row Labels	Count of Age	% of Total	
Middle-aged	508	50.80%	
Senior	205	20.50%	
Young	287	28.70%	
Grand Total	1000	100.00%	



Calculating Descriptive Statistics in Excel

- The preparation time for the examination of 12 randomly selected students (in days):

5 21 18 9 4 17 11 28 19 2 18 22

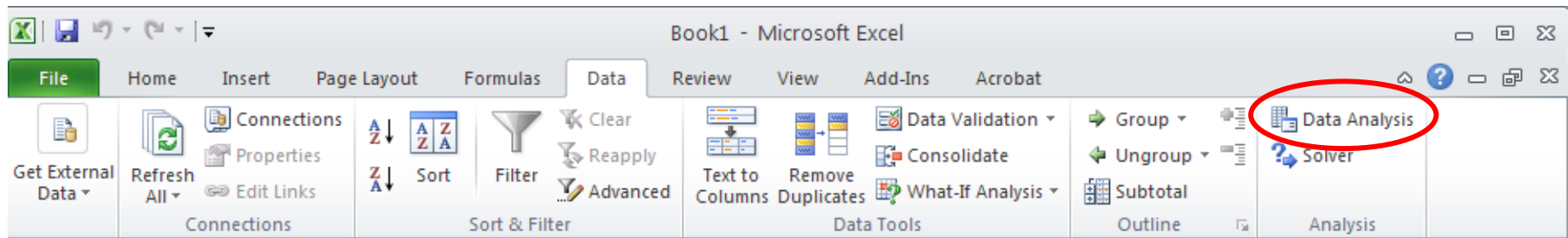
D5			f_x	=STDEV(A1:A12)
	A	B	C	D
1	5			
2	21		Mean	14.5
3	18		Median	17.5
4	9		Mode	18
5	4		Sample Standard Deviation	8.151966
6	17		Sample Variance	66.45455
7	11		Minimum	2
8	28		Maximum	28
9	19		Range	26
10	2		Sum	174
11	18		Count	12
12	22			

Mean	=average(A1:A12)
Median	=median(A1:A12)
Mode	=mode(A1:A12)
Sample Standard Deviation	=stdev.s(A1:A12)
Sample Variance	=var.s(A1:A12)
Population Standard Deviation	=stdev.p(A1:A12)
Population Variance	=var.p(A1:A12)
Maximum	=max(A1:A12)
Minimum	=min(A1:A12)
Range	=max(A1:A12)-min(A1:A12)
Sum	=sum(A1:A12)
Count	=count(A1:A12)

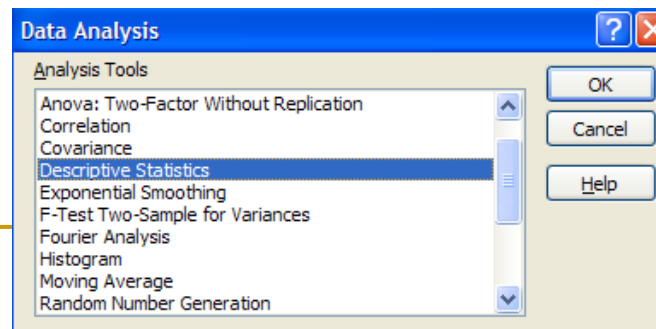
Calculating Descriptive Statistics in Excel

Cont'd

- Use of Excel “Data Analysis” Add-Ins tool to find descriptive measures
 - File → Options → Add-Ins → Click “Go” at the bottom → Check “Analysis ToolPak” and click “OK”
 - You can find “Data Analysis” in the “Data” menu bar



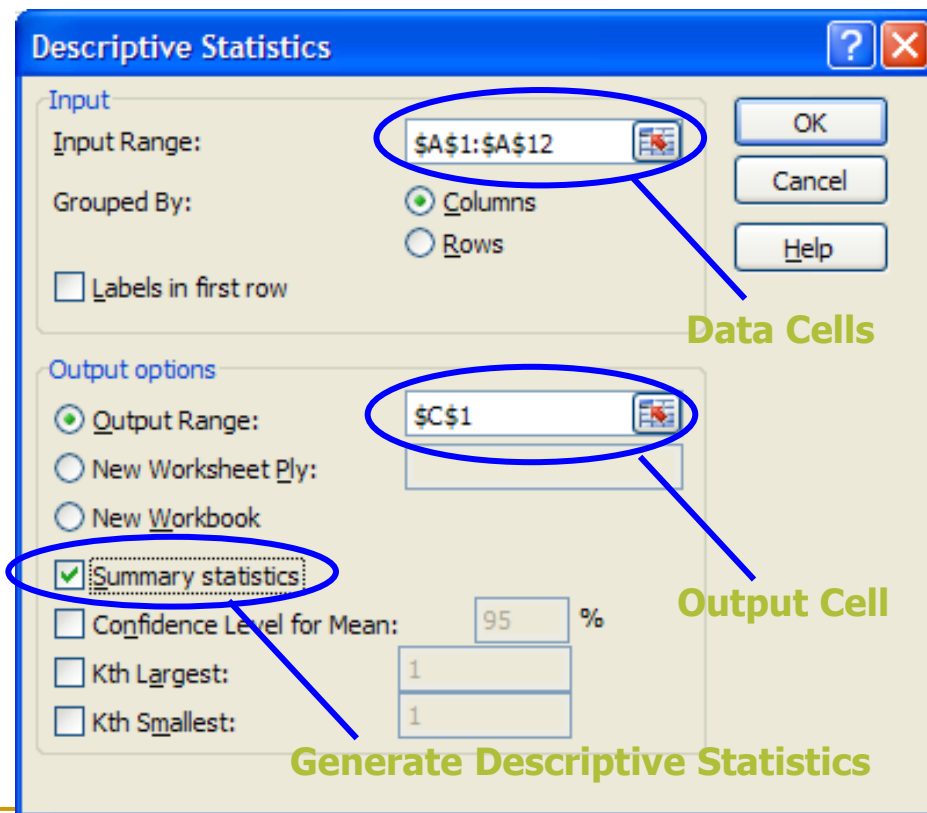
- Choose “Descriptive Statistics” at “Data Analysis” browser



Calculating Descriptive Statistics in Excel

Cont'd

- Use of Excel “Data Analysis” Add-Ins tool to find descriptive measures



C	D
Column1	
Mean	14.5
Standard Error	2.353269808
Median	17.5
Mode	18
Standard Deviation	8.151965742
Sample Variance	66.45454545
Kurtosis	-1.011390427
Skewness	-0.166736953
Range	26
Minimum	2
Maximum	28
Sum	174
Count	12

Calculating Descriptive Statistics in Calculator (For Casio fx-50F)

Date Set:

163.6 156.2 166.3 179.3 157.8 165.4 159.5 161.7 160.4

1. Change to "Lin" mode

MODE **MODE** 5 1

2. Clear previous data

SHIFT **CLR** 1 **EXE**

3. Input data

163.6 **M+** 156.2 **M+** 166.3 **M+** 179.3 **M+**
157.8 **M+** 165.4 **M+** 159.5 **M+** 161.7 **M+**
160.4 **M+**

4. Calculate descriptive statistics

Mean: **SHIFT** 2 1 1 **EXE** = 163.3555556

Population standard deviation: **SHIFT** 2 1 2 **EXE** = 6.459637417

Sample standard deviation: **SHIFT** 2 1 3 **EXE** = 6.851480132

No. of Data Input: **SHIFT** 1 3 **EXE** = 9