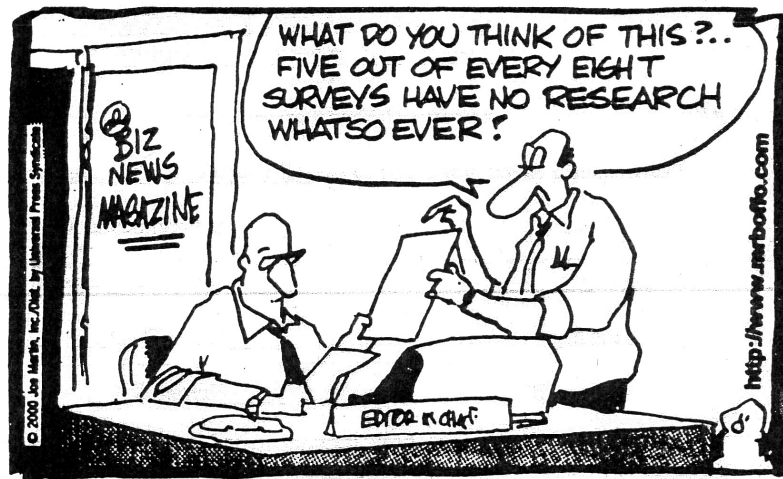


# Chapter 2

## Data Collection

MDM4U

**MISTER BOFFO** By Joe Martin



## Unit Outline

| Section | Subject                         | Homework Notes | Lesson and Homework Complete (initial) |
|---------|---------------------------------|----------------|--|
| 2.1     | Thesis Development              |                |  |
| 2.2     | Characteristics of Data         |                |  |
| 2.3     | Random Sampling                 |                |  |
| 2.4     | Survey Design and Types of Bias |                |  |
| 2.5     | Experiment Design               |                |  |

### Unit Performance

**Homework Completion:**    None            Some            Most            All

**Days absent:**\_\_\_\_\_

**Test Review Complete?**    None            Some            All

**Assignment Mark (%):**\_\_\_\_\_

**Test Mark (%):**\_\_\_\_\_

Notes to yourself to help with exam preparation:

## 2.1 - Developing a Thesis

MDM4U  
Jensen

### Part 1: ISU Intro

This chapter will prepare you to begin your ISU that is worth 10% of your final grade. For the ISU you will be required to choose a topic that interests you and conduct a study that analyses large amounts of data using:

- one-variable statistics tools (chapter 3)
- two variable statistics tools (chapter 1)
- probability (chapter 4/5)

### Part 2: Mind-Map

Before you can begin your project, you must create a thesis:

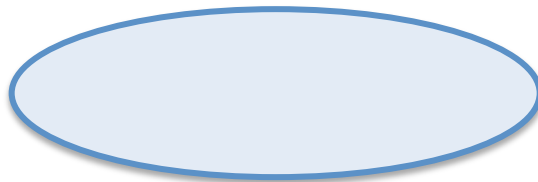
**thesis:** a formal statement or question that your project will answer or discuss

To begin creating a thesis, you must first determine what topics interest you and then determine what concepts related to that topic you want to study. A useful brainstorming tool that can illustrate how a topic relates to other concepts is a *mind map*.

**mind map:** a visual display used in brainstorming to illustrate relationships

### Constructing a Mind Map

1. Start by making a mind map of your interests with you at the centre. Start off as simple as possible and draw arrows to show how topics are connected. Work from the inside out.



## Extended Mind Map

2. Pick one of the topics from your mind map and extend it with sub-topics.



### Part 3: Thesis Question Development

Once you have narrowed down your topic, you will need to pose a problem that you plan to investigate.

#### Money in Sports

3. Brainstorm and create number of questions that can be explored with the use of statistical information

**a)** How do people at my school feel about high salaries in professional sports?

**b)** How have salaries paid to professional hockey players changed from 1960 to present?

**c)** Is there a relationship between a very large salary increase to an athlete and his or her subsequent performance?

**d)** Does the amount a country spends to prepare its athletes for the Olympics correspond to the country's success at the games?

## **Thesis Question Analysis**

Questions to ask of your Thesis:

- i.** What are the main variables in my question?
- ii.** Can these variables be measured statistically?
- iii.** Is there enough data to make an interesting analysis

**4.** Once you have chosen your thesis, analyse it using the three questions above to make sure your study will be able to provide an insightful answer.

**Thesis:** Is there a relationship between a very large salary increase to an athlete and his or her subsequent performance?

**Analysis:**

- i.** player salaries, performance statistics (goals, home-runs, etc.)
- ii.** yes; however it may be difficult to choose which performance statistics to use
- iii.** yes there would be lots of available data for professional athletes and their salaries and performance.

**Project tips:**

One way of posing a problem is to generate questions from data. For example, once a topic has been identified, do a preliminary data search. The type and quantity of available data may indicate some possible questions. Data from print sources, the Internet, and E-Stat are some resources that may be used.

## 2.2 - Characteristics of Data

MDM4U  
Jensen

### Part 1: Population vs. Sample

**Data** are any collection of numbers, characters, images, or other items that provide information about something.

The entire group of individuals that we want information about is called the **population**.

A **census** is an attempt to gather information about every individual member of the population. Problems with census—**costs**; **time** needed to complete; sometimes testing can **destroy** items.

A **sample** is a part of the population that we actually examine in order to gather information.

**Note:** It usually isn't practical to collect data from the entire population; instead you should take a representative sample and study it.

**Example 1:** Determine the population of each of the following questions

a) Whom will you plan to vote for in the next Ontario election

All legal voters in Ontario

b) What is your favourite brand of hockey stick?

All hockey players

c) Do women prefer to wear ordinary glasses or contact lenses?

All women who wear glasses and/or contacts

Once you have identified the population, you need to decide how you will obtain your data. If the population is **small**, it may be possible to survey the entire group (census). For **larger** populations, you need to use appropriate sampling technique.

We will discuss different sampling techniques next lesson.

## Part 2: Types of Studies

### **Cross Sectional:**

a study that considers individuals from different groups at the same time

(specific time frame, range of people)

### **Longitudinal:**

a study that considers individuals over a long period of time.

(extended period, small group of people)

### **Example 2:**

#### **For the thesis question:**

*How do the opinions about the cafeteria change among students from Grade 9 to Grade 12?*

**a)** How could you conduct a cross-sectional study?

Ask students from each grade about their opinions of the cafeteria

**b)** How could you conduct a longitudinal study?

Interview a selection of grade 9 students and then return to ask them again each year

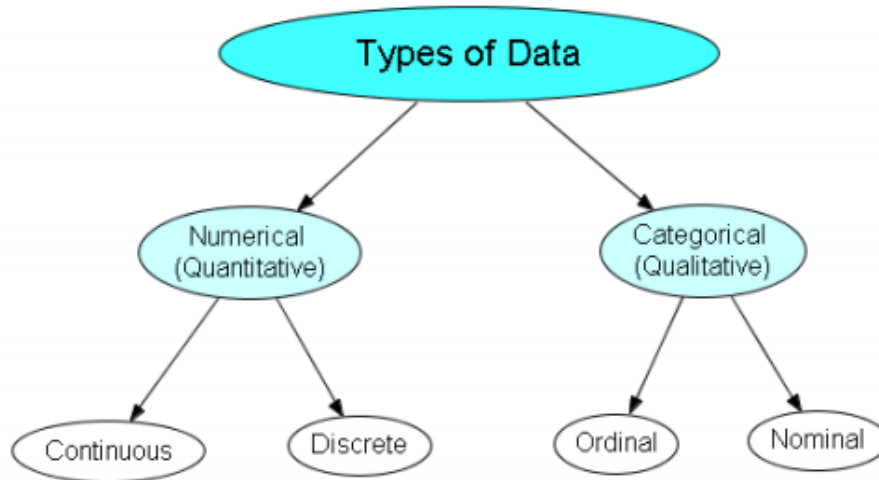
**c)** Which study would be more time efficient?

Cross-sectional study would be more practical; especially since you won't go to this school next year.

**d)** Re-write the thesis question to reflect a cross-sectional study

How do the opinions about the cafeteria among a random sample of students in Grades 9 and 12 differ?

### Part 3: Types of Variables



**Quantitative/Numeric Variable:** A quantitative variable that takes numerical values for which it makes sense to find an average. These variables can be either continuous or discrete

**Qualitative/Categorical Variable:** A variable that places an individual into one of several groups or categories (also known as qualitative variables). Categorical variables may have categories that are naturally ordered (ordinal variables) or have no natural order (nominal variables).

**Example 3:** Identify whether each of the following questions measures a qualitative or quantitative variable.

a) How tall are you?

QUANTITATIVE

b) What conference are the Leafs in?

QUALITATIVE

c) What colour is your hair?

QUALITATIVE

d) How many students are in this class?

QUANTITATIVE

e) What is your favourite school subject?

QUALITATIVE



## Part 4: Types of Quantitative Variables

**Continuous Variable:** A numeric variable that can have an infinite number of values in a given interval. Measurable with all real numbers.

Examples: **temperature, height, weight, speed**

**Discrete Variable:** A numeric variable that can take on only a finite number of values within a given range. (usually measured with integer values only)

Examples: **number of dogs, number of goals scored, number of siblings**

**Example 4:** Classify each quantitative variable as either continuous or discrete

**a)** Temperature outside

**CONTINUOUS**

**b)** Number of goals scored by Crosby

**DISCRETE**

**c)** Number of songs on your iPod

**DISCRETE**

**d)** Speed of Zdeno Chara's slapshot (108.8 mph) <https://www.youtube.com/watch?v=vZssDq7lJus>

**CONTINUOUS**

## 2.3 – Sampling Principles

MDM4U  
Jensen

### Part 1: Random Rectangles Activity

1. a. Guess the average area of all rectangles on the page: **(guess)** \_\_\_\_\_
- b. Choose six rectangles (before you calculate any areas) that you think represent the entire population of rectangles well.

6 rectangles – subjective – “rectangle **expert**”:

rectangle number

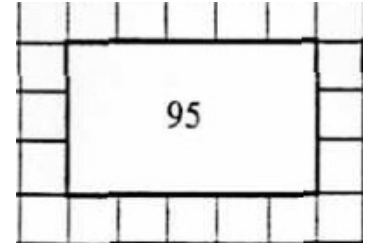
\_\_\_\_\_  
\_\_\_\_\_  
\_\_\_\_\_  
\_\_\_\_\_  
\_\_\_\_\_  
\_\_\_\_\_

area

\_\_\_\_\_  
\_\_\_\_\_  
\_\_\_\_\_  
\_\_\_\_\_  
\_\_\_\_\_  
\_\_\_\_\_

average:

\_\_\_\_\_



2. a. After setting a new seed value on your calculator, use the randint function to choose six random rectangles for you.

6 rectangles – **random**:

rectangle number

\_\_\_\_\_  
\_\_\_\_\_  
\_\_\_\_\_  
\_\_\_\_\_  
\_\_\_\_\_  
\_\_\_\_\_

b. area

\_\_\_\_\_  
\_\_\_\_\_  
\_\_\_\_\_  
\_\_\_\_\_  
\_\_\_\_\_  
\_\_\_\_\_

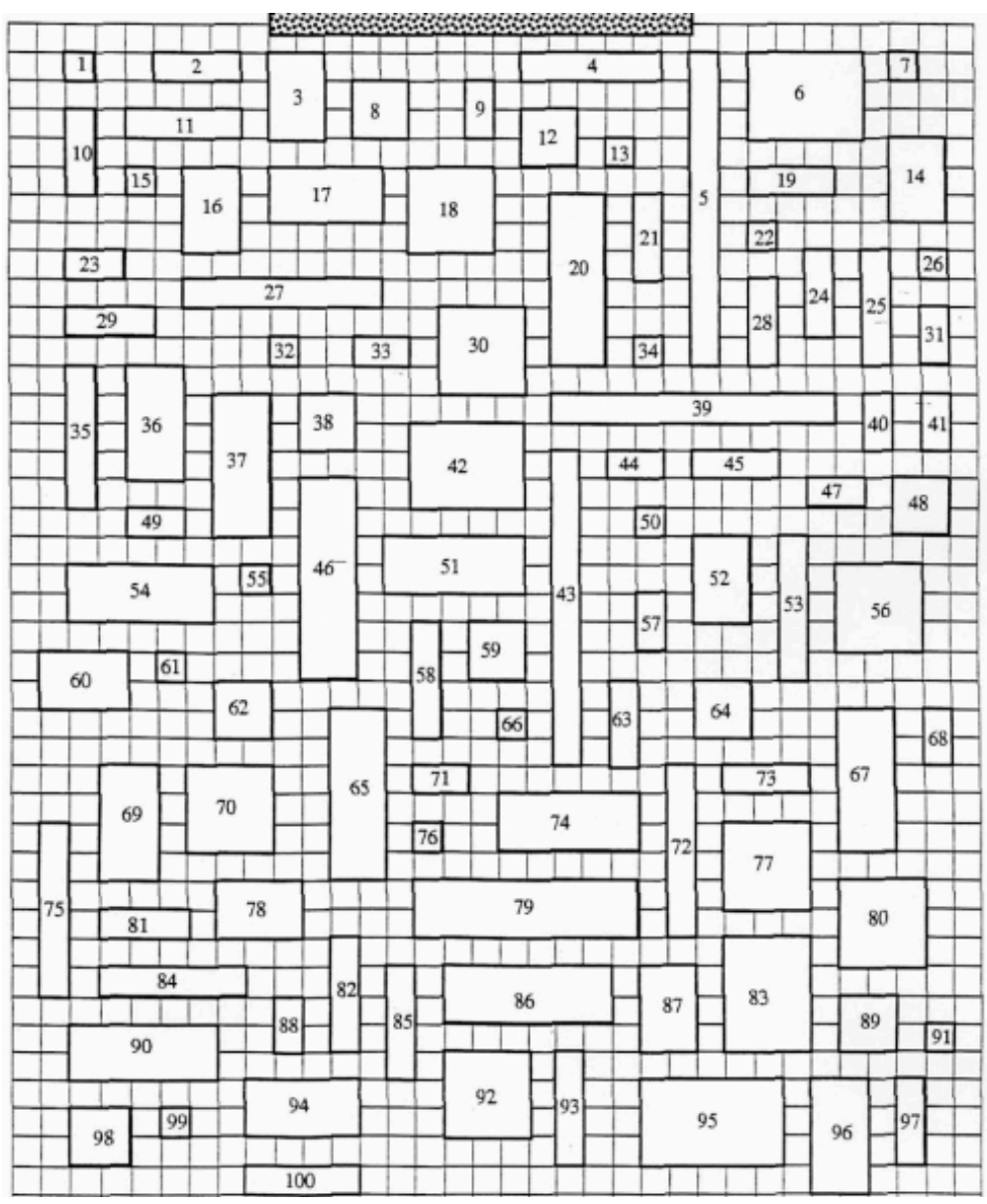
average:

\_\_\_\_\_

3. a. mean of sample averages:
- guesses \_\_\_\_\_
  - subjective (expert) \_\_\_\_\_
  - random \_\_\_\_\_
- c. actual area of 100 rectangles (population): \_\_\_\_\_

Wrap-up (what have you learned?):

The design of a study is biased if it systematically favors certain outcomes. The design of a study shows bias if it consistently over or under estimates the value you want to know. Random sampling is necessary to get a representative sample.



## Part 2: Random Sampling Methods

### 1. Simple Random Sampling

A sample is a **simple random sample** if it is selected so that:

- each member of the population is **equally** likely to be chosen and the members of the sample are chosen independently of one other;

OR

- every set of  $n$  units has an **equal** chance to be the sample actually selected.

**Example:** Put names in hat and draw until have desired sample size; more commonly, number names and use random number generator or other source of random numbers to select sample. Notice that some type of unbiased method must be used; haphazard  $\neq$  random.

### 2. Systematic Random Sampling

A sample is a systematic random sample if you randomly choose some **starting point**; then select every  **$n^{\text{th}}$**  element in the population, where  $n$  is the sampling interval. This guarantees that the sample is taken from throughout the **population** but it requires an ordered list of everyone in the population.

**Example:** If we wanted to get a systematic random sample of 10% of the students from King's which has approximately 600 students...

- Calculate number of students required for sample:  $600 \times 0.10 = 60$
- Calculate the sampling interval:  $\text{sampling interval} = \frac{\text{population size}}{\text{sample size}} = \frac{600}{60} = 10$
- Choose a random starting point using a random number generator
- Include every 10<sup>th</sup> student from the randomly chosen starting point in your sample

### 3. Stratified Random Sampling

When using a stratified random sample, the population is divided into **groups** called **strata** (e.g. age, geographical areas, grade, etc.)

A **simple random sample** of the members of **each** stratum is then taken. The size of the sample for each stratum is **proportionate** to the stratum's size (you must survey the same **percentage** of people from each stratum).

**Example:** If we want a stratified random sample of 10% of the 600 King's students, we can divide the population into four groups based on grade (9, 10, 11, 12) and then take a simple random sample of 10% of the students in each grade.

## 4. Cluster Random Sampling

When using a cluster random sampling method, divide the population into **groups** or **clusters**; randomly select a few of those groups and then sample **all** members from the selected groups.

**Example:** **Randomly** select 5 block C classes—survey **all** students in each class selected.

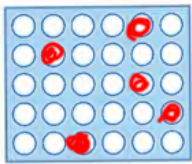
## 5. Multi-Stage Random Sampling

When using multi-stage random sampling, the population is organized in to groups, a simple random sample of groups is chosen, and then a simple random sample of people within the chosen groups is taken.

**Example:** **Randomly** select 5 block C classes—survey **a random sample of 10%** of the students in each class selected.

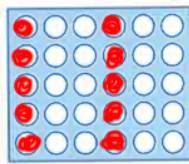
## Review of Different Random Sampling Techniques:

### *Simple Random*



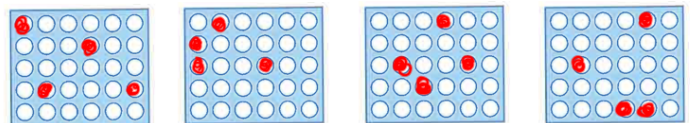
- all selections are equally likely

### *Systematic Random*



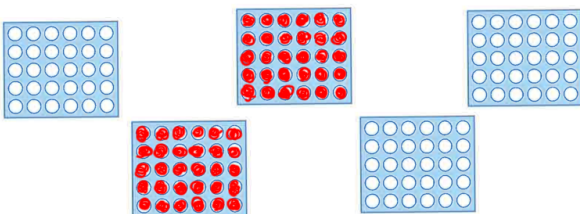
- random starting point  
choose individuals at interval (every  $n$ th person)

### *Stratified Random*



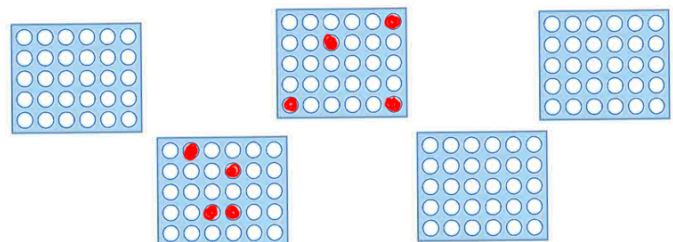
- Divide population into groups then survey an equal percentage of each group.

### *Cluster Random*



- Divide the population into groups. Choose a random sample of groups and then survey **every member** of the groups chosen.

### *Multi-Stage Random*



- Divide the population into groups. Choose a random sample of groups and then choose a random sample of members of the chosen groups.

### Part 3: Types of Non-Random Samples

#### 1. Convenience sampling

The easiest way to obtain a sample is to choose it without any random mechanism (also called haphazard sampling). Choosing individuals from the population who are easy to reach results in a convenience sample. Convenience sampling often produces unrepresentative data.

**Example:** Suppose we want to know how long students at a large high school spent doing homework last week. We might go to the school library and ask the first 30 students we see about their homework time.

#### 2. Voluntary Response Sampling

A voluntary response sample consists of people who choose themselves by responding to a general invitation. Voluntary response samples attract people who feel strongly about an issue, and who often share the same opinion. This leads to bias.

**Example:** A radio host invites listeners to call in to give opinions on a new band.

### Part 4: River Activity

A farmer has just cleared a new field for corn. It is a unique plot of land in that a river runs along one side. The corn looks good in some areas of the field but not others. The farmer is not sure that harvesting the field is worth the expense. He has decided to harvest 10 plots and use this information to estimate the total yield. Based on this estimate, he will decide whether to harvest the remaining plots.



Part I.

#### A. Method Number One: Convenience Sample

The farmer began by choosing 10 plots that would be easy to harvest. They are marked on the grid below:

|   |  |  |  |  |  |  |  |  |  |
|---|--|--|--|--|--|--|--|--|--|
| X |  |  |  |  |  |  |  |  |  |
| X |  |  |  |  |  |  |  |  |  |
| X |  |  |  |  |  |  |  |  |  |
| X |  |  |  |  |  |  |  |  |  |
| X |  |  |  |  |  |  |  |  |  |
| X |  |  |  |  |  |  |  |  |  |
| X |  |  |  |  |  |  |  |  |  |
| X |  |  |  |  |  |  |  |  |  |
| X |  |  |  |  |  |  |  |  |  |
| X |  |  |  |  |  |  |  |  |  |

River

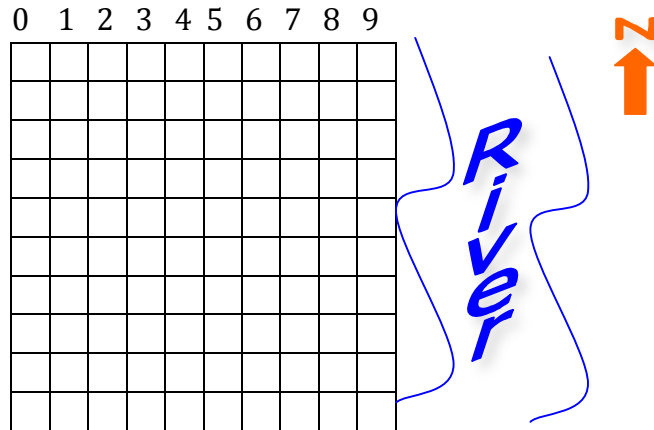


Since then, the farmer has had second thoughts about this selection and has decided to come to you (knowing that you are an AP statistics student, somewhat knowledgeable, but far cheaper than a professional statistician) to determine the approximate yield of the field.

You will still be allowed to pick 10 plots to harvest early. Your job is to determine which of the following methods is the best one to use – and to decide if this is an improvement over the farmer’s original plan.

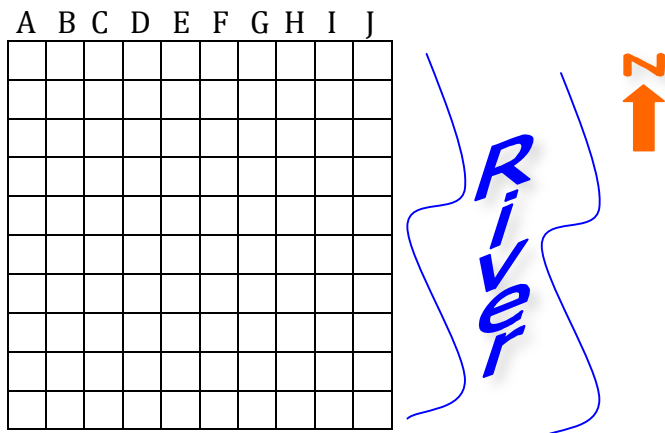
**B. Method Number Two: Simple Random Sample**

Use your calculator or a random number table to choose 10 plots to harvest. Mark them on the grid below, and describe your method of selection.



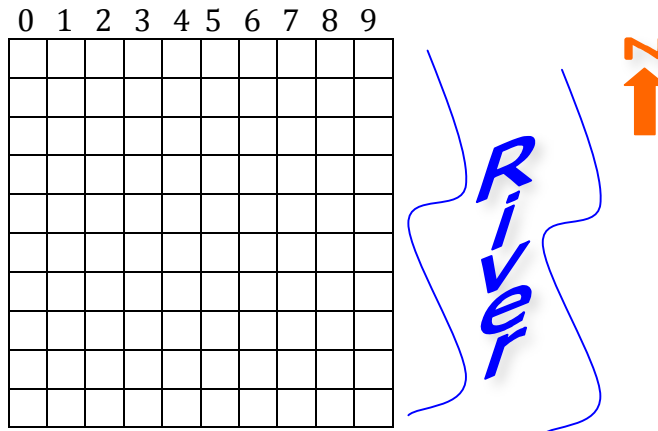
**C. Method Number Three: Stratified Sample**

You and the farmer think the river might have a strong influence on corn production so you decide to consider the field as grouped in vertical columns (called strata—remember you can only stratify data your sample when you think a factor will have a strong influence on the outcome.). Using your random number table, randomly choose one plot from each vertical column and mark on the grid. (Label your columns A through J, rows 0 through 9.)



### D. Method Number Four: Stratified Sample

You and the farmer rethink the plan and decide that direction (north—south) may have a strong influence on corn production. You decide to consider the field as grouped in horizontal rows (also called strata). Using your random number table, randomly choose one plot from each horizontal row and mark them on the grid. (Label your rows A through J, columns 0 through 9.)



OK, the crop is ready! Below is a grid with the yield per plot. Estimate the average yield per plot based on each of the four sampling techniques.

|   |    |    |    |    |    |    |    |    |    |
|---|----|----|----|----|----|----|----|----|----|
| 6 | 17 | 20 | 38 | 47 | 55 | 69 | 76 | 82 | 97 |
| 7 | 14 | 23 | 34 | 43 | 56 | 63 | 75 | 81 | 92 |
| 2 | 14 | 28 | 30 | 50 | 50 | 62 | 80 | 85 | 96 |
| 9 | 15 | 27 | 34 | 43 | 51 | 65 | 72 | 88 | 91 |
| 4 | 15 | 28 | 32 | 44 | 50 | 64 | 76 | 82 | 97 |
| 5 | 16 | 27 | 31 | 48 | 59 | 69 | 72 | 86 | 99 |
| 5 | 18 | 28 | 34 | 50 | 60 | 62 | 75 | 90 | 90 |
| 8 | 15 | 20 | 38 | 40 | 54 | 62 | 77 | 88 | 93 |
| 7 | 17 | 29 | 39 | 44 | 53 | 61 | 77 | 80 | 90 |
| 7 | 19 | 22 | 33 | 49 | 53 | 67 | 76 | 86 | 97 |

| Sampling Method               | Mean yield per plot | Estimate of total yield |
|-------------------------------|---------------------|-------------------------|
| Convenience Sample (farmer's) |                     |                         |
| Simple Random Sample          |                     |                         |
| Vertical Strata               |                     |                         |
| Horizontal Strata             |                     |                         |



## Observations:



- 1) You have looked at four different methods of choosing plots. Is there a reason, other than convenience, to choose one method over another?

One needs to choose a method that will give the best estimate of the yield. This can be affected by factors that cannot be controlled: e.g. the placement of the river. That's why one shouldn't choose the ten plots chosen by the farmer.

- 2) How did your estimates vary according to the different sampling methods you used?

The student will see that the farmer's sample yields a very low estimate compared to the other methods used.

- 3) Compare your results to someone else in the class. Were your results similar?

Comparing results with a peer helps the student verify that the sampling was done correctly. This does not mean the students will have the same sample, but each student should use the same process of drawing a sample for a given method. Some methods will produce highly variable results while others are much more consistent.

- 4) When we compare the class boxplots for each sampling method. What do you see?

The variability of the means of the sample yields, as shown by the length of the boxplot and the width of the middle 50%, will reduce drastically once the student has stratified appropriately. Thus the strata that are effective are the vertical ones, in which the values in each stratum are similar. This stratification reduces the variation in the sample means since the values chosen for a particular stratum vary little from sample to sample relative to the variability in the population.

- 5) Which sampling method should you use? Why do you think this method is best?

Vertical stratification should be used since the sample would then include higher yielding plots as well as lower yielding ones.

- 6) What was the actual yield of the farmer's field? How did the boxplots relate to this actual value?

The actual yield is 5004. The class boxplot for the means resulting from the vertical stratification should be centered near  $5004/100$  or about 50.

## 2.4 – Bias and Survey Design

MDM4U  
Jensen

If you conduct a survey and collect information firsthand, this is called **primary** data. This type of data is easy to work with because you control how it is collected.

Information obtained from similar studies conducted by OTHER researchers is called **secondary** data.



### Part 1: Principles of Survey Design

#### **Basic Principle #1:**

A survey is not merely a collection of questions, thrown together without purpose—surveys should be designed around specific needs for information about a **relevant** topic.

#### **Basic Principle #2:**

Both parties to the survey have responsibilities:

- The interviewer's work must be mostly done in advance; identify relevant variables, craft questions, design the flow of the survey.
- The interviewee's task is to—having agreed to answer questions—be **truthful**.

#### **Basic Principle #3:**

A prime task of the interviewer at the question design stage is to help the interviewee be honest.

### Part 2: Open vs. Closed Questions

#### **1. Open Questions**

- answered in respondents own words
- wide variety of possibilities
- answers sometimes difficult to interpret

#### **Examples:**

*How do you feel about the salaries paid to professional athletes?*

*What is the most important issue for King's students?*

#### **2. Closed Questions**

- respondents select from a given list of responses or the question requires an exact response
- answers are easily analyzed
- options present may bias results

### Part 3: Types of Closed Questions

#### **i) Information**

Circle the appropriate response:

a) Gender:    M        F

b) Age:            under 14                    15 or 16  
                      17 or 18                    19 and over

#### **ii) Checklist**

Which of the following sports do you enjoy watching? (check all that apply)

- |                                     |                                   |
|-------------------------------------|-----------------------------------|
| <input type="checkbox"/> Basketball | <input type="checkbox"/> UFC      |
| <input type="checkbox"/> Baseball   | <input type="checkbox"/> Lacrosse |
| <input type="checkbox"/> Hockey     | <input type="checkbox"/> Soccer   |

**iii) Rating** – asks survey respondents to compare different items using a common scale. It can also be used just to rate one item using a scale.

How satisfied were you with your grade from the first unit test? (check the one that applies)

\_\_\_ Very dissatisfied

\_\_\_ Dissatisfied

\_\_\_ Satisfied

\_\_\_ Very Satisfied

Using a scale of 0 = not at all to 4 = very important, please rate the importance of each of the following aspects of service in a restaurant

|                       | 0                        | 1                        | 2                        | 3                        | 4                        |
|-----------------------|--------------------------|--------------------------|--------------------------|--------------------------|--------------------------|
| Speed of service      | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> |
| Friendliness of staff | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> |
| Helpfulness of staff  | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> |
| Value for money       | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> |
| Taste of food         | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> |

**iv) Ranking** – asks survey respondents to compare a list of objects to one another by ORDERING them

When choosing a restaurant to eat at, please rank the following in order of importance from 1 to 4 where 1 is the most important to you and 4 is the least important to you

\_\_\_ Speed of Service    \_\_\_ Ease of parking    \_\_\_ Cleanliness            \_\_\_ Friendliness of staff

## Part 4: Good vs. Bad Questions

Good Questions are: **simple, specific, relevant, readable**

Good Questions avoid: **jargon, abbreviations, negatives, leading respondents, insensitivity**

**Example 1:** What's wrong with each of the following questions?

1. Given the increasing problem of obesity amongst teenagers in North America, do you agree that King's should make physical education a mandatory class for every grade?

**Leading respondents**

2. Do you think the NHLPA should have agreed to the last CBA?

**Abbreviations**

3. Which player would you not select first in a fantasy hockey draft?

- Ovechkin
- Crosby
- Malkin
- Stamkos

**Negatives, possibly jargon**

## Part 3: Types of Bias

The results of a survey can be accurate only if the sample is **representative** of the population and the measurements are objective. The methods used for choosing the sample and collecting the data must be free from **bias**. Statistical bias is any factor that favours certain outcomes or responses and hence systematically **skews** the survey results.



**Sampling Bias:** When the chosen sample does not accurately represent the population

**Household Bias:** When one type of respondent is overrepresented because groupings of different sizes are polled equally instead of proportionately

**Non-response Bias:** Occurs when an individual chosen for the sample can't be contacted or refuses to participate

**Response/Measurement Bias:** Refers to anything in the survey design that influences the responses. This includes but is not limited to:

- tendency of respondents to tailor responses to try to please the interviewer
- natural unwillingness of respondent to reveal personal facts or admit to bad behavior
- the wording of questions can influence responses

## Example 2: Identifying Bias

You are the campaign manager for your best friend, Rebecca, who is running for student council Prime Minister. You have been asked to determine the overall level of support for Rebecca among the 1500 students at your school. Design a sampling method that will provide the least **sampling bias**.

### Potential Solution - Plan A

To save time, you have decided that a sample of about 50 students will provide a good picture of the school's political landscape. Students have lunch periods 2, 3, or 4. By random draw from a hat, you have decided to conduct the survey in the cafeteria during period 4. The first 50 students who enter the cafeteria are given the questionnaire, and you instruct them to fill it out and return it to you before the end of lunch.

**What is wrong with this scenario?**

**Non-response bias** - some student may not complete or return the survey

**Sampling bias** - perhaps more seniors were let out of class early (seniors are over-represented)

- only 50 out of 1500 students were surveyed (should survey at least 10% of population)

### Plan B

To fix the problems with Plan A, you have decided to provide a questionnaire to one person from each homeroom (your sample size is now 73). You can wait until the respondent finishes with the questionnaire to collect it. This will eliminate the non-response bias.

**What is wrong with this scenario?**

**Sampling bias** - still only 73 students out of 1500 (less than 10%)

**Response bias** - some students may just rush the survey to get through it or answer dishonestly

**Household bias** - some homerooms are bigger than others

**Create a Plan C that is free from as much bias as possible:**

**Sample Answer:** A stratified random sampling technique could be used to ensure a suitable sample of the student body. Students in each grade could be assigned a number. The appropriate number of students from each grade could then be selected by using a random number generator. The table below shows how a sample of 150 students could be selected to ensure that each grade is represented proportionately to its population. Interviews with each student selected would eliminate non-response bias.

- 10% of grade 9's – 42
- 10% of grade 10's – 42
- 10% of grade 11's – 36
- 10% of grade 12's – 30

**Example 3: Identifying Sources of Response Bias**

Consider the questionnaire below developed by Rebecca's friends. Identify examples of response bias.

**Election Survey**  
(brought to you by the friends of Rebecca committee)

Circle the appropriate response

Gender:      Male      Female

Grade:            9    10    11    12

On Election Day, I intend to vote for:  
**Rebecca**      Mable      Jacob

Circle what you would like:

- more dances
- more theme-dress days
- more holidays
- more fun

Brought to you by friends of Rebecca – may lead to respondents trying to please interviewer with answers

Grade 9, 10, 11, 12 – may confuse students taking classes in different levels

**Rebecca** – bolding the name may lead to more people choosing that name

More fun – not specific enough; won't generate any useful information

## 2.5 - Experiment Design

MDM4U  
Jensen

### Part 1: Experiment Design Video

<http://www.learner.org/courses/againstalldds/unitpages/unit15.html>

While watching the video, answer the following questions

1. Why is the study of the effect of humans on the coral reefs not an experiment?

The study did not impose human populations on the various coral reefs. Instead, scientists simply observed the health of the coral reefs in four areas where human interaction with the areas was varied from no humans living in the area to a sizable population of humans currently living in the area.

2. Who were the subjects in the Glucosamine/Chondroitin study? What did researchers want to find out?

The subjects were patients suffering with osteoarthritis of the knee. Researchers wanted to compare the effects on joint pain of the dietary supplements of Glucosamine or Chondroitin compared to a prescription medication or a placebo.

3. Why were subjects randomly assigned to the treatments?

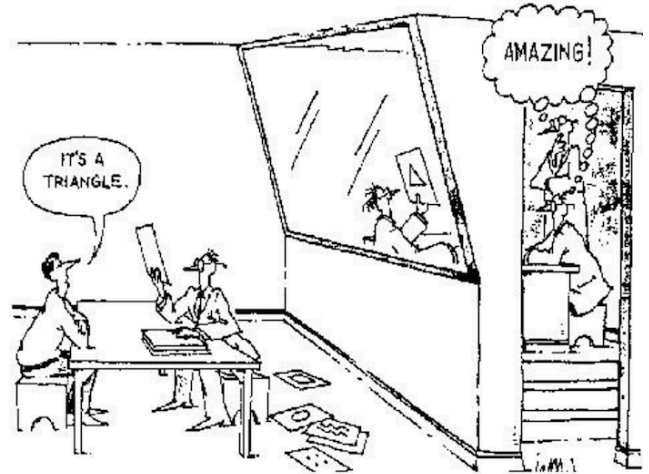
Randomization produces groups of subjects that should be similar in all respects before the treatments are applied. It allows us to equalize the effect from unknown or uncontrollable sources of variation.

4. Dr. Confound conducted a very badly designed experiment on mood-altering medication. List some of the problems with his experiment.

Sample answer: His sample size was extremely small (the last two he called 7 and 8, so there were 8 subjects total). He treated the two subjects differently – one was allowed to sit and the other had to stand for over an hour. The treatment and having stand are now confounding variables. This difference in treatment would certainly affect subjects' moods. He didn't randomly assign the medications. He interacted with the patients sympathizing with their responses. He didn't record exactly what one of his patients said and instead recorded only the higher ranking of mood.

## Part 2: Observational Studies vs. Experiments

A **sample survey** aims to gather information about a population without disturbing the population in the process. Sample surveys are one kind of **observational** study. Other observational studies watch the behavior of animals in the wild or the interactions between teacher and students in the classroom. This section is about statistical designs for **experiments**, a very different way to produce data.



In contrast to observational studies, experiments don't just observe individuals or ask them questions. They actively impose some **treatment** to measure the response. The purpose of an experiment is to determine whether the treatment **causes** a change in the response.

When our goal is to understand **cause and effect**, randomized experiments are the only source of fully convincing data.

An experimenter must identify at least one **independent** variable to manipulate (this is the treatment) and at least **one dependent** variable (response) to measure. The experimenter deliberately manipulates the treatments and must assign subjects to treatments at random.

**Experimental units** (subjects) are the collection of individuals to which treatments are applied.

### **Example 1: Observation vs. Experiment**

Should women take hormones such as estrogen after menopause, when natural production of these hormones ends? Several major medical organizations thought yes because women who took hormones seemed to reduce their risk of a heart attack 35 to 50%. The evidence in favour of hormone replacement came from a number of observational studies that compared women who were taking hormones with other who were not. But the women who chose to take hormones were richer and better educated and saw doctors more often than women who didn't take hormones. It isn't surprising that they had fewer heart attacks. In this scenario, wealth, education level, and number of doctor visits are **confounding** (we don't know if it was the hormone or any of these variables that caused a reduce in heart attacks)

To get convincing data on the link between hormone replacement and heart attacks, we should do an experiment. Experiments don't let women decide what to do. They assign women to either hormone replacement pills or to placebo pills that look and taste the same as hormone pills. The assignment is done by a coin toss, so that all kinds of women are equally likely to get either treatment.

By 2002, several experiments with women of different ages agreed that hormone replacement does not reduce the risk of heart attacks. In fact, some studies concluded that hormone replacement with estrogen carried increase risk of stroke.



**Example 2:** In 2007, deaths of a large number of pet dogs and cats were ultimately traced to contamination of some brands of pet food. The manufacturer now claims that the food is safe, but before it can be released, it must be tested. In an experiment to test whether the food is now safe for dogs to eat, what would be the treatments and what would be the response variable measured?

**Treatments:** ordinary sized portions of two dog food: the new one from the company, and one other type that is known to be safe

**Response:** a veterinarian's assessment of the health of the test subjects

*Note: the test subjects (dogs) must be randomly assigned to either treatment*

### **Part 3: Experimental Design**

#### **4 Principles of Experimental Design**

1. **Comparison** – use a design that compares two or more treatments
2. **Random Assignment** – Use chance to assign experimental units to different treatments.
3. **Control** – Keep other variables (besides the ones you are testing) that might affect the response of the subject the same for all groups.
4. **Replication** – use enough experimental units in each group so that any differences in the effects of the treatments can be distinguished from chance differences between groups

**Example 3:** We're planning an experiment to see if the new dog food is safe to eat. We have established that we will feed some dogs the new food and some dogs food that is known to be safe (principle of comparison). In this experiment, how could you implement the principles of control, random assignment, and replication?

**Control:**

- control portion sizes
- control environment (pen, amount of water drank, amount exercise and sleep)
- restrict experiment to single breed of dog

**Random Assignment:**

- assign dogs to the two different treatments randomly by flipping a coin

**Replication:**

- Assign more than one dog to each treatment to allow for variability among dogs.

## Strategies to Improve Experiments

**1. Use a control group** – researchers vary the independent variable (treatment) for the **experimental** group but not for the **control** group. Any differences in the dependent variable (response) for the two groups can be attributed to the changes in the independent variable.

Example: A medical researcher wants to test a new drug believed to help smokers quit. 50 people volunteer for the study. The researcher randomly divides the smokers into two groups. One group is given nicotine patches with the new drug, while the second group uses ordinary nicotine patches. The researcher then measures how many in each group quit smoking.

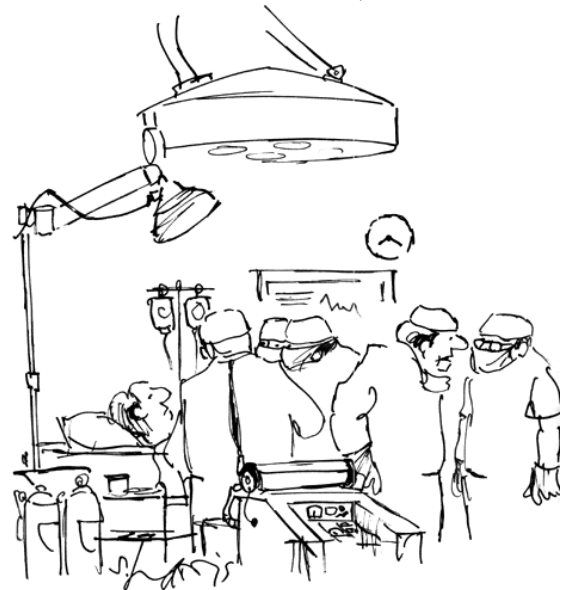
**2. Blinding** – keep anyone who could affect the outcome of the response from knowing which **subjects** have been assigned to which **treatments**. A **double-blind** experiment is when both the subject and experimenter don't know which treatment the subject has been given.

Example: in the earlier pet food example, the vet should not be told which dogs ate which food.

**3. Use a placebo** – often, simply applying **any** treatment can induce an improvement. A **fake** treatment that looks just like the treatments being tested is called a placebo. Placebos are the best way to blind subjects from knowing whether they are receiving the treatment or not.

**4. Blocking** – group **similar** experimental units together. Then random assignment of subjects to treatments is carried out separately within each block.

Example: in the previous dog food example, different breeds of dogs may respond differently to the foods. Blocking by breeds can remove that variation.



*"We'll just mill around till he's asleep, and then send him back up. This operation is actually for a placebo effect."*

### Example 4: Tire Blocking

A firm wishes to test the durability of four tire types that we'll call A, B, C, and D for convenience. Here are four possible studies they might perform. In all cases, the cars are to be driven on a track under controlled conditions until its tires are deemed "worn out". The response variable for each experimental unit (a car) is the number of miles the car drove with the tires. Each of the first three designs contains at least one serious weakness. Comment briefly on them. The fourth design is called a blocked design. State what the blocks are and explain what the advantage is of this design over design number 3.

**1.** Four Cadillacs of the same type are purchased new from four dealers. One gets tire A (i.e., gets outfitted with four type A tires), one gets B, one gets C, and one gets D.

**This design involves no replication. Without replication, you can't tell whether any difference in wear is due to tire type or to car differences.**

2. Twelve Cadillacs of the same type are purchased new from four dealers. Three get tire A, three get B, three get C, and three get D.

You can't infer to all cars what you observe only on Cadillacs. This was true in design 1 as well.

3. Twelve vehicles of different types are randomly selected from a list of many vehicle types and then are randomly allocated into four groups of three. One group gets tire A, one group gets tire B, one group gets tire C, and one group gets tire D.

The differences in wear on the tires may be due to the types of car in the group and not the tire type.

4. Four Cadillacs, four Fords, and four Volkswagens are purchased. One of each type of car gets tire A, one gets tire B, one gets tire C, and one gets tire D.

The blocks in this design are the car types. If there is a difference in tire types, it would be most easily detected with this design.