

## **To Do List**

**Do Problems 1-20 in Chapter 1**

**Assignment 1 is due Friday  
September 23 at NOON.**

**Tutorial Test 1 is on Wednesday  
September 28 – see detailed  
instructions posted on Learn**

# Today's Class: Finish Graphical Summaries

- 1) Histograms
- 2) Empirical Cumulative Distribution Function (e.c.d.f.)
- 3) Boxplots
- 4) Run Charts
- 5) Scatterplots (**Bivariate Data**)
- 6) **Bar Charts, Pie Charts**

# Clicker Question 1

**Which one of the following statements is true for you?**

**A: I am a CS student and my hometown is in Canada.**

**B: I am not a CS student and my hometown is in Canada.**

**C: I am a CS student and my hometown is not in Canada.**

**D: I am not a CS student and my hometown I not in Canada.**

# CS/Hometown Data

**What type of data are we collecting in this example?**

# Bivariate Categorical Data

PROGRAM/ HOMETOWN	Canadian Hometown	Non-Canadian Home town	Total
Computer Science	35	43	78
Non-Computer Science	18	69	87
Total	53	112	165

**Does there appear to be a relationship between hometown and program?**

**Can't use a scatterplot!**

# **Bivariate Categorical Data**

**Proportion of CS students with Canadian hometown =  $35/78 = 0.448$ .**

**Proportion of Non-CS students with Canadian hometown =  $18/87 = 0.206$ .**

**The relative risk of a Canadian hometown among CS students as compared to non-CS students =  $(35/78) / (18/87) = 2.17$ .**

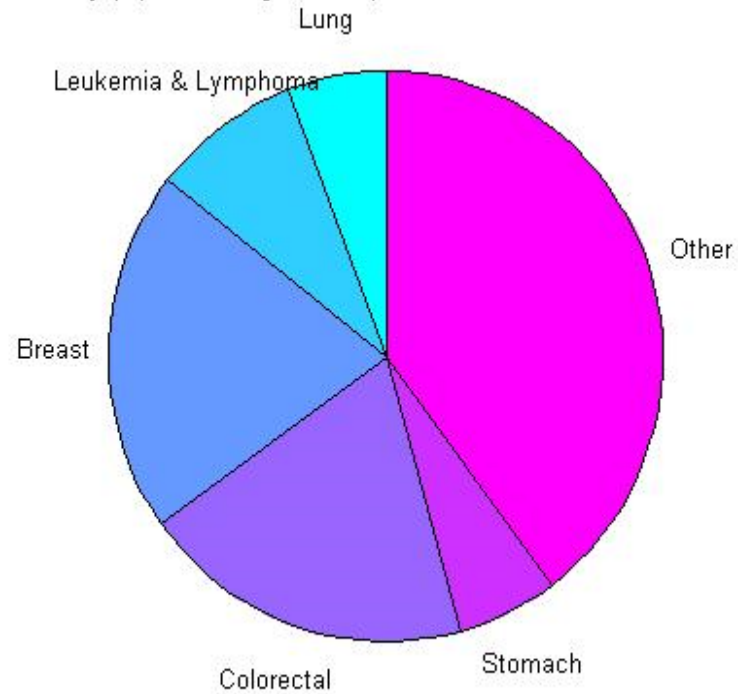
**Example: Mortality (%) from malignant neoplasms for females in Ontario, 1970 and 2000.**

<b>Females</b>	<b>1970</b>	<b>2000</b>
<b>Lung</b>	<b>5.8</b>	<b>21.6</b>
<b>Leukemia &amp; Lymphoma</b>	<b>8.6</b>	<b>9.0</b>
<b>Breast</b>	<b>20.5</b>	<b>17.4</b>
<b>Colorectal</b>	<b>19.5</b>	<b>12.6</b>
<b>Stomach</b>	<b>5.8</b>	<b>2.5</b>
<b>Other</b>	<b>39.8</b>	<b>37.0</b>

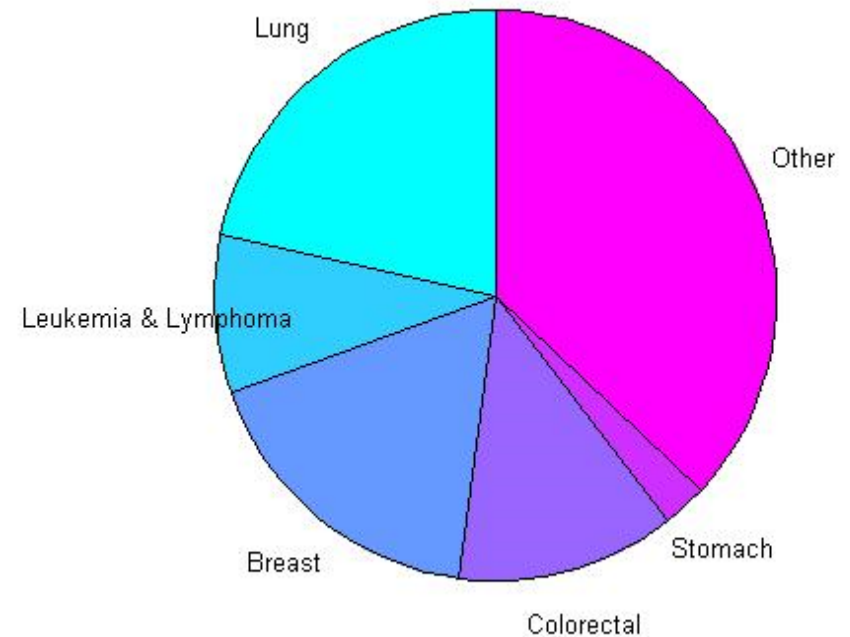
**How can we best display these data graphically?**

# Pie Charts

Mortality (%) from malignant neoplasms for females in Ontario 1970

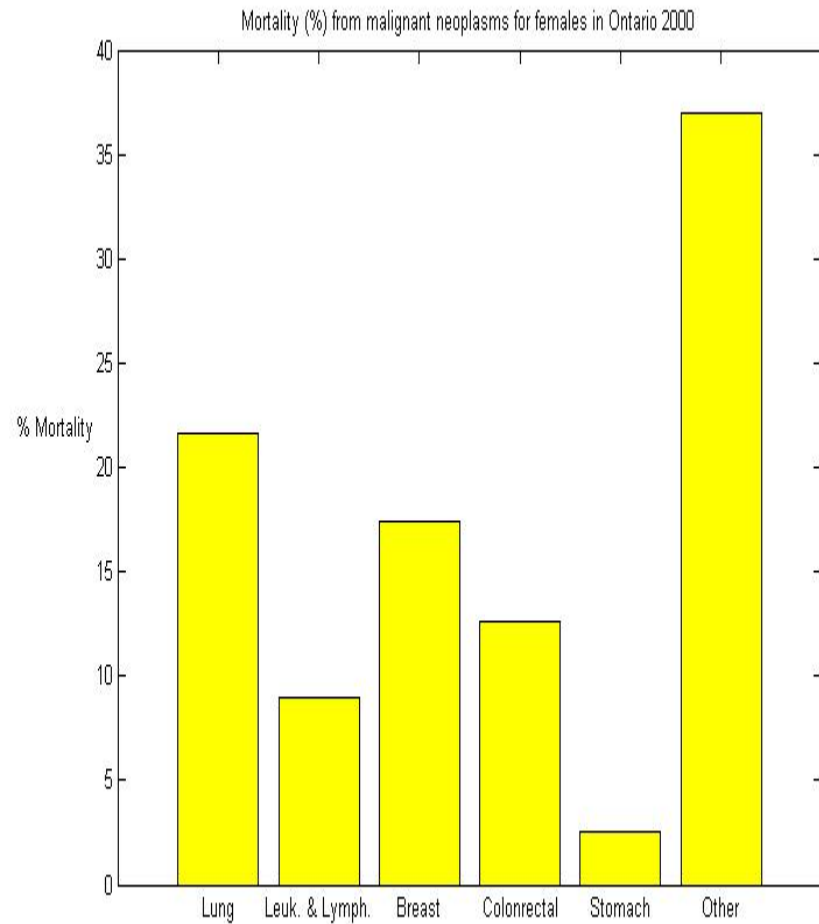
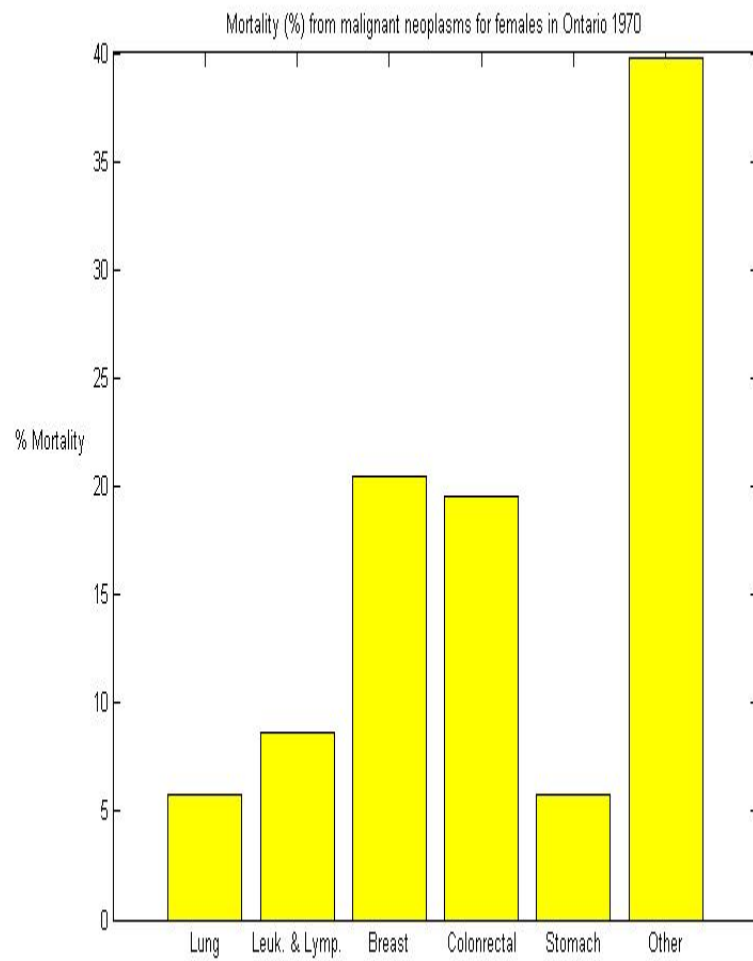


Mortality (%) from malignant neoplasms for females in Ontario 2000

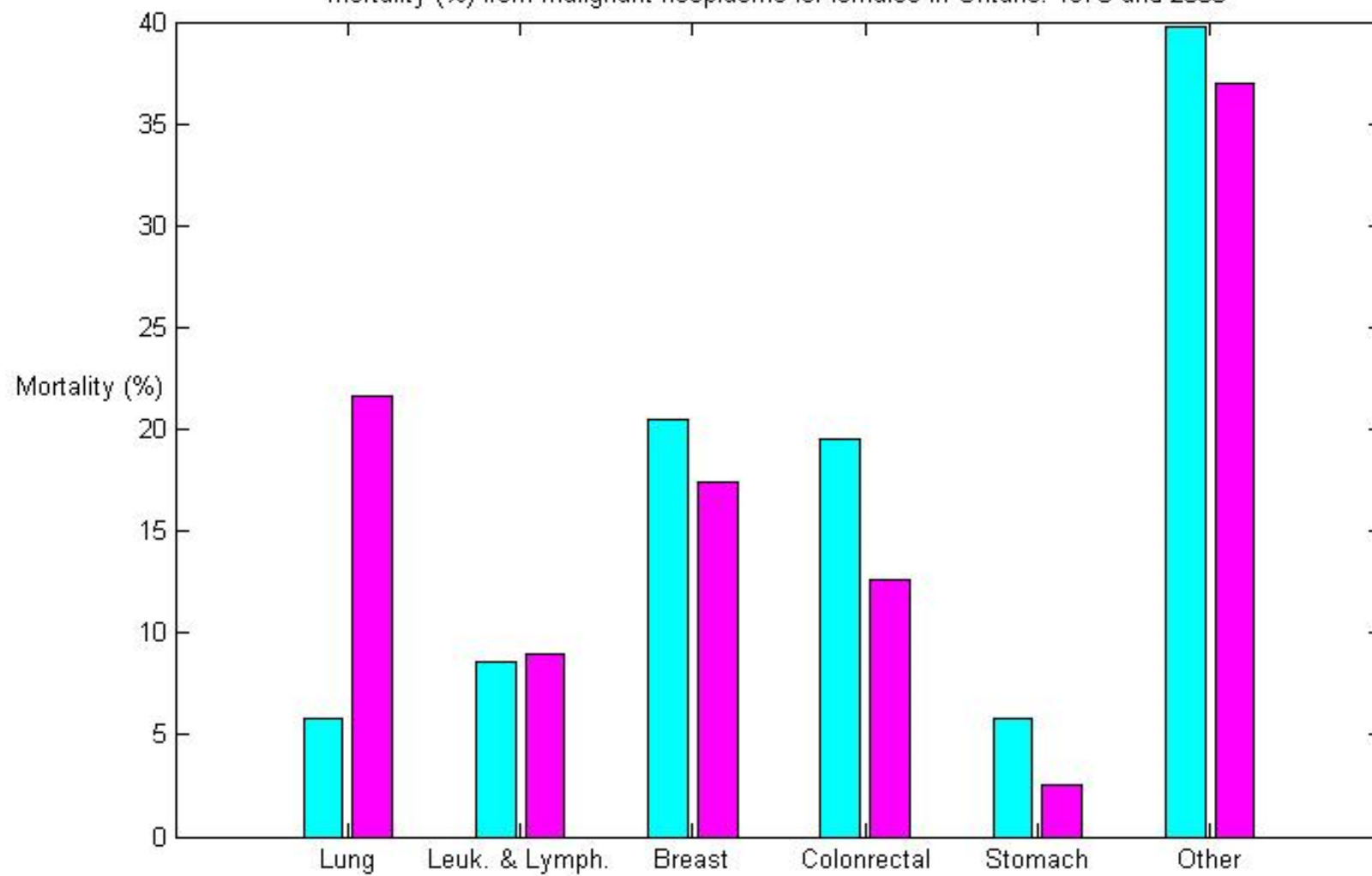




# Bar Charts



Mortality (%) from malignant neoplasms for females in Ontario: 1970 and 2000

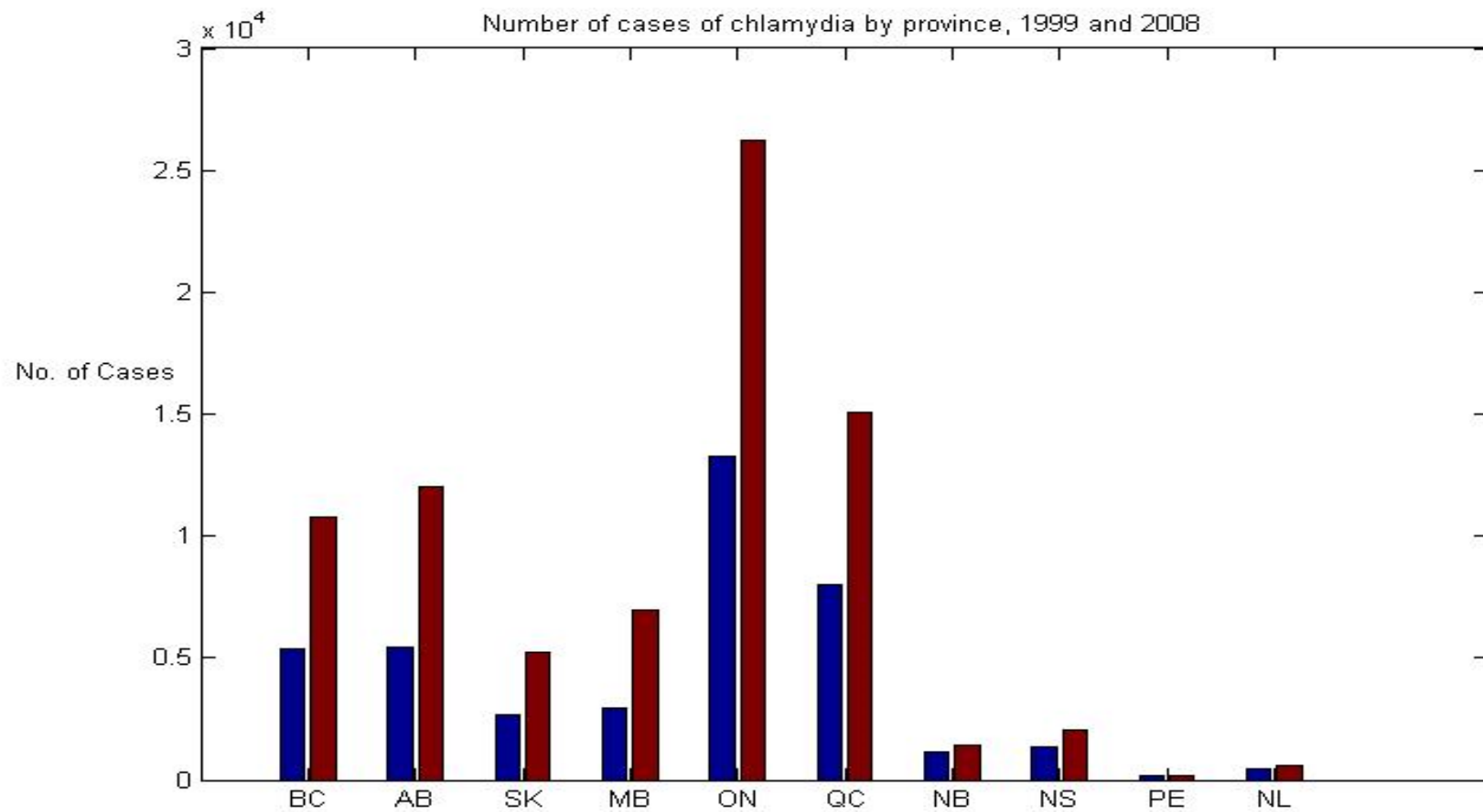


# Example: Cases and Rates of Chlamydia

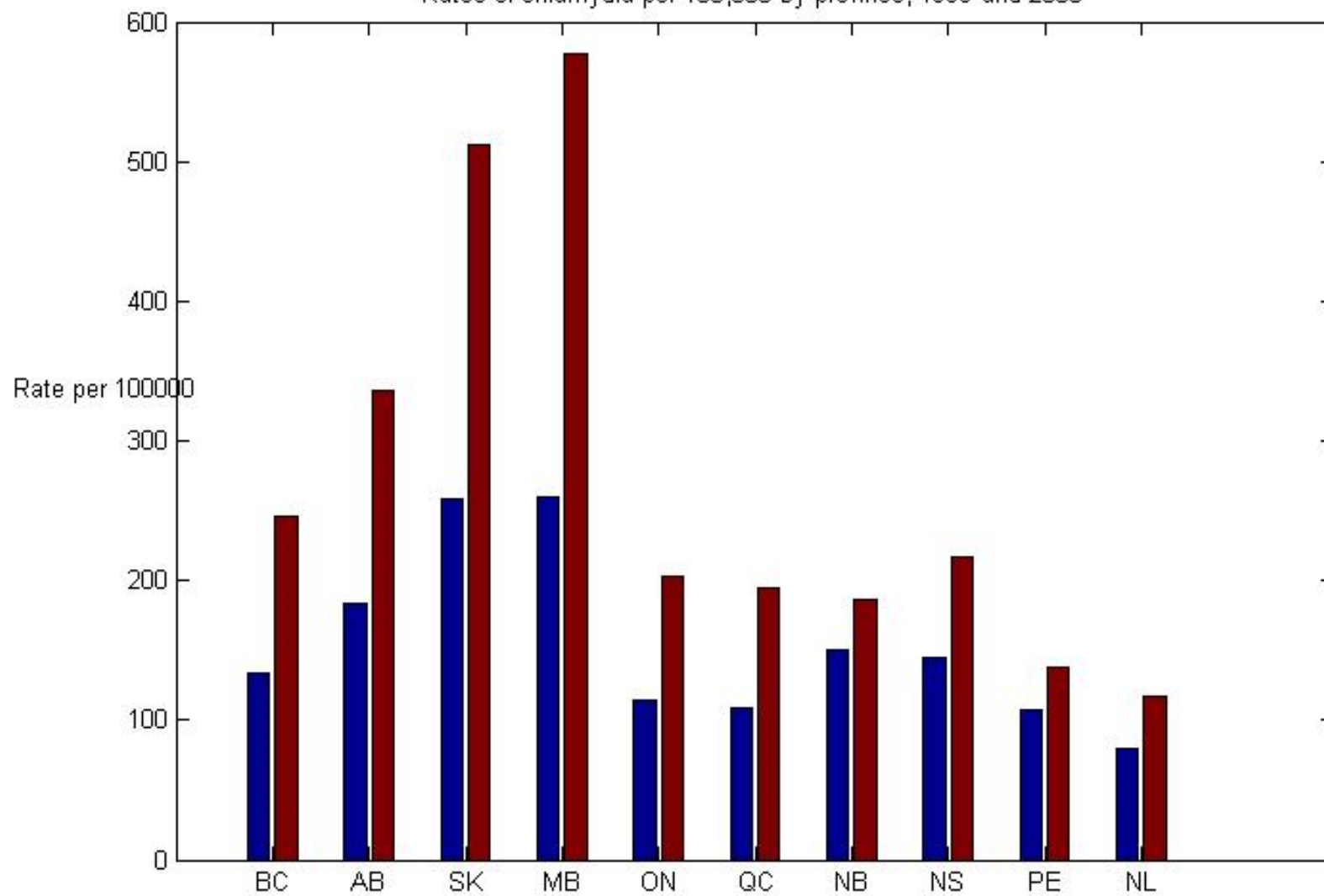
**Table 1: Reported Cases and Rates<sup>1</sup> of Chlamydia by Province/Territory, 1999 and 2008, Canada**

Jurisdiction	Number of Cases		Rates per 100,000 <sup>3</sup>		Rate Change (%)
	1999	2008	1999	2008	1999–2008
<i>Canada</i>	42,141	82,919	138.2	248.9	80.2
<b>BC</b>	5,402	10,766	134.1	245.7	83.2
<b>AB</b>	5,416	12,047	<b>183.0</b>	<b>336.0</b>	83.6
<b>SK</b>	2,656	5,203	<b>259.0</b>	<b>512.1</b>	97.7
<b>MB</b>	2,967	6,965	<b>259.7</b>	<b>576.6</b>	122.0
<b>ON</b>	13,256	26,245	115.0	203.0	76.5
<b>QC</b>	7,968	15,043	108.4	194.1	79.0
<b>NB</b>	1,136	1,389	<b>150.6</b>	185.9	23.4
<b>NS</b>	1,364	2,033	<b>145.1</b>	216.7	49.3
<b>PE</b>	148	193	107.6	138.0	28.3
<b>NL</b>	433	596	80.1	117.3	46.5
<b>YT</b>	176	232	<b>567.0</b>	<b>700.0</b>	23.4
<b>NT</b>	1,219	870	<b>1,796.4</b>	<b>2,010.0</b>	N/A
<b>NU<sup>2</sup></b>	N/A	1,337	N/A	<b>4,251.5</b>	N/A

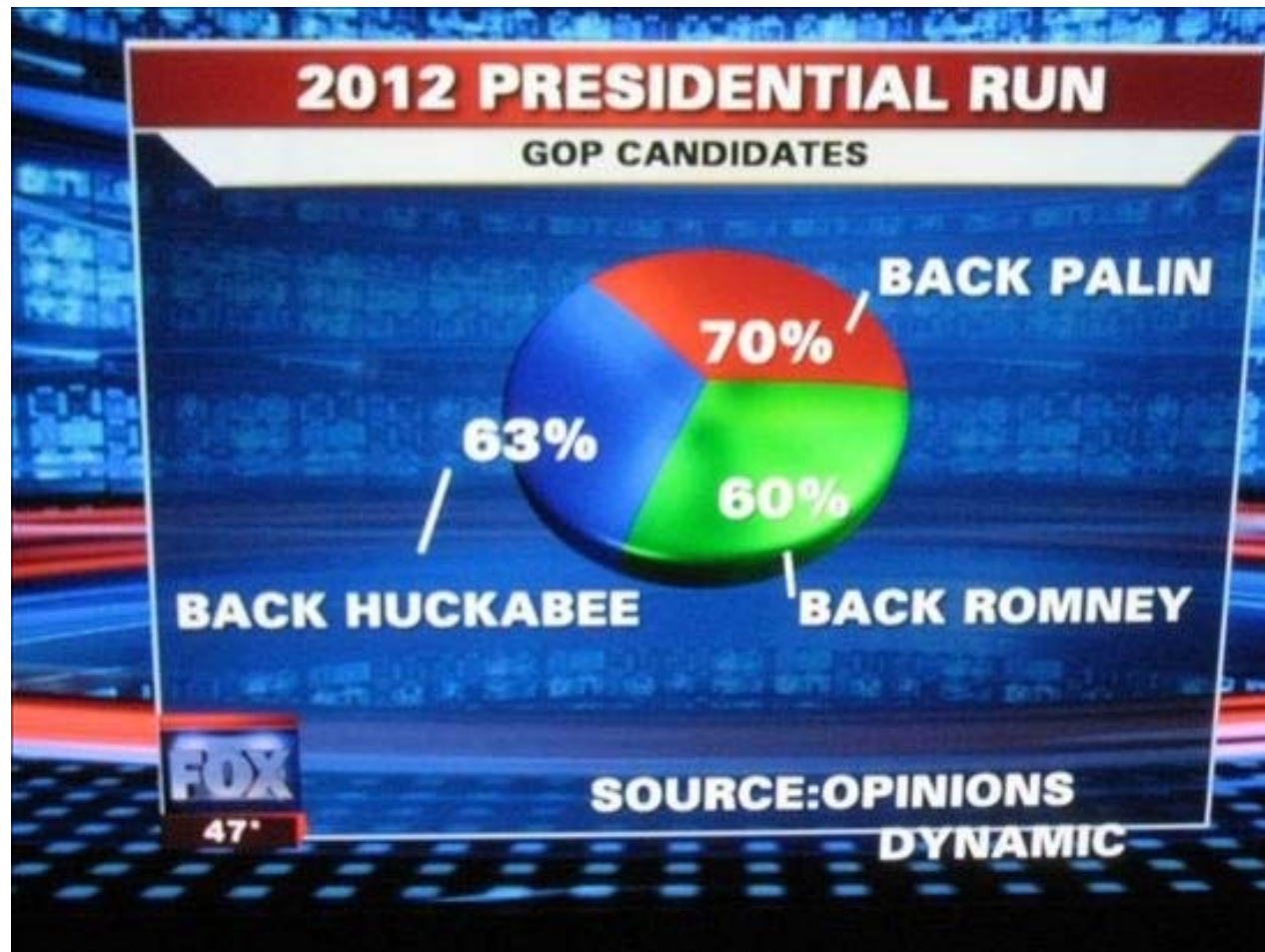
# Number of Cases by Province



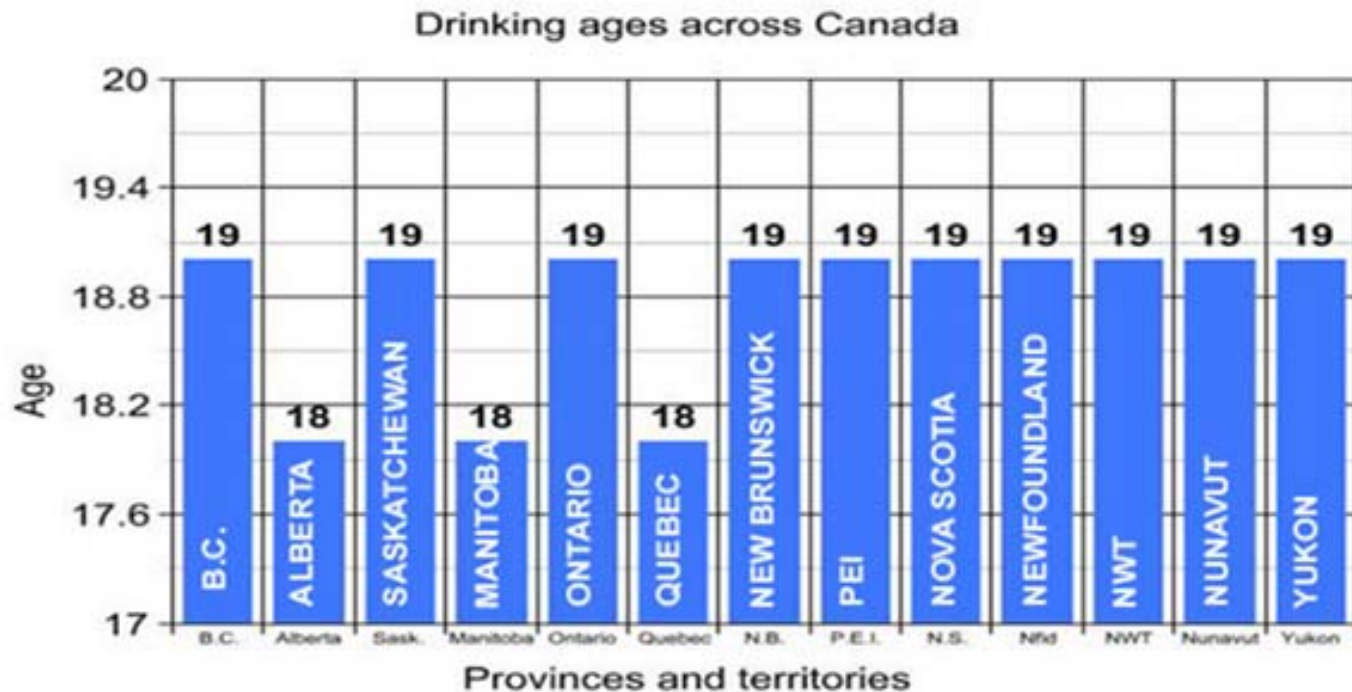
Rates of chlamydia per 100,000 by province, 1999 and 2008



# Not all graphical representations are good



# One of my favourites:



Canadian Centre on Substance Abuse

You have to be 19 in Saskatchewan to have a drink, while in Alberta and Manitoba, the drinking age 18. (CBC)

# **Remainder of Today's Class and Friday's Class:**

- 1) Descriptive Statistics and Statistical Inference  
(Inductive versus Deductive Reasoning)**
- 2) Types of Statistical Problems:**
  - (i) Estimation Problems**
  - (ii) Hypothesis Testing Problems**
  - (ii) Prediction Problems**
- 3) Statistical Models – Why and How to Choose**
- 4) Families of Statistical Models**
- 5) Unknown Parameters in a Statistical Model**
- 6) Estimates of Unknown Parameters in Statistical Models**



# **Data Analysis and Statistical Inference**

**We collect data to increase our knowledge or as a basis for making decisions.**

**Proper analysis of the data is crucial.**

**Two broad aspects of the analysis and interpretation of data are:**

- 1) Descriptive Statistics**
- 2) Statistical Inference**

# Descriptive Statistics

**Descriptive statistics** is the portrayal of the data, or parts of it, in numerical and graphical ways to show features of interest.

The numerical and graphical summaries we have examined are all examples of descriptive statistics.

A more complex example is “knowledge discovery” and “data mining” (KDD) which refers to exploratory data analysis of databases. The emphasis is on descriptive statistics for very large data bases. The goal is to find interesting patterns and relationships.

# Statistical Inference

**When the data obtained in the study of a population or process are used to draw general conclusions about the population or process itself we call this process **statistical inference**.**

# Inductive versus Deductive Reasoning

When we reason from the specific (the observed data on a sample of units) to the general (the target population or process) we call this **inductive reasoning**. Therefore statistical inference is a form of inductive reasoning.

Note: in mathematics when we use general results (axioms) to prove theorems this is called **deductive reasoning**.

# Methods of Statistical Inference

The methods of statistical inference that we will look at in STAT 231 will be used to examine 3 main types of problems:

- 1) Estimation problems
- 2) Hypothesis testing problems
- 3) Prediction problems

# Estimation Problems

In an **estimation problem** we are interested in estimating one or more attributes of a process or population.

## Examples:

- 1) Estimate the proportion of Ontario residents aged 14 - 20 who smoke.
- 2) Estimate the distribution of survival times for certain types of AIDS patients.
- 3) “Fit” or select a probability distribution for a process.

# Hypothesis Testing Problems

In a **hypothesis testing problem** we use the data to assess the truth of some question or hypothesis.

## Examples:

- 1) Is it true that in the 14-20 age group a higher proportion of females than males smoke?
- 2) Will the use of a new treatment with unpleasant side effects increase the average survival time of AIDS patients by at least 20 percent?

# Prediction Problems

In a **prediction problem** we use the data to predict a future value of a variate for a unit to be selected from the population or process.

## Examples:

- 1) Given the past performance of a stock and other data, predict the value of the stock at some point in the future.
- 2) Based on the results of a clinical trial predict how much an individual's blood pressure would drop for a given dosage of a drug.



# Statistical Analyses

**See other examples in Section 1.5 of the Course Notes.**

**We will be discussing methods of estimation, tests of hypotheses, and methods for prediction in much more detail in future lectures.**

# What we have done so far:

We have looked at **numerical** and **graphical summaries** of a dataset both univariate  $\{y_1, y_2, \dots, y_n\}$  and bivariate  $\{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$ .

We will use these summaries to help us chose a **statistical model** for the data.

A statistical model is a mathematical model that incorporates probability.

# Why are Statistical Models Useful?

**Statistical models can be used to describe many different processes such as:**

- **the daily closing value of the Canadian dollar**
- **the number of hits on a government website describing how to apply for unemployment insurance**
- **catastrophic events such as oil spills, nuclear meltdowns, epidemics, diseases etc.**

# Why are Statistical Models Useful?

**Statistical models can also be used to describe:**

- **the selection of units from a population or process (Labour Force Survey)**
- **the measurement of variates**
- **the variability in the variate values for a population or process**

# Why are Statistical Models Useful?

- Drawing conclusions from data involves some degree of uncertainty. Statistical models can be used to quantify this uncertainty.
- Often, questions of interest can be formulated in terms of parameters of the statistical model.  
e.g.  $\text{Poisson}(\theta)$ ,  $\theta$  = mean number of events

# Why use a Statistical Model?

- **Procedures for making decisions based on observed data can often be formulated in terms of statistical models.**
- **Models allow us to characterize processes and to simulate them via computer experiments.**