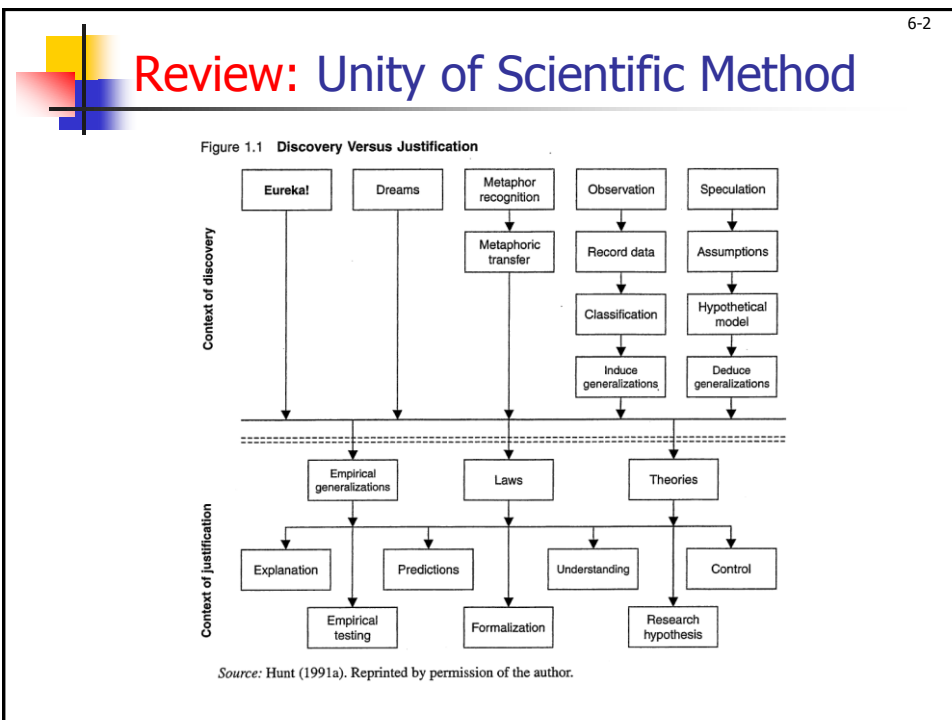
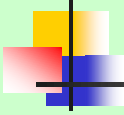


## Chapter Eleven

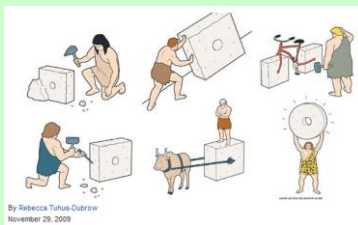
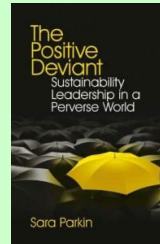
### Sampling: Design and Procedures





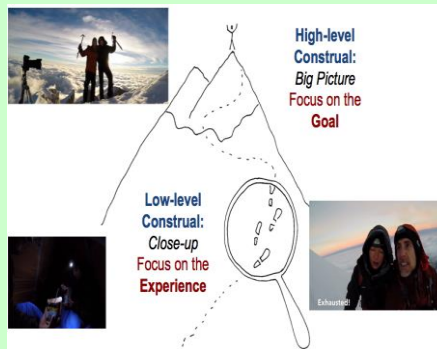
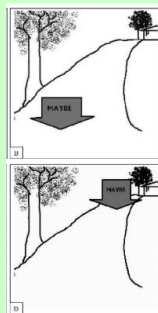
## Review: The Observation Route

Inductive 歸納: Observations to theory



## Review: The Speculation route

Deductive 演繹: General to specific (acquire knowledge)





## Review: Metaphor Recognition

Kathleen Eisenhardt (1989), 從建構個案展開理論之步驟

Step	Activity	研究流程
Getting Started	<ul style="list-style-type: none"> <li>Definition of research question</li> <li>Possibly a priori constructs</li> <li>Neither theory nor hypotheses</li> </ul>	<ul style="list-style-type: none"> <li>電信業者與行動終端設備商何以在合作關係中，潛藏競爭的態勢</li> <li>Co-petition(競合), strategic inflection point, substitute + complement, prison dilemma</li> </ul>
Selecting Cases	<ul style="list-style-type: none"> <li>Specified population</li> <li>Theoretical, not random, sampling</li> </ul>	<ul style="list-style-type: none"> <li>行動通信產業，電信營運商與行動終端設備商</li> <li>Old and New key players at the inflection point of open innovations (from Android)</li> </ul>
Crafting Instruments & Protocols	<ul style="list-style-type: none"> <li>Multiple data collection methods</li> <li>Qualitative and quantitative data combined</li> <li>Multiple investigators</li> </ul>	<ul style="list-style-type: none"> <li>質性與量化數據資料（行動數據流量、智慧手機銷售量、各國電信業者手機綁約方案）relates to your construct</li> <li>多元的資料來源，有助建立不同觀點視角，並強化證據力</li> </ul>
Entering the Field	<ul style="list-style-type: none"> <li>Overlap data collection and analysis, including field notes</li> <li>Flexible and opportunistic data collection methods</li> </ul>	<ul style="list-style-type: none"> <li>資料蒐集與分析併行：從初始資料蒐集並分析（互補品關係），而後延續形成下階段資料蒐集方向（競合成因：補貼）</li> <li>依浮現的命題概念，彈性調整資料蒐集方式，以幫助理論形成</li> </ul>
Analyzing Data	<ul style="list-style-type: none"> <li>Within-case analysis</li> <li>Cross-case pattern search using divergent techniques</li> </ul>	<ul style="list-style-type: none"> <li>綜合相關資料與產業現象，著手探討行動通信產業個案</li> <li>跨個案與個人電腦 PC 產業相較，尋求相似的態樣分析，以了解隱藏於個案表象下的成因 (take away)</li> </ul>



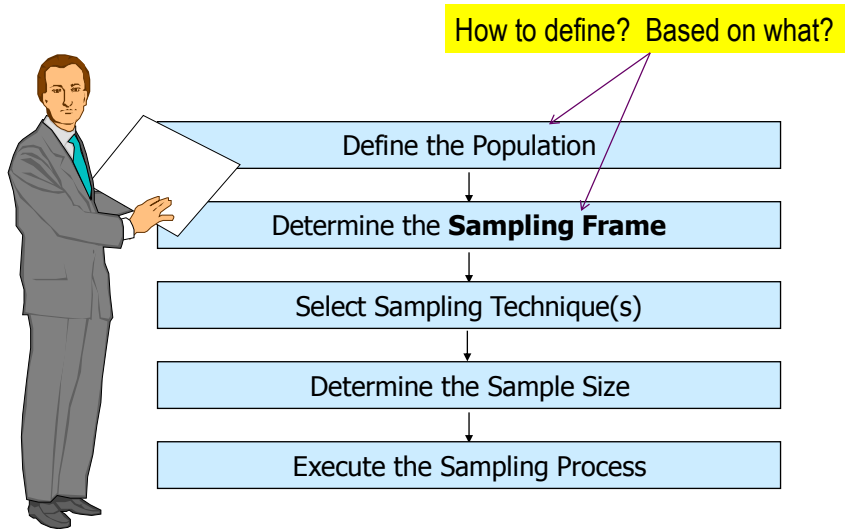
## Sample vs. Census

Table 11.1

Type of Study	Conditions Favoring the Use of	
	Sample	Census
1. Budget	Small	Large
2. Time available	Short	Long
3. Population size	Large	Small
4. Variance in the characteristic	Small	Large
5. Cost of sampling errors	Low	High
6. Cost of nonsampling errors	High	Low
7. Nature of measurement	Destructive	Nondestructive
8. Attention to individual cases	Yes	No

## The Sampling Design Process

Fig. 11.1



## Define the Target Population

The target population is the collection of elements or objects that possess the information sought by the researcher and about which inferences are to be made. The target population should be defined in terms of elements, sampling units, extent, and time.

- An **element** is the object about which or from which the information is desired, e.g., the respondent.
- A **sampling unit** is an element, or a unit containing the element, that is available for selection at some stage of the sampling process.
- **Extent** refers to the geographical boundaries.
- **Time** is the time period under consideration.

Example: Procumbent survey on B2B customer, 臍帶血的  
認知學習, social networking, smart Grid, FinTech, structural holes



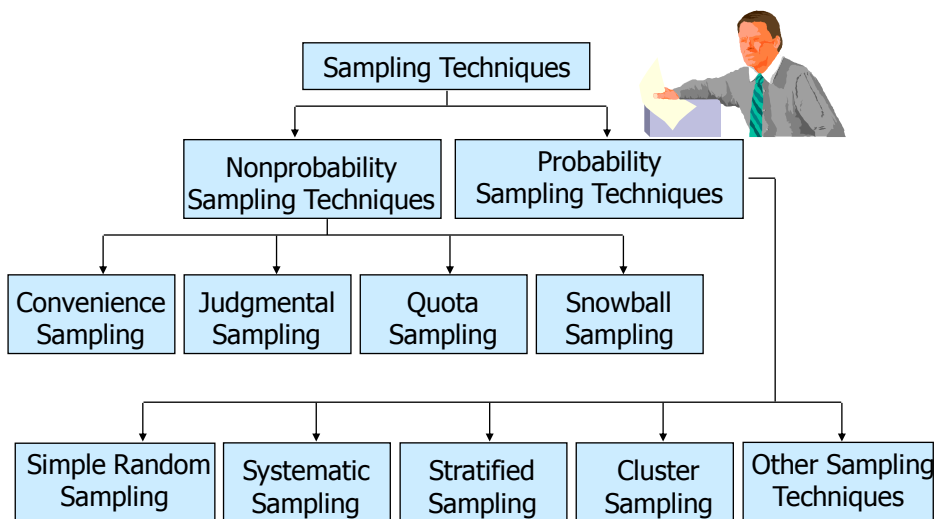
## Define the Target Population

<b>(element)</b>	purchasing agents of
<b>(sampling unit)</b>	multinational companies in China that have
<b>(extent)</b>	bought any of our products
<b>(time)</b>	In the last two years



## Classification of Sampling Techniques

Fig. 11.2





## Convenience Sampling

**Convenience sampling** attempts to obtain a sample of convenient elements. Often, respondents are selected because they happen to be in the right place at the right time.

- use of students, and members of social organizations
- mall intercept interviews without qualifying the respondents
- department stores using charge account lists
- “people on the street” interviews



## Judgmental Sampling

**Judgmental sampling** is a form of convenience sampling in which the population elements are selected based on the judgment of the researcher.

- test markets
- purchase engineers selected in industrial marketing research
- bellwether precincts selected in voting behavior research
- expert witnesses used in court
- IMM, OB, and OT journals?



## Quota Sampling

**Quota sampling** may be viewed as two-stage restricted judgmental sampling.

- The first stage consists of developing control categories, or quotas, of population elements.
- In the second stage, sample elements are selected based on convenience or judgment.

	<u>Population composition</u>	<u>Sample composition</u>
Control		
Characteristic	Percentage	Number
Sex		
Male	48	480
Female	52	520
	<hr/> 100	<hr/> 1000



## Snowball Sampling

In **snowball sampling**, an initial group of respondents is selected, usually at random.

- After being interviewed, these respondents are asked to identify others who belong to the target population of interest.
- Subsequent respondents are selected based on the referrals.
- Examples: social study on minority groups such as gay, delinquent, foreign workers



## Simple Random Sampling 簡單隨機抽樣

- Each element in the population has a known and equal probability of selection.
- Each possible sample of a given size ( $n$ ) has a known and equal probability of being the sample actually selected.
- This implies that every element is selected independently of every other element.
- Thesis Examples: TaiPower households list? primary school teachers or students? Online bloggers?



## Systematic Sampling

- The sample is chosen by selecting a random starting point and then picking every  $i$ th element in succession from the sampling frame.
- The sampling interval,  $i$ , is determined by dividing the population size  $N$  by the sample size  $n$  and rounding to the nearest integer.
- When the ordering of the elements is related to the characteristic of interest, systematic sampling increases the representativeness of the sample.
- If the ordering of the elements produces a cyclical pattern, systematic sampling may decrease the representativeness of the sample.  
 For example, there are 100,000 elements in the population and a sample of 1,000 is desired. In this case the sampling interval,  $i$ , is 100. A random number between 1 and 100 is selected. If, for example, this number is 23, the sample consists of elements 23, 123, 223, 323, 423, 523, and so on.
- Example: IRA tax payers list? Google log file?





## Stratified Sampling

- A two-step process in which the population is partitioned into subpopulations, or strata.
- The strata should be mutually exclusive and collectively exhaustive in that every population element should be assigned to one and only one stratum and no population elements should be omitted.
- Next, elements are selected from each stratum by a random procedure, usually SRS.
- A major objective of stratified sampling is to increase precision without increasing cost.



## Stratified Sampling

- The elements within a stratum should be as **homogeneous** as possible, but the elements in different strata should be as heterogeneous as possible.
- The stratification variables should also be closely related to the characteristic of interest.
- Finally, the variables should decrease the cost of the stratification process by being easy to measure and apply.
- In proportionate stratified sampling, the size of the sample drawn from each stratum is proportionate to the relative size of that stratum in the total population.
- In disproportionate stratified sampling, the size of the sample from each stratum is proportionate to the relative size of that stratum and to the standard deviation of the distribution of the characteristic of interest among all the elements in that stratum.

## Stratified Sampling



- Low variance in each stratum
- 1, 1, 1, 6, 6, 6, 6, 4, 4, 4
- $\binom{14}{7} = 3432 > \binom{3}{1}\binom{5}{2}\binom{6}{4} = 450$
- Example, accounting systems such as months, types of accounts, locations, activities, operations etc.
- Example, OEM, ODM, EMS, software, network, service, etc. (assumption?)

## Stratified Sampling

	Mean	Variance	Standard Deviation
Without Stratification	22.6	8.3	2.88
Within Stratum 1	25.5	0.5	0.71
Within Stratum	20.7	2.35	1.53



## Stratified Sampling (Optimal Allocation)

National store chains	25%		47%
Large independent store	12%		
Medium independent store	33%		26%
Small independent store	30%		19%
			8%

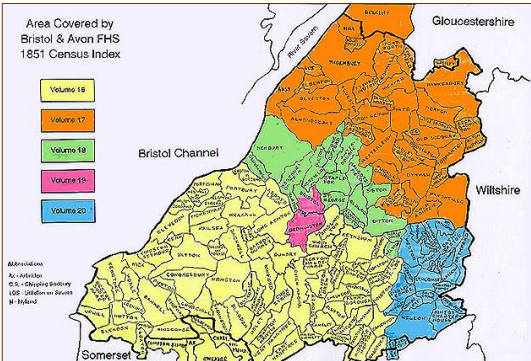


## Cluster Sampling

- The target population is first divided into mutually exclusive and collectively exhaustive subpopulations, or clusters.
- Then a random sample of clusters is selected, based on a probability sampling technique such as SRS.
- For each selected cluster, either all the elements are included in the sample (one-stage) or a sample of elements is drawn probabilistically (two-stage).
- Elements within a cluster should be as **heterogeneous** as possible, but clusters themselves should be as homogeneous as possible. Ideally, each cluster should be a small-scale representation of the population.
- In **probability proportionate to size sampling**, the clusters are sampled with probability proportional to size. In the second stage, the probability of selecting a sampling unit in a selected cluster varies inversely with the size of the cluster.



# Area Sampling



$$\binom{M}{m} \binom{N_k}{n_k} m$$

One-stage: N=n

Two-stage or multiple-level stage

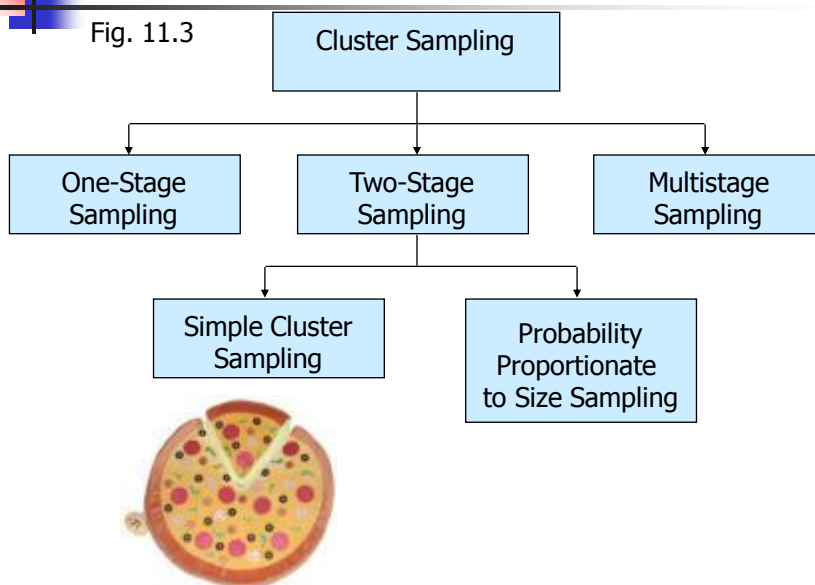


# Examples on kinds of Clusters

Adult Population in China	Cities, Hsiens, census tracts, street blocks, households
Manufacturing firms	Regions, cities, plants
Airline travelers	Airports, planes
Hospital patients	?

## Types of Cluster Sampling


Fig. 11.3



## Strengths and Weaknesses of Basic Sampling Techniques

Table 11.3

Technique	Strengths	Weaknesses
<b>Nonprobability Sampling</b>		
Convenience sampling	Least expensive, least time-consuming, most convenient	Selection bias, sample not representative, not recommended for descriptive or causal research
Judgmental sampling	Low cost, convenient, not time-consuming	Does not allow generalization, subjective
Quota sampling	Sample can be controlled for certain characteristics	Selection bias, no assurance of representativeness
Snowball sampling	Can estimate rare characteristics	Time-consuming
<b>Probability sampling</b>		
Simple random sampling (SRS)	Easily understood, results projectable	Difficult to construct sampling frame, expensive, lower precision, no assurance of representativeness.
Systematic sampling	Can increase representativeness, easier to implement than SRS, sampling frame not necessary	Can decrease representativeness
Stratified sampling	Include all important subpopulations, precision	Difficult to select relevant stratification variables, not feasible to stratify on many variables, expensive
Cluster sampling	Easy to implement, cost effective	Imprecise, difficult to compute and interpret results

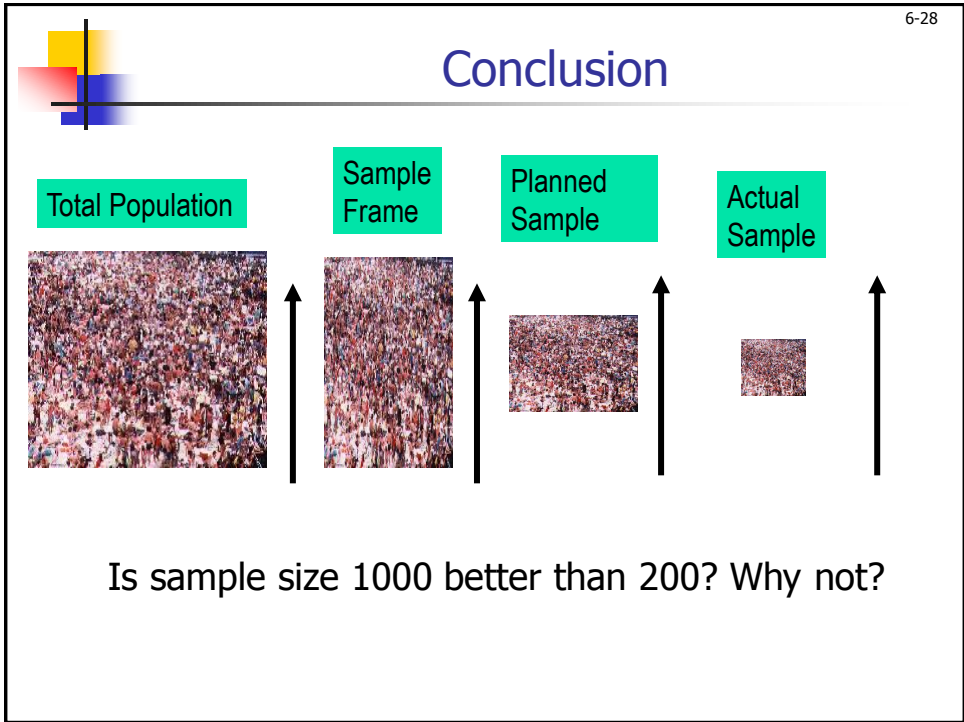


## Choosing Nonprobability vs. Probability Sampling

6-27

Table 11.4 cont.

Factors	Conditions Favoring the Use of	
	Nonprobability sampling	Probability sampling
Nature of research	Exploratory	Conclusive
Relative magnitude of sampling and nonsampling errors	Nonsampling errors are larger	Sampling errors are larger
Variability in the population	Homogeneous (low)	Heterogeneous (high)
Statistical considerations	Unfavorable	Favorable
Operational considerations	Favorable	Unfavorable






## Chapter Twelve

---

### Sampling: Final and Initial Sample Size Determination



## Definitions and Symbols

6-30

- **Parameter:** A **parameter** is a summary description of a fixed characteristic or measure of the target **population**. A parameter denotes the true value which would be obtained if a census rather than a sample was undertaken.
- **Statistic:** A **statistic** is a summary description of a characteristic or measure of the **sample**. The sample statistic is used as an estimate of the population parameter.
- **Finite Population Correction:** The **finite population correction** (fpc) is a correction for overestimation of the variance of a population parameter, e.g., a mean or proportion, when the sample size is 10% or more of the population size.



## Definitions and Symbols

- **Precision level:** When estimating a population parameter by using a sample statistic, the **precision level** is the desired size of the estimating interval. This is the maximum permissible difference between the sample statistic and the population parameter.
- **Confidence interval:** The **confidence interval** 信賴區間 is the range into which the true population parameter will fall, assuming a given level of confidence.
- **Confidence level:** The **confidence level** is the probability that a confidence interval will include the population parameter.



## Symbols for Population and Sample Variables

Table 12.1

Variable	Population	Sample
Mean	$\mu$	$\bar{x}$
Proportion	$\Pi$	$p$
Variance	$\sigma^2$	$s^2$
Standard deviation	$\sigma$	$s$
Size	$N$	$n$
Standard error of the mean	$\sigma_{\bar{x}}$	$S_{\bar{x}}$
Standard error of the proportion	$\sigma_p$	$S_p$
Standardized variate (z)	$(X-\mu)/\sigma$	$(\bar{x}-\mu)/S$
Coefficient of variation (C)	$\sigma/\mu$	$S/\bar{x}$



## Mean, Variance, and S.D.

- The **mean**,

$$\bar{X} = \sum_{i=1}^n X_i / n$$

Where,

$X_i$  = Observed values of the variable  $X$

$n$  = Number of observations (sample size)

- The **variance** is the mean squared deviation from the mean. The variance can never be negative.
- The **standard deviation** is the square root of the variance.

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

## Sampling distribution of the sample means

- Suppose that in a company the retirement fund is invested in five corporate stocks with the following returns:

Stock.....	Return
A.....	7%
B.....	12%
C.....	-3%
D.....	21%
E.....	3%

In this example, the **population** mean  $\mu$  is equal to 8%, and the population standard deviation  $\sigma$  is equal to 8.15%.



## Sampling distribution of the sample means

- Suppose we take a **random sample** of three stocks, there are ten possibilities  $C(5_3)$
- Sample Stocks.....Returns.....Mean
- |                 |                    |        |
|-----------------|--------------------|--------|
| 1) A, B, C..... | 7%..12%..-3%.....  | 5.33%  |
| 2) A, B, D..... | 7%..12%..21%.....  | 13.33% |
| 3) A, B, E..... | 7%..12%..3%.....   | 7.33%  |
| 4) A, C, D..... | 7%..-3%..21%.....  | 8.33%  |
| 5) A, C, E..... | 7%..-3%..3%.....   | 2.33%  |
| 6) A, D, E..... | 7%..21%..3%.....   | 10.33% |
| 7) B, C, D..... | 12%..-3%..21%..... | 10.00% |
| 8) B, C, E..... | 12%..-3%..3%.....  | 4.00%  |
| 9) B, D, E..... | 12%..21%..3%.....  | 12.00% |
| 0) C, D, E..... | -3%..21%..3%.....  | 7.00%  |

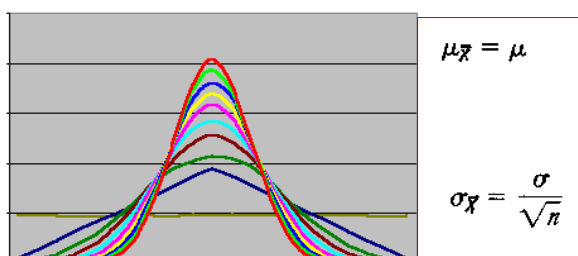


## Sampling distribution of the sample means

- If a population is *normally distributed*, then:
- The mean of the sampling distribution of means equals the population mean.
- 2. The standard deviation of the sampling distribution of means (or standard error of the mean) is smaller than the population standard deviation.

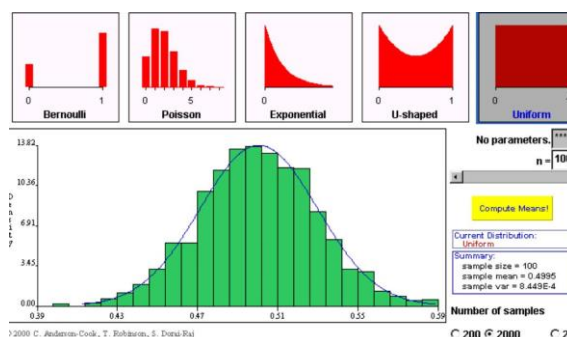
## Central Limit Theorem

- Given a distribution with a mean  $\mu$  and variance  $\sigma^2$ , the sampling distribution of the mean approaches a normal distribution with a mean ( $\mu$ ) and a variance  $\sigma^2/N$  as  $N$  (the sample size) increases.
- $\bar{X} \sim \text{N.I.D}(\mu, \sigma^2/N)$



## Central Limit Theorem

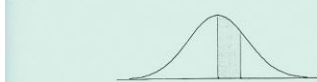
No matter what the shape of the original distribution, the sampling distribution of the mean approaches a normal distribution (不論母體的分佈如何，其平均數的分佈都會傾向常態分佈). Furthermore, for most distributions, a normal distribution is approached very quickly as  $N$  increases.





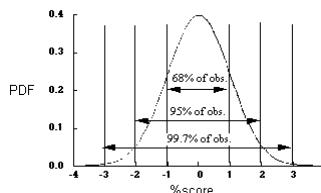
# Standard Normal Distribution

Appendix 4 - Values of the Standard Normal distribution



The table gives the probability that a standard Normal variable lies between 0 and x (which is equivalent to the shaded area on the figure).

z	0.00	0.01	0.02	0.03	0.04	0.05	0.06	0.07	0.08	0.09
0.0	0.0000	0.0040	0.0080	0.0120	0.0160	0.0199	0.0239	0.0279	0.0319	0.0359
0.1	0.0398	0.0438	0.0478	0.0517	0.0557	0.0596	0.0636	0.0675	0.0714	0.0754
0.2	0.0793	0.0832	0.0871	0.0910	0.0948	0.0987	0.1025	0.1064	0.1103	0.1141
0.3	0.1179	0.1217	0.1255	0.1293	0.1331	0.1368	0.1406	0.1443	0.1480	0.1517
0.4	0.1554	0.1591	0.1628	0.1664	0.1700	0.1736	0.1772	0.1808	0.1844	0.1879
0.5	0.1915	0.1950	0.1985	0.2019	0.2054	0.2088	0.2123	0.2157	0.2190	0.2224
0.6	0.2258	0.2291	0.2324	0.2357	0.2389	0.2422	0.2454	0.2486	0.2518	0.2549
0.7	0.2580	0.2612	0.2643	0.2673	0.2703	0.2734	0.2764	0.2794	0.2823	0.2852
0.8	0.2881	0.2910	0.2939	0.2967	0.2996	0.3023	0.3051	0.3078	0.3106	0.3133
0.9	0.3159	0.3186	0.3212	0.3238	0.3264	0.3289	0.3315	0.3340	0.3365	0.3389
1.0	0.3413	0.3438	0.3461	0.3485	0.3508	0.3531	0.3554	0.3577	0.3599	0.3621
1.1	0.3643	0.3665	0.3686	0.3708	0.3729	0.3749	0.3770	0.3790	0.3810	0.3829
1.2	0.3849	0.3869	0.3888	0.3907	0.3925	0.3944	0.3962	0.3980	0.3997	0.4015
1.3	0.4032	0.4049	0.4066	0.4082	0.4099	0.4115	0.4131	0.4147	0.4162	0.4177
1.4	0.4192	0.4207	0.4222	0.4236	0.4251	0.4265	0.4279	0.4292	0.4306	0.4319
1.5	0.4332	0.4345	0.4357	0.4370	0.4382	0.4394	0.4406	0.4418	0.4429	0.4441
1.6	0.4452	0.4463	0.4474	0.4484	0.4495	0.4505	0.4515	0.4525	0.4535	0.4545
1.7	0.4554	0.4564	0.4573	0.4582	0.4591	0.4599	0.4608	0.4616	0.4625	0.4633
1.8	0.4641	0.4648	0.4656	0.4664	0.4671	0.4678	0.4685	0.4693	0.4699	0.4706
1.9	0.4713	0.4719	0.4726	0.4732	0.4738	0.4744	0.4750	0.4756	0.4761	0.4767
2.0	0.4772	0.4778	0.4783	0.4788	0.4793	0.4798	0.4803	0.4808	0.4812	0.4817
2.1	0.4821	0.4826	0.4830	0.4834	0.4838	0.4842	0.4846	0.4850	0.4854	0.4857
2.2	0.4861	0.4864	0.4868	0.4871	0.4875	0.4878	0.4881	0.4884	0.4887	0.4890
2.3	0.4893	0.4896	0.4898	0.4901	0.4904	0.4906	0.4908	0.4911	0.4913	0.4916
2.4	0.4918	0.4920	0.4922	0.4925	0.4927	0.4929	0.4931	0.4932	0.4934	0.4936
2.5	0.4938	0.4940	0.4941	0.4943	0.4945	0.4946	0.4948	0.4949	0.4951	0.4952
2.6	0.4953	0.4955	0.4956	0.4957	0.4959	0.4960	0.4961	0.4962	0.4963	0.4964
2.7	0.4965	0.4966	0.4967	0.4968	0.4969	0.4970	0.4971	0.4972	0.4973	0.4974
2.8	0.4974	0.4975	0.4976	0.4977	0.4977	0.4978	0.4979	0.4979	0.4980	0.4981
2.9	0.4981	0.4982	0.4982	0.4983	0.4984	0.4984	0.4985	0.4985	0.4986	0.4986
3.0	0.4987	0.4987	0.4988	0.4988	0.4989	0.4989	0.4989	0.4989	0.4989	0.4990



$$Z = (X - \mu) / \sigma$$

Z score ~ N.I.D (0, 1)

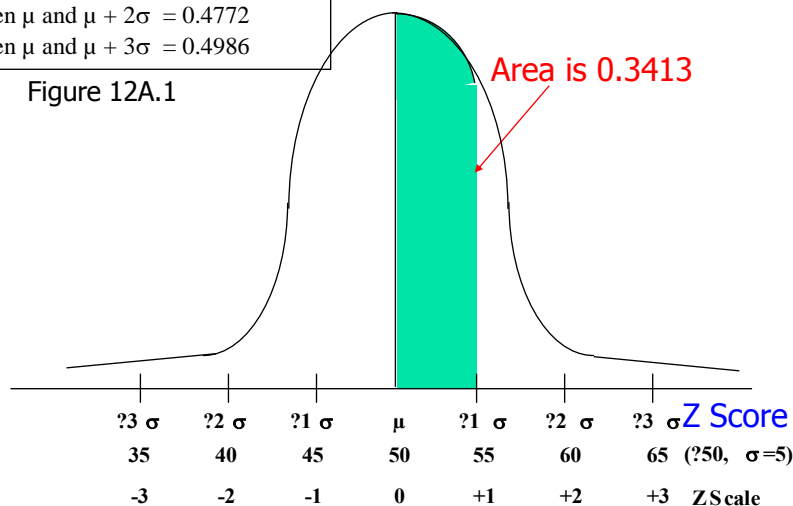
$$P(x) = \frac{1}{\sigma \sqrt{2\pi}} e^{-(x-\mu)^2 / (2\sigma^2)}$$



# Standard Normal Distribution

Area between  $\mu$  and  $\mu + 1\sigma = 0.3431$   
 Area between  $\mu$  and  $\mu + 2\sigma = 0.4772$   
 Area between  $\mu$  and  $\mu + 3\sigma = 0.4986$

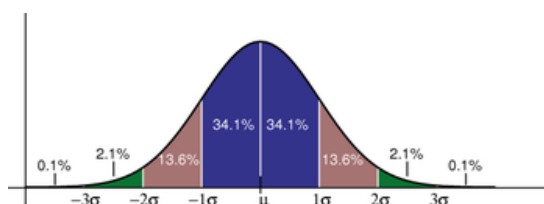
Figure 12A.1



## Application of CLT: Z test

- CLT states that sample means are normally distributed regardless of the shape of the population for large samples and for any sample size with normally distributed population, thus sample means can be analyzed by using Z scores:  

$$Z = (X - \text{Mean}) / \text{Standard deviation} = (X - \mu) / \sigma$$
- If sample means are normally distributed, the Z score equation applied to sample means would be:



## Application of CLT: Confidence Interval

The confidence interval is given by  $\bar{X} \pm z\sigma_{\bar{x}}$   
 Sample mean  $\bar{X}$  is known. We use CLT to estimate the population mean  $\mu$

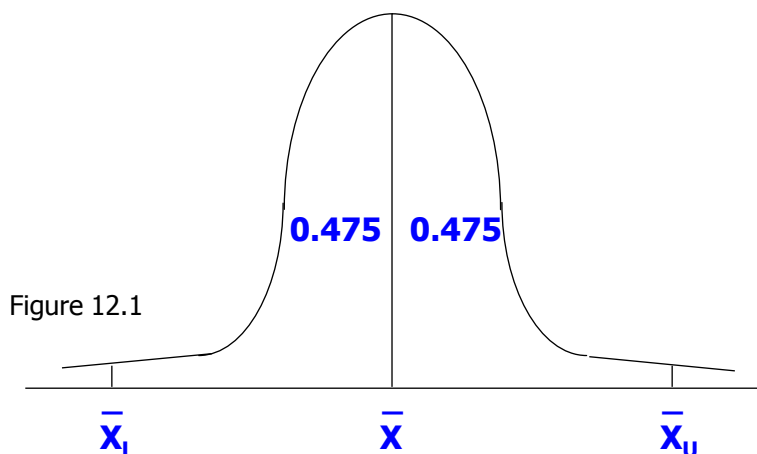
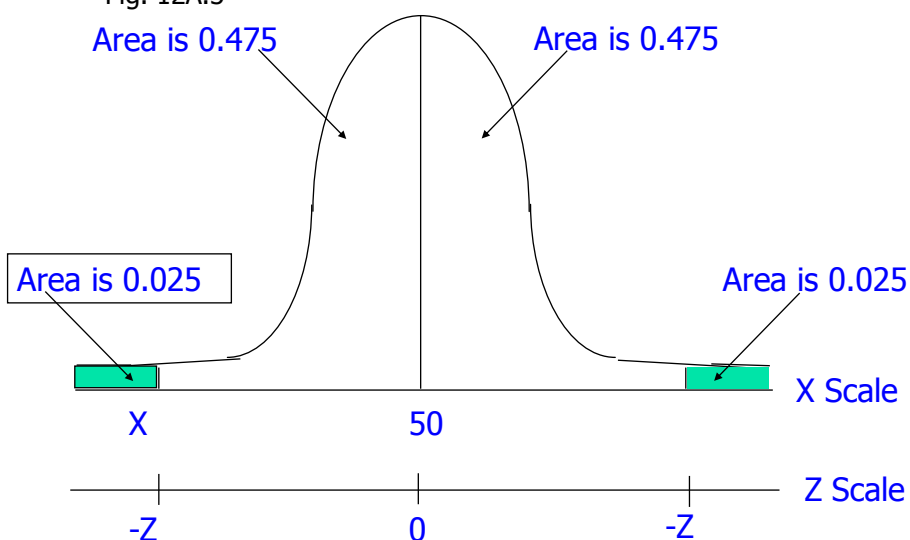


Figure 12.1

## Application of CLT: Confidence Interval

Fig. 12A.3

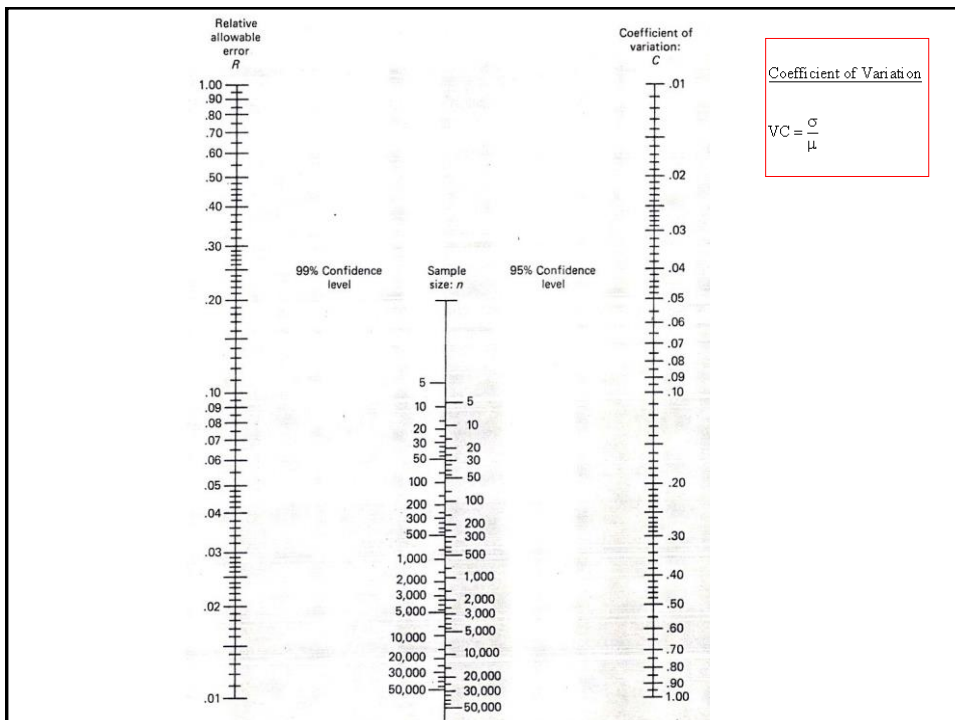


## Application of CLT: Sample Size

- Recall precision is the difference between the sample statistic and the population parameter,  $D = \bar{X} - \mu$

- Re-arrange z score equation. We get  


$$n = z^2 \sigma^2 / D^2$$



6-46

## Application of CLT: Sample Size

- If sample size  $n$  is greater than 10% of the population size  $N$ . i.e.,  $n/N > 0.10$ , then  $z$  is
 
$$Z = \frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}} \sqrt{\frac{N-n}{N-1}}}$$
- Z score for Sampling Distribution of Sample Proportion is:
 
$$Z = \frac{\hat{p} - P}{\sqrt{\frac{P \cdot Q}{n}}}$$

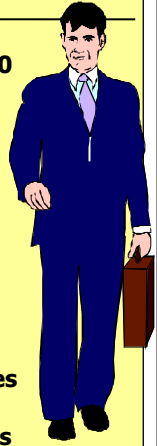


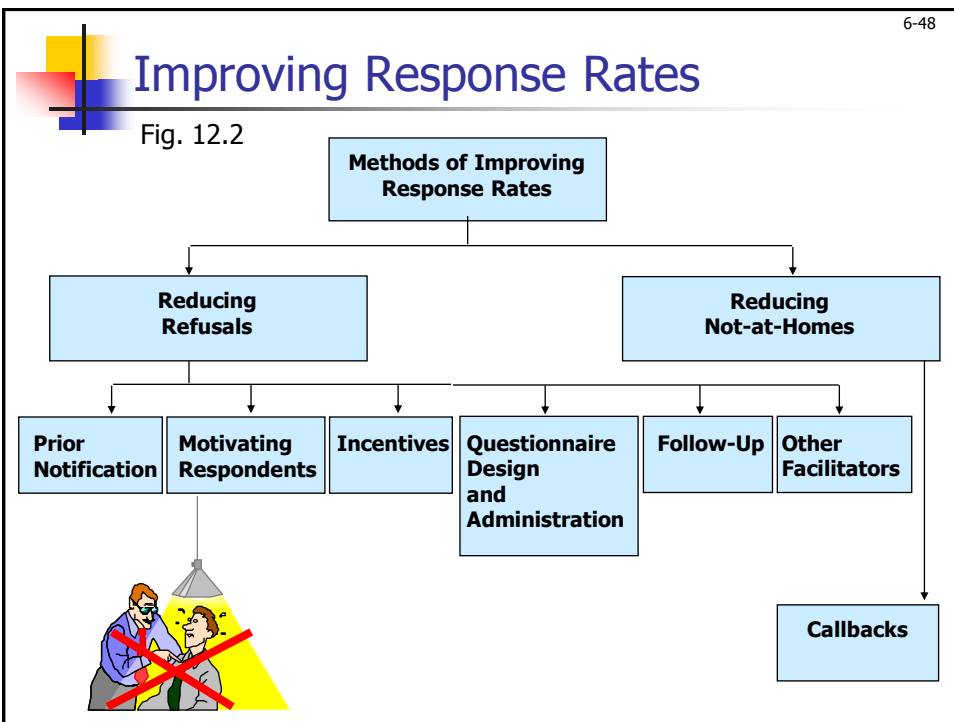
## Sample Sizes Used in Marketing Research Studies

6-47

Table 11.2

Type of Study	Minimum Size	Typical Range
Problem identification research (e.g. market potential)	500	1,000-2,500
Problem-solving research (e.g. pricing)	200	300-500
Product tests	200	300-500
Test marketing studies	200	300-500
TV, radio, or print advertising (per commercial or ad tested)	150	200-300
Test-market audits	10 stores	10-20 stores
Focus groups	2 groups	4-12 groups









## Adjusting for Nonresponse

- **Subsampling of Nonrespondents** – the researcher contacts a subsample of the nonrespondents, usually by means of telephone or personal interviews.
- In **replacement**, the nonrespondents in the current survey are replaced with nonrespondents from an earlier, similar survey. The researcher attempts to contact these nonrespondents from the earlier survey and administer the current survey questionnaire to them, possibly by offering a suitable incentive.



## Adjusting for Nonresponse

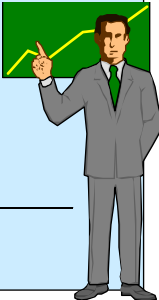
- In **substitution**, the researcher substitutes for nonrespondents other elements from the sampling frame that are expected to respond. The sampling frame is divided into subgroups that are internally homogeneous in terms of respondent characteristics but heterogeneous in terms of response rates. These subgroups are then used to identify substitutes who are similar to particular nonrespondents but dissimilar to respondents already in the sample.
- **Subjective Estimates** – When it is no longer feasible to increase the response rate by subsampling, replacement, or substitution, it may be possible to arrive at subjective estimates of the nature and effect of nonresponse bias. This involves evaluating the likely effects of nonresponse based on experience and available information.
- **Trend analysis** is an attempt to discern a trend between early and late respondents. This trend is projected to nonrespondents to estimate where they stand on the characteristic of interest.

Use of Trend Analysis in  
Adjusting for Non-response

6-51

Table 12.4

	Percentage Response	Average Dollar Expenditure	Percentage of Previous Wave's Response
First Mailing	12	412	—
Second Mailing	18	325	79
Third Mailing	13	277	85
Nonresponse	(57)	(230)	91
Total	100	275	



Chapter Thirteen

Fieldwork





## Chapter Outline

- 1) Overview
- 2) The Nature of Fieldwork
- 3) Fieldwork/Data Collection Process
- 4) Selection of Field Workers
- 5) Training of Field Workers
  - i. Making the Initial Contact
  - ii. Asking the Questions
  - iii. Probing
  - iv. Recording the Answers
  - v. Terminating the Interview



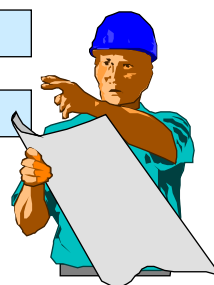
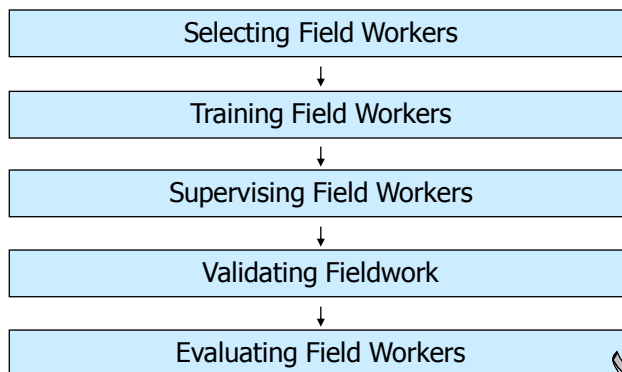
## Chapter Outline

- 6) Supervision of Field Workers
  - i. Quality Control and Editing
  - ii. Sampling Control
  - iii. Control of Cheating
  - iv. Central Office Control
- 7) Validation of Fieldwork
- 8) Evaluation of Field Workers
  - i. Cost and Time
  - ii. Response Rates
  - iii. Quality of Interviewing
  - iv. Quality of Data



## Fieldwork/Data Collection Process

Fig. 13.1



## Chapter Fifteen

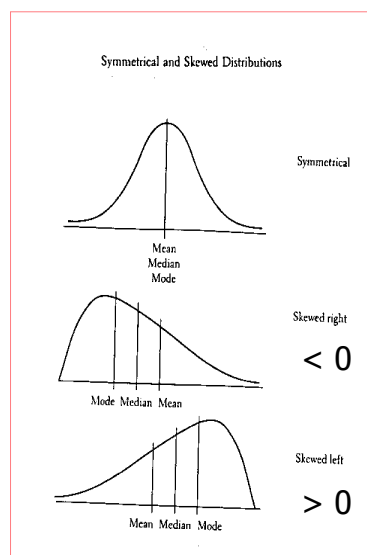
Frequency Distribution,  
Cross-Tabulation,  
and Hypothesis Testing

## Frequency Distribution

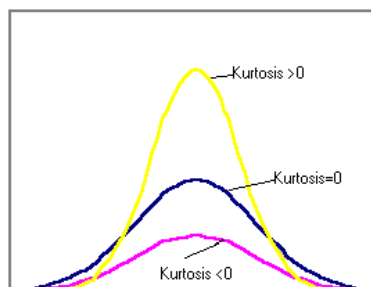
- In a **frequency distribution**, one variable is considered at a time.
- A frequency distribution for a variable produces a table of frequency counts, percentages, and cumulative percentages for all the values associated with that variable.
- Statistics Associated with Frequency Distribution: mean, median, mode, range, interquartile, variance and standard deviation (central tendency & dispersion)
- **Coefficient of variation** is the ratio of the standard deviation to the mean expressed as a percentage, and is a measure of relative variability commonly used measure of central tendency

## Statistics Associated with Frequency Distribution Measures of Shape

- **Skewness.** The tendency of the deviations from the mean to be larger in one direction than in the other. It can be thought of as the tendency for one tail of the distribution to be heavier than the other.

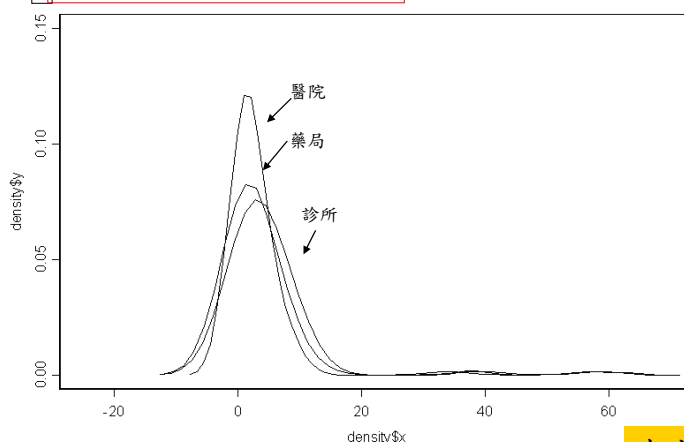


- **Kurtosis** is a measure of the relative peakedness or flatness of the curve defined by the frequency distribution. The kurtosis of a normal distribution is zero. If the kurtosis is positive, then the distribution is more peaked than a normal distribution. A negative value means that the distribution is flatter than a normal distribution.



# 口服降血糖藥品市場之品牌、通路與價格之關係研究

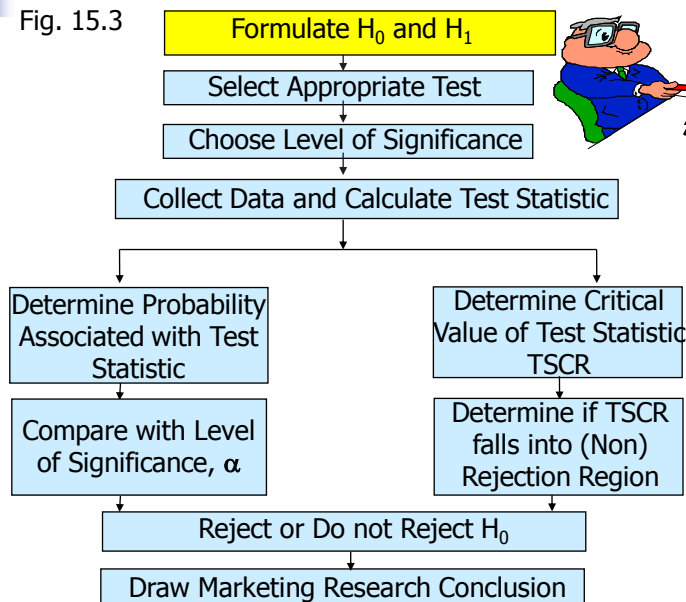
2003年通路價格密度圖



朱庭萱 (2005)

## Steps Involved in Hypothesis Testing

Fig. 15.3




## New Product Introduction



1979: 45% in Taiwan  
PepsiCo global sales: \$6 billion



Half empty or half full?



Formulate H0 and H1:  
**Hypothesis Testing**


6-63

Observed data: 45%  
Parameter  $\mu = ?$

$H_0: \mu = 45\%$   
 $H_0: \mu > 45\%$   
 $H_0: \mu < 45\%$

Type	Clear Cola
Manufacturer	PepsiCo, Inc.
Country of origin	United States
Introduced	1992
Discontinued	1993
Variants	Diet Crystal Pepsi, PepsiClear
Related products	Pepsi Blue, New Coke, Tab Clear, Pepsi

Source: Wikipedia



*A General Procedure for Hypothesis Testing*  
Step 1: Formulate the Hypothesis

6-64

- A **null hypothesis** is a statement of the status quo, one of no difference or no effect. If the null hypothesis is not rejected, no changes will be made.
- An **alternative hypothesis** is one in which some difference or effect is expected. Accepting the alternative hypothesis will lead to changes in opinions or actions.
- The null hypothesis refers to a specified value of the population parameter (e.g.,  $\mu$ ,  $\sigma$ ,  $\pi$ ), not a sample statistic (e.g.,  $\bar{X}$ ).





## *A General Procedure for Hypothesis Testing*

### Step 1: Formulate the Hypothesis

6-65

- A null hypothesis may be rejected, but it can never be accepted based on a single test. In classical hypothesis testing, there is no way to determine whether the null hypothesis is true.
- In marketing research, the null hypothesis is formulated in such a way that its rejection leads to the acceptance of the desired conclusion. The alternative hypothesis represents the conclusion for which evidence is sought.

$$H_0: \pi \leq 0.40$$

$$H_1: \pi > 0.40$$



## *A General Procedure for Hypothesis Testing*

### Step 1: Formulate the Hypothesis

6-66

- The test of the null hypothesis is a **one-tailed test**, because the alternative hypothesis is expressed directionally. (The one-tailed test provides more power to detect an effect in one direction by not testing the effect in the other direction.)
- If that is not the case, then a **two-tailed test** would be required, and the hypotheses would be expressed as:

$$H_0: \pi = 0.40$$

$$H_1: \pi \neq 0.40$$



## *A General Procedure for Hypothesis Testing*

### Step 2: Select an Appropriate Test

6-67

- The **test statistic** measures how close the sample has come to the null hypothesis.
- The test statistic often follows a well-known distribution, such as the normal,  $t$ , or chi-square distribution.
- In our example, the  $z$  statistic, which follows the standard normal distribution, would be appropriate.

$$z = \frac{p - \pi}{\sigma_p}$$

where

$$\sigma_p = \sqrt{\frac{\pi (1 - \pi)}{n}}$$



## *A General Procedure for Hypothesis Testing*

### Step 3: Choose a Level of Significance $\alpha$


6-68

#### **Type I Error**

- **Type I error** occurs when the sample results lead to the rejection of the null hypothesis when it is in fact true.
- The probability of type I error ( $\alpha$ ) is also called the **level of significance**.

#### **Type II Error**

- **Type II error** occurs when, based on the sample results, the null hypothesis is not rejected when it is in fact false. e.g., a woman is not pregnant, when in reality, she is. (sensitivity of the test)
- The probability of type II error is denoted by  $\beta$ .
- Unlike  $\alpha$ , which is specified by the researcher, the magnitude of  $\beta$  depends on the actual value of the population parameter (proportion).



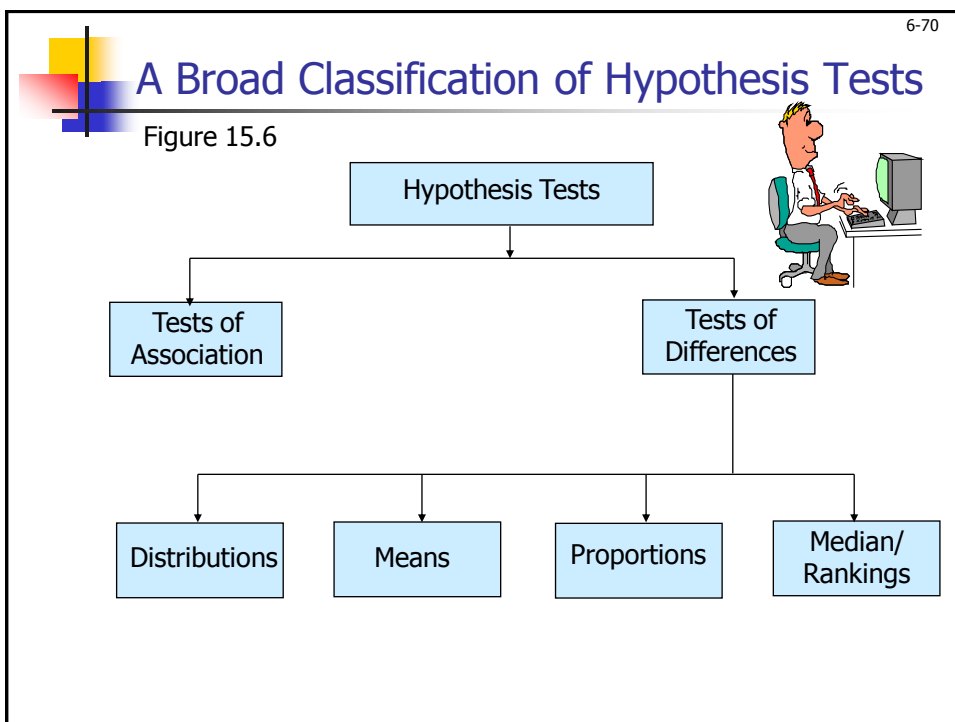
## A General Procedure for Hypothesis Testing

### Step 3: Choose a Level of Significance $\alpha$

6-69

#### Power of a Test

- The **power of a test** is the probability ( $1 - \beta$ ) of rejecting the null hypothesis when it is false and should be rejected (i.e. that it will not make a Type II error). As power increases, the chances of a Type II error decrease.
- Although  $\beta$  is unknown, it is related to  $\alpha$ . An extremely low value of  $\alpha$  (e.g.,  $= 0.001$ ) will result in intolerably high  $\beta$  errors.
- Probability of not committing a type II error in hypothesis testing.
- So it is necessary to balance the two types of errors.





## Cross-Tabulation

- While a frequency distribution describes one variable at a time, a **cross-tabulation** describes two or more variables simultaneously.
- Cross-tabulation results in tables that reflect the joint distribution of two or more variables with a limited number of categories or distinct values, e.g., Table 15.3.



## Gender and Internet Usage

Table 15.3

Internet Usage	Gender		Row Total
	Male	Female	
Light (1)	5	10	15
Heavy (2)	10	5	15
Column Total	15	15	

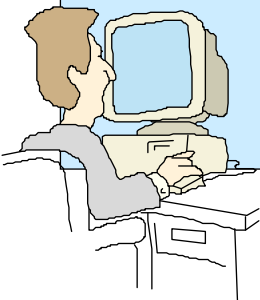




# Internet Usage by Gender

Table 15.4

Internet Usage	Gender	
	Male	Female
Light	33.3%	66.7%
Heavy	66.7%	33.3%
Column total	100%	100%




# Gender by Internet Usage

Table 15.5

Gender	Internet Usage		
	Light	Heavy	Total
Male	33.3%	66.7%	100.0%
Female	66.7%	33.3%	100.0%



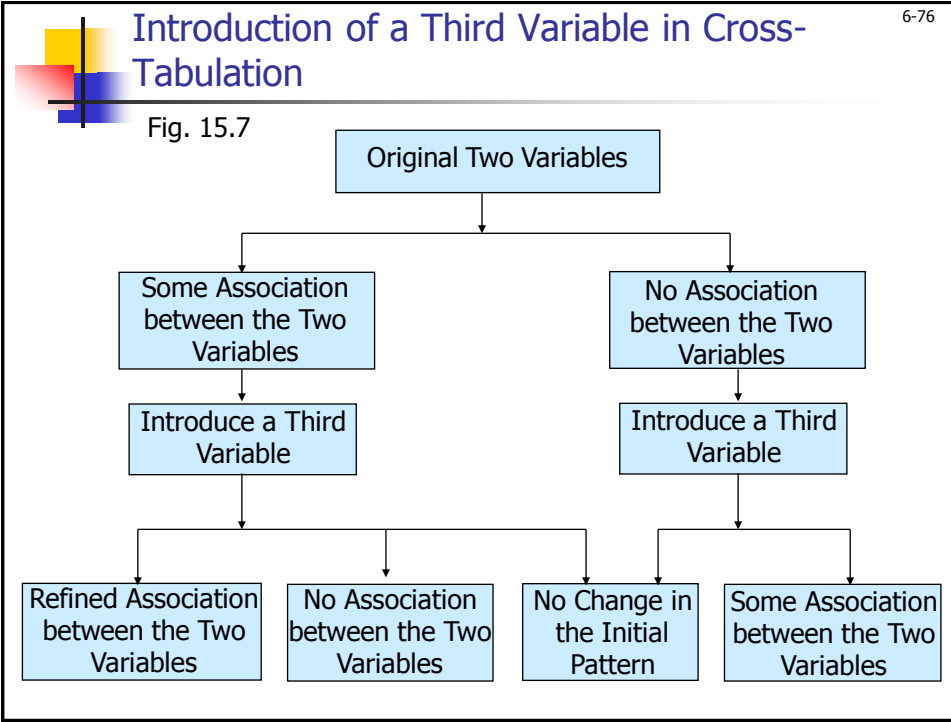


## Internet Usage by Gender

Table 15.4 (page 456)

6-75

Internet Usage	Gender	
	Male	Female
Light	33.3%	66.7%
Heavy	66.7%	33.3%
Column total	100%	100%





## Purchase of Fashion Clothing by Marital Status

Table 15.6

Purchase of Fashion Clothing	Current Marital Status	
	Married	Unmarried
High	31%	52%
Low	69%	48%
Column	100%	100%
Number of respondents	700	300



## Purchase of Fashion Clothing by Marital Status

Table 15.6

Purchase of Fashion Clothing	Current Marital Status	
	Married	Unmarried
High	31%	52%
Low	69%	48%
Column	100%	100%
Number of respondents	700	300

## Purchase of Fashion Clothing by Marital Status

Table 15.7

Purchase of Fashion Clothing	Sex			
	Male		Female	
	Married	Not Married	Married	Not Married
High	35%	40%	25%	60%
Low	65%	60%	75%	40%
Column totals	100%	100%	100%	100%
Number of cases	400	120	300	180

## Two Variables Cross-Tabulation

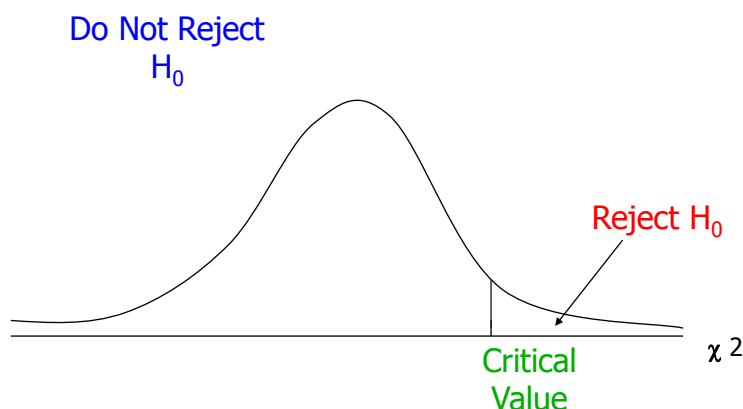
- Since two variables have been cross classified, percentages could be computed either columnwise, based on column totals (Table 15.4), or rowwise, based on row totals (Table 15.5).
- The general rule is to compute the percentages in the direction of the independent variable, across the dependent variable. The correct way of calculating percentages is as shown in Table 15.4.





## Chi-square Distribution

Figure 15.8



## Statistics Associated with Cross-Tabulation

### Phi Coefficient

- The **phi coefficient** ( $\phi$ ) is used as a measure of the strength of association in the special case of a table with two rows and two columns (a 2 x 2 table).
- The phi coefficient is proportional to the square root of the chi-square statistic

$$\phi = \sqrt{\frac{\chi^2}{n}}$$

- It takes the value of 0 when there is no association, which would be indicated by a chi-square value of 0 as well. When the variables are perfectly associated, phi assumes the value of 1 and all the observations fall just on the main or minor diagonal.



## Statistics Associated with Cross-Tabulation

### Contingency Coefficient

6-83

- While the phi coefficient is specific to a 2 x 2 table, the **contingency coefficient** ( $C$ ) can be used to assess the strength of association in a table of any size.

$$C = \sqrt{\frac{\chi^2}{\chi^2 + n}}$$

- The contingency coefficient varies between 0 and 1.
- The maximum value of the contingency coefficient depends on the size of the table (number of rows and number of columns). For this reason, it should be used only to compare tables of the same size.



## Statistics Associated with Cross-Tabulation

### Cramer's V

6-84

- **Cramer's V** is a modified version of the phi correlation coefficient,  $\Phi$ , and is used in tables larger than 2 x 2.

$$V = \sqrt{\frac{\phi^2}{\min(r-1), (c-1)}}$$

or

$$V = \sqrt{\frac{\chi^2/n}{\min(r-1), (c-1)}}$$



## Cross-Tabulation in Practice

While conducting cross-tabulation analysis in practice, it is useful to

proceed along the following steps.

1. Test the null hypothesis that there is ***no association*** between the variables using the chi-square statistic. If you fail to reject the null hypothesis, then there is no relationship.
2. If  $H_0$  is rejected, then determine the ***strength*** of the association using an appropriate statistic (phi-coefficient, contingency coefficient, Cramer's  $V$ , lambda coefficient, or other statistics), as discussed earlier.
3. If  $H_0$  is rejected, interpret the ***pattern of the relationship*** by computing the percentages in the direction of the independent variable, across the dependent variable.
4. If the variables are treated as ordinal rather than nominal, use *tau b*, *tau c*, or Gamma as the test statistic. If  $H_0$  is rejected, then determine the strength of the association using the magnitude, and the direction of the relationship using the sign of the test statistic.



## Hypothesis Testing Related to Differences

- **Parametric tests** assume that the variables of interest are measured on at least an interval scale.
- **Nonparametric tests** assume that the variables are measured on a nominal or ordinal scale.
- These tests can be further classified based on whether one or two or more samples are involved.
- The samples are **independent** if they are drawn randomly from different populations. For the purpose of analysis, data pertaining to different groups of respondents, e.g., males and females, are generally treated as independent samples.
- The samples are **paired** when the data for the two samples relate to the same group of respondents (such as repeat measurements)



## Non-Parametric Tests One Sample

- The **chi-square test** can also be performed on a single variable from one sample. In this context, the chi-square serves as a goodness-of-fit test.
- The **runs test** 連檢定法 is a test of randomness for the dichotomous variables (研究序列分佈規律, 即數值的發生順序, 是否為隨機). This test is conducted by determining whether the order or sequence in which observations are obtained is **random**.
- The **binomial test** is also a goodness-of-fit test for dichotomous variables. It tests the goodness of fit of the observed number of observations in each category to the number expected under a specified binomial distribution.



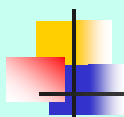
## Non-Parametric Tests Paired Samples

- The **Wilcoxon matched-pairs signed-ranks test** 配對符號檢定 analyzes the differences between the paired observations (related samples or repeated measurements), taking into account the magnitude of the differences.
- It computes the differences between the pairs of variables and ranks the absolute differences.
- The next step is to sum the positive and negative ranks. The test statistic,  $z$ , is computed from the positive and negative rank sums.
- Under the null hypothesis of no difference,  $z$  is a standard normal variate with mean 0 and variance 1 for large samples.



## Non-Parametric Tests Paired Samples

- The example considered for the paired  $t$  test, whether the respondents differed in terms of attitude toward the Internet and attitude toward technology, is considered again. Suppose we assume that both these variables are measured on ordinal rather than interval scales. Accordingly, we use the Wilcoxon test. The results are shown in Table 15.18.
- The **sign test** is not as powerful as the Wilcoxon matched-pairs signed-ranks test as it only compares the signs of the differences between pairs of variables without taking into account the ranks.
- In the special case of a binary variable where the researcher wishes to test differences in proportions, the McNemar test can be used. Alternatively, the chi-square test can also be used for binary variables.



## Non-Parametric Tests on Likert scales

- Data from Likert scales can be analyzed a number of ways using non-parametric tests, including the Mann-Whitney test, the Wilcoxon signed-rank test, and the Kruskal-Wallis test. These tests use the median rather than the mean because of the ordinal quality of the scale and the lack of a true zero.
- Data from Likert scales are sometimes reduced to the nominal level by combining all agree and disagree responses into two categories of "accept" and "reject". The Cochran Q, or McNemar-Test are common statistical procedures used after this transformation.