

STAT206

Andy Zhang

Fall 2014

# Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
1.1	Statistics . . . . .	2
1.1.1	Definitions . . . . .	2
1.1.2	Process . . . . .	2
1.1.3	Data Types . . . . .	3
1.1.4	(Grouped) Frequency Tables . . . . .	3
1.1.5	Stem and Leaf Plot . . . . .	3
1.1.6	Bar Chart . . . . .	3
1.1.7	Histogram . . . . .	4
1.1.8	Measures of Centrality . . . . .	4
1.1.9	Measures of Variability . . . . .	4
1.1.10	Box Plot . . . . .	4
<b>2</b>	<b>Probability</b>	<b>5</b>
2.1	Definitions . . . . .	5
<b>3</b>	<b>Random Variables</b>	<b>7</b>
3.1	Definitions . . . . .	7
<b>4</b>	<b>Discrete Probability Distributions</b>	<b>8</b>
4.1	Definitions . . . . .	8
<b>5</b>	<b>Continuous Probability Distributions</b>	<b>10</b>
5.1	Definitions . . . . .	10
<b>6</b>	<b>Normal Distribution</b>	<b>12</b>
6.1	Definitions . . . . .	12
<b>7</b>	<b>Confidence Intervals</b>	<b>14</b>
7.1	Definitions . . . . .	14
<b>8</b>	<b>Confidence Interval II</b>	<b>16</b>
8.1	Definitions . . . . .	16
8.1.1	Words . . . . .	16
8.1.2	Methods . . . . .	16
8.2	Confidence Intervals . . . . .	17
8.2.1	Difference of Means for large samples . . . . .	17
8.2.2	Difference of Means for small samples . . . . .	17

8.2.3	Difference of means, not independent . . . . .	18
8.2.4	Exact Confidence Intervals . . . . .	18
8.2.5	Approximate Confidence Intervals . . . . .	18
<b>9</b>	<b>Hypothesis Testing</b>	<b>19</b>
9.1	Steps . . . . .	19
9.2	Difference between One sided and Two sided . . . . .	19
9.3	Formulas . . . . .	19
9.4	Contingency tables . . . . .	20

# Chapter 1

## Introduction

### 1.1 Statistics

#### 1.1.1 Definitions

**Statistics** Collection, organization, analysis, interpretation and presentation of data. It is also defined as the quantification of uncertainty.

**Unit** A single element, usually a person or object, whose characteristics are of interest. Ex: A student enrolled in the course.

**Population** The set of all units which are of interest. Ex: All students enrolled in the course

**Variable** A measurement of the characteristic of interest from a unit. Ex: Number of Canadian provinces visited by a student

**Sample** A subset of units from the population for which measurements of the desired variable are actually made. Ex: 29 students chosen from the class

**Descriptive Statistics** Summarize the data in the sample, both graphically and numerically

**Inferential statistics** Use the sample data to estimate an attribute of the population. Include a quantification of uncertainty

**Sampling Error** An error which occurs due to the uncertainty in randomly selecting a sample.

**Study error** A systematic error which occurs because the sample does not accurately represent the population

#### 1.1.2 Process

Identify the problem of interest

- Who or what do you want to learn about?
  - Define the **population** of interest

- Individual elements of the population are called **units**
- What research question would you like answered?
  - Define your **hypothesis**

Plan the data collection

- How will you select a subset of **units** from the **population** to be in your **sample**?
  - How large will the **sample** be?
- What is (are) the **variable (s)** of interest?
  - How will you measure it (them)?

Analyze the data

- Graph the data — histogram, scatter-plot, etc
- Compute **Descriptive statistics** — e.g. sample mean, sample variance, etc.
- Compute **Inferential statistics** — e.g. confidence intervals, hypothesis tests about population **parameters**
  - Inferential statistics include a quantification of the sampling error

Draw conclusions

- Use the results of your analysis to address the original research question
- Address limitations of the study, especially any potential systematic **study errors**

### 1.1.3 Data Types

**Categorical Variable** A qualitative measure. Each unit belongs to **one of K** possible classes.

**Discrete variable** A quantitative measure. Each unit's measurement can take on one of a **countable** number of possible values

**Continuous variable** A quantitative measure. Each unit's measurement can take on an **uncountable** number of possible values, usually some interval of real numbers

### 1.1.4 (Grouped) Frequency Tables

- Display the number of units which are in each class
- Discrete / Continuous variables are grouped into classes
- In the case of numerical variables, there is a loss of information

See more: [http://en.wikipedia.org/wiki/Stem-and-leaf\\_display](http://en.wikipedia.org/wiki/Stem-and-leaf_display)

### 1.1.5 Stem and Leaf Plot

- A **stem-and-leaf plot** is a way to summarize a relatively **small** data set, without the loss of information that occurs with a frequency table
- Left is possible **first** digits, right is remaining digits in ascending order

See more: [http://en.wikipedia.org/wiki/Stem-and-leaf\\_display](http://en.wikipedia.org/wiki/Stem-and-leaf_display)

### 1.1.6 Bar Chart

- Bar charts are used to graphically display information from categorical variables

See more: [http://en.wikipedia.org/wiki/Bar\\_chart](http://en.wikipedia.org/wiki/Bar_chart)

### 1.1.7 Histogram

- A histogram is similar to a bar chart, but it's for numerical data
- The range is divided in distinct classes, and each observation is assigned to exactly one class
- Histogram shows frequency of observations in each class

See more: <http://en.wikipedia.org/wiki/Histogram>

- If class ranges are not same length, we can use density histogram instead
- When interpreting a density histogram, it is the area that is meaningful
- Height is  $height = \frac{relative\ frequency}{width} = \frac{frequency}{width * n}$

See more <http://en.wikipedia.org/wiki/Histogram>

### 1.1.8 Measures of Centrality

- The **sample mean** of a set of  $n$  values,  $x_1, x_2, x_3, \dots, x_n$  denoted by  $\bar{x}$  is  $\bar{x} = \frac{\sum_{i=1}^n x_i}{n}$
- The **median** is the number  $x^*$  such that half of the observed values are below  $x^*$  and half are above
- If after writing our values in ascending order, we denote the  $i^{th}$  value as  $x_{(i)}$ , then

$$x^* = \begin{cases} x_{(\frac{n+1}{2})} & \text{if } n \text{ is odd} \\ x_{(\frac{n}{2})} + x_{(\frac{n+2}{2})} & \text{if } n \text{ is even} \end{cases}$$

### 1.1.9 Measures of Variability

Measures of variability

- The **sample variance** of a set of values  $x_1, x_2, x_3, \dots, x_n$  denoted by  $s^2$  is

$$s^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1}$$

- The **sample standard deviation** denoted  $s$ , is the square root of the sample variance
- The **range** of the set is the difference between the maximum and minimum value

$$range = x_{(n)} - x_{(1)}$$

### 1.1.10 Box Plot

- The box indicates the middle 50% of the observations, i.e. the second and third quartiles
- The line through the box indicates the median observation
- The whiskers indicate the highest and lowest observations

See more: [http://en.wikipedia.org/wiki/Box\\_plot](http://en.wikipedia.org/wiki/Box_plot)

# Chapter 2

## Probability

### 2.1 Definitions

**Probability** measure the uncertainty associated with an event. An event is something that might occur

- Classical:  $\frac{\text{Number of ways event can occur}}{\text{Total number of equally likely outcomes}}$
- Relative Frequency: Proportion of times the event occurs, as the number of trials approaches infinity
- Subjective: Estimates of probability that the event occurs, based on subjective opinion

**Experiment** is a repeatable phenomenon or process

**Trial** is a single repetition of an experiment

**Sample Space** ,  $S$ , is the set of distinct outcomes for an experiment or process

**Discrete** A sample space is discrete if it has a finite or countably infinite number of simple events. Otherwise it is non discrete or continuous

**Mutually Exclusive** means two events never occur simultaneously

**Complement** of an event and an event are always mutually exclusive

**Uniform distribution** The total probability is uniformly distributed among all possible outcomes

**Permutation** is the number of ways to arrange  $r$  out of  $n$  objects:  $n^{(r)} = \frac{n!}{(n-r)!} = n(n-1)(n-2)\dots(n-r+1)$

**Combinations** If we don't care about the order of objects, but just which objects are chosen, the number of ways to choose  $r$  out of  $n$  items is  $\binom{n}{r} = \frac{n!}{r!(n-r)!}$

### Set Operations

- $AB$  or  $A \cap B$  is the intersection of two events.
- $A \cup B$  is the union of two events.
- $\bar{A}$  is the complement of  $A$ , not event  $A$

**Conditional** The probability of event  $A$ , conditional on the occurrence of event  $B$ , denoted by  $P(A|B)$  is  $P(A|B) = \frac{P(A \cap B)}{P(B)}$ ,  $P(B) \neq 0$

**Independent** Two events are said to be independent iff  $P(A \cap B) = P(A)P(B)$ . This implies that  $P(A|B) = P(A)$ ,  $P(B|A) = P(B)$ . In other words, events  $A$  and  $B$  are independent if whether  $B$  occurs does not influence whether  $A$  occurs, and vice versa

**Bayes' Theorem** Suppose  $A$  and  $B$  are any two events in  $S$ , then  $P(B|A) = \frac{P(A|B)P(B)}{P(A|B)P(B) + P(\bar{A}|B)P(B)}$



# Chapter 3

## Random Variables

### 3.1 Definitions

**Random Variable**  $X$  is a function from the sample space  $S$  to the real numbers  $X : S \rightarrow \mathbb{R}$ . Its range  $R(X)$  is the set of possible real values it can take. We use  $X$  to denote a random variable and  $x$  to denote its observed value

- Discrete if takes on a countable number of possible values
- Continuous if it takes on all values in some interval of the real line

**Indicator or binary** variables take values of 0 or 1

**Probability Function (pf)** of  $X$ , denoted  $f(x)$ , denotes the probability that  $X$  takes on the value  $x$ .  $f(x) = P(X = x)$  defined for all  $x$  in the range of  $X$

**Probability Distribution** The set of pairs  $\{(x, f(x)) | x \in R(X)\}$  is called the probability distribution of  $X$

- $0 \leq f(x) \leq 1$  for any  $x$
- $\sum f(x) = 1$

**Cumulative distribution function (cdf)** of  $X$ , denoted  $F(x)$  denotes the probability that  $X$  takes on a value  $\leq x$ :  $F(x) = P(X \leq x)$ . Very useful for continuous random variables

**Mean or Expected Value** of  $X$  is defined as  $\mu = E(X) = \sum_x x f(x)$ . We can understand  $E(X)$  as the average value that  $X$  would assume over a theoretically infinite number of trials.  $E(X)$  is not a random variable, it is a constant.

We also define the expected value as  $E[g(X)] = \sum_x g(x) f(x)$

**Variance** of a random variable  $X$  is the expected squared difference from the mean, that is  $Var(X) = E[(X - E(X))^2] = \sum_{all x} f(x)(x - \mu)^2$ . An alternative would be  $Var(X) = E[X^2] - (E(X))^2$

# Chapter 4

## Discrete Probability Distributions

### 4.1 Definitions

**Bernoulli Distribution** Repeated trials of an experiment

- Each trial can be a success or a failure
- The probability of a success is the same for each trial
- The outcomes of different trials are **independent**
- Let  $X$  record success or failure

We say that  $X$  follows a **Bernoulli distribution** ( $X \sim \text{Bernoulli}(p)$ ), where  $p$  is the probability of success

$$f(x) = \begin{cases} p & \text{if } x = 1 \\ (1 - p) & \text{if } x = 0 \end{cases}$$

$$E(X) = p$$

$$\text{Var}(X) = E[X^2] - (E(X))^2 = p - p^2 = p(1 - p)$$

**Binomial Distribution** Physical setup: We perform a sequence of  $n$  independent Bernoulli trials

- Each trial has two possible outcomes: success or failure
- Trials are independent
- Each trial has probability of success equal to  $p$

$$f(x) = \binom{n}{x} p^x (1 - p)^{n-x}, \quad x = 0, 1, 2, \dots, n$$

$$E(X) = np$$

$$\text{Var}(X) = np(1 - p)$$

**Poisson Process** Physical setup: Events occur randomly in time (or space) according to the following conditions

- Independence: The number of occurrences in disjoint (non overlapping) intervals are independent
- Individuality: Events occur singly i.e  $P(\text{two or more events occur simultaneously}) = 0$
- Homogeneity: Events occur according to a uniform (constant) rate or intensity ( $\lambda$ )

If events occur with an average rate of  $\lambda$  per unit of time and  $X$  is the number of events which occur in  $t$  units of time, then  $X \sim \text{Poisson}(\lambda t)$

$$f(x) = \frac{e^{-\lambda t} (\lambda t)^x}{x!}, \quad x = 0, 1, \dots$$

$$E(X) = \lambda t$$

$$Var(X) = \lambda t$$

**Hypergeometric Distribution** Physical setup: We have a collection of  $N$  objects which can be classified into two distinct types, called success and failure. There are  $r$  and  $N - r$  failures. A sample of  $n$  objects is selected without replacement.

Let  $X$  be the number of successes selected, then  $X$  is said to follow a hypergeometric distribution ( $X \sim Hyper(N, r, n)$ ). To compute the probability function, note that, for  $X = x$

- There are  $\binom{N}{n}$  points in the sample space (if we do not consider the order of selection)
- There are  $\binom{r}{x}$  ways too select the  $x$  success objects from the  $r$  available
- There are  $\binom{N-r}{n-x}$  ways to select the remaining  $n - x$  failure objects from the  $N - r$  available

$$f(x) = \frac{\binom{r}{x} \binom{N-r}{n-x}}{\binom{N}{n}}, x = 0, 1, \dots, \min(r, n)$$

$$E(X) = \frac{nr}{N}$$

$$Var(X) = \frac{nr(N-r)(N-n)}{N^2(N-1)}$$

**Geometric Distribution** Physical Setup: Bernoulli Trials are repeated until the first success. Let  $X$  be the number of independent Bernoulli ( $p$ ) trials until the first success (including the first success), then  $X$  follows a geometric distribution,  $X \sim Geom(p)$

$$f(x) = p(1 - p)^{x-1}, x = 1, 2, \dots$$

$$E(X) = \frac{1}{p}$$

$$Var(X) = \frac{1-p}{p^2}$$

# Chapter 5

## Continuous Probability Distributions

### 5.1 Definitions

**Continuous Random Variable**  $X$  is a function from the sample space to the real numbers:  
 $X : S \rightarrow \mathbb{R}$

The range  $R(X)$  is continuous. Individual points in  $\mathbb{R}$  must have a 0 probability since the interval length is 0

**Probability Density Function (pdf)** of a continuous random variable  $X$ , denoted  $f(x)$  assigns a probability to an  $x \in R(X)$ . If we have the probability density function for  $X$ , then we can define the probability that  $X$  takes a value in an interval  $(a, b) \subseteq R(X)$  as  $P(A < X < b) = \int_a^b f(x)dx, (a, b) \subseteq R(X)$

Properties of the probability density function:

$$f(x) \geq 0, \forall x \in R(X)$$

$$\int_{x \in R(X)} f(x)dx = 1$$

**Cumulative Distribution Function** of a continuous random variable  $X$ , denoted  $F(x)$ , gives the probability that  $X$  takes on a value less than or equal to  $x$

$$F(x) = P(X \leq x) = P(X < x)$$

Properties

- $F(-\infty) = 0$
- $F(\infty) = 1$
- $F(x)$  is non decreasing

**Continuous Uniform Distribution** Physical setup: The probability of any subinterval of the range is proportional to the length of the interval. For  $a < b$ , if  $X$  is uniformly distributed on the interval  $(a, b)$  then we write  $X \sim U(a, b)$

$$f(x) = \frac{1}{b-a}, a \leq x \leq b$$

$$F(x) = \begin{cases} 0 & \text{if } x < a \\ \frac{x-a}{b-a} & \text{if } a \leq x \leq b \\ 1 & \text{if } b < x \end{cases}$$

**Mean or Expected Value** of a continuous random variable  $X$  is defined as

$$E(X) = \int_{-\infty}^{\infty} xf(x)dx = \mu$$

Properties:

- $E[g(X)] = \int_{-\infty}^{\infty} g(x)f(x)dx$
- $E[aX + bY] = aE(X) + bE(Y)$

**Variance** is the expected squared difference from the mean, that is  $Var(X) = E[(X - E(X))^2]$ . If  $X_1$  and  $X_2$  are independent random variables, and  $a, b \in \mathbb{R}$ , then  $Var[aX_1 + bX_2] = a^2Var(X_1) + b^2Var(X_2)$

**Exponential Distribution** Physical setup: Events occur according to a Poisson process, and we measure the inter arrival times between events. If  $X$  is the amount of time until the next event in a Poisson process, then  $X \sim Exp(\theta)$ , where  $\theta = \frac{1}{\lambda}$

$$f(x) = \left(\frac{1}{\theta}\right)e^{-\frac{x}{\theta}}, x > 0$$

$$F(x) = 1 - e^{-\frac{x}{\theta}}, x > 0$$

$$E(X) = \theta$$

$$Var(X) = \theta^2$$

**Simulation** The most common use of the continuous uniform distribution to simulate other random variables.

**Theorem:** If  $F(x)$  is an arbitrary cdf, and  $Y \sim U(0, 1)$ , then  $X = F^{-1}(Y)$  has cdf  $F(x)$ .

We can use  $Y$  to generate an observation of the random variable  $X$ :

- Generate an observation  $y$  from  $Y \sim U(0, 1)$  using your favourite software
- Compute  $x = F^{-1}(y)$

# Chapter 6

## Normal Distribution

### 6.1 Definitions

**Normal Distribution** A continuous random variable  $X$  with range  $(-\infty, \infty)$  has a normal distribution denoted  $X \sim N(\mu, \sigma^2)$ , if its pdf has the form  $f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2}(\frac{x-\mu}{\sigma})^2}$ ,  $x \in \mathbb{R}$  where the mean  $\mu$  and variance  $\sigma^2$  are parameters.  
 $E(X) = \mu$ ,  $Var(X) = \sigma^2$

**Properties** Let  $X_1 \sim N(\mu_1, \sigma_1^2)$  and  $X_2 \sim N(\mu_2, \sigma_2^2)$  be independent  
 $Y = aX_1 + bX_2 + c \sim N(a\mu_1 + b\mu_2 + c, a^2\sigma_1^2 + b^2\sigma_2^2)$   
If  $X \sim N(\mu, \sigma^2)$ , then  $Z = (\frac{1}{\sigma})X - (\frac{\mu}{\sigma}) = \frac{X-\mu}{\sigma} \sim N(0, 1)$   
We also have  $P(Z > z) = P(Z < -z)$  for any  $z \in \mathbb{R}$

**Central Limit Theorem** : We use the normal distribution to approximate probabilities for non normal distributions. This is possible because the normal distribution tends to approximate sums of random variable. Although this is a Theorem about limits, we will use it when  $n$  is large, but finite to approximate the distribution of  $\sum_i X_i$ , or  $\bar{X}$  by a normal distribution

**Independent** We say that  $X$  and  $Y$  are independent if, for all  $x$  and  $y$ , we have  
 $f(x, y) = P(X = x \cap Y = y) = P(X = x)P(Y = y) = f_X(x)f_Y(y)$

**Central Limit Theorem — Sum** Let  $X_1, X_2, X_3, \dots, X_n$  be independent variables all having the same distribution with  $E(X_i) = \mu$  and  $Var(X_i) = \sigma^2$

As  $n \rightarrow \infty$ , the cumulative distribution function of the random variable  $\sum_i X_i$  approaches the cumulative distribution function for  $N(n\mu, n\sigma^2)$

The cumulative distribution function of the random variable  $\frac{\sum_i X_i - n\mu}{\sigma\sqrt{n}}$  approaches the cumulative distribution function for  $N(0, 1)$

**Central Limit Theorem — Average** Let  $X_1, X_2, X_3, \dots, X_n$  be independent variables all having the same distribution with  $E(X_i) = \mu$  and  $Var(X_i) = \sigma^2$

As  $n \rightarrow \infty$ , the cumulative distribution function of the random variable  $\bar{X}$  approaches the cumulative distribution function for  $N(\mu, \frac{\sigma^2}{n})$

The cumulative distribution function of the random variable  $\frac{\bar{X}-\mu}{\frac{\sigma}{\sqrt{n}}}$  approaches the cumulative distribution function for  $N(0, 1)$

**Continuity Correction** can improve the approximation to a sum or average of discrete random variables using a normal random variable. We think of the center of a bar with width 1 as an integer value and the bar actually covering  $(x - 0.5, x + 0.5)$ . So instead of integrating from  $(0, 5)$  for example, we integrate on  $(-0.5, 5.5)$

# Chapter 7

## Confidence Intervals

### 7.1 Definitions

**Introduction to Estimation** Suppose that a probability distribution which serves as a model for some random process depends on an unknown parameter  $\theta$ . In order to use the model we have to estimate or specify a value for  $\theta$  using some data sets collected for the random variable.

**Estimate** of a parameter  $\theta$  is the value of a function of the observed data  $y_1, y_2, \dots, y_n$  and other known quantities such as the sample size  $n$ .

**Likelihood function** for  $\theta$  is defined as  $L(\theta) = L(\theta; y) = P(Y = y; \theta)$  for  $\theta \in \omega$  where the parameter space  $\omega$  is the set of possible values for  $\theta$

Suppose that  $\theta$  is the success probability in a binomial model, so that  $Y \sim Bi(n, \theta)$ . Suppose that we ran the experiment once and recorded  $y$  successes in  $n$  trials. Then  $L(\theta) = P(Y = y; \theta) = \binom{n}{y} \theta^y (1 - \theta)^{n-y}$

- Hold  $y$  constant and vary  $\theta$ .
- This makes  $L(\theta)$  into a function of  $\theta$ .
- Then we can use calculus to choose  $\theta$  to maximize  $L(\theta)$
- This choice for  $\theta$  is what we call  $\hat{\theta}$ , the maximum likelihood estimate for  $\theta$ .
- Here we would get  $\hat{\theta} = \frac{y}{n}$  which agrees with our intuition

**Maximum Likelihood Estimate** The value of  $\theta$  which maximizes  $L(\theta)$  for given data  $y$  is called the maximum likelihood estimate of  $\theta$ . It is denoted as  $\hat{\theta}$ .

**Relative Likelihood Function** is defined as  $R(\theta) = \frac{L(\theta)}{L(\hat{\theta})}$  for  $\theta \in \omega$ .

Note that  $0 \leq R(\theta) \leq 1$  for all  $\theta \in \omega$

**Log Likelihood Function** is  $l(\theta) = \ln L(\theta)$  for  $\theta \in \omega$ . Note that  $\hat{\theta}$  maximizes  $R(\theta)$  and  $l(\theta)$



**Estimator** We call  $\tilde{\theta}$  the estimator of  $\theta$  corresponding to  $\hat{\theta}$ . We will always use

- $\hat{\theta}$  to denote an estimate, that is, a numerical value
- $\tilde{\theta}$  to denote the corresponding estimator

An estimator  $\tilde{\theta}$  is a random variable which is a function  $\tilde{\theta} = g(Y_1, \dots, Y_n)$  of the random variables  $Y_1, \dots, Y_n$ . The distribution of  $\tilde{\theta}$  is called the sampling distribution of the estimator.

List of estimators:

- $Y_i \sim \text{Bernoulli}(p)$ ,  $\tilde{p} = \bar{Y}$
- $Y_i \sim \text{Poisson}(\lambda)$ ,  $\tilde{\lambda} = \bar{Y}$
- $Y_i \sim \text{Exponential}(\theta)$ ,  $\tilde{\theta} = \bar{Y}$
- $Y_i \sim \text{Normal}(\mu, \sigma^2)$ ,  $\tilde{\mu} = \bar{Y}$
- $Y_i \sim \text{Normal}(\mu, \sigma^2)$ ,  $\tilde{\sigma}^2 = \bar{s}^2$

**Unbiased** An estimator is said to be unbiased if its expected value equals the parameter being estimated:  $E(\tilde{\theta}) = \theta$ .

The standard deviation of an estimator is called its standard error:  $SE(\tilde{\theta}) = \sqrt{\text{Var}(\tilde{\theta})}$

If we have two unbiased estimators for a parameter, the one with the smaller standard error is preferred.

**Interval Estimate** for  $\theta$  based on observed data  $y$  takes the form  $[L(y), U(y)]$ . We assume that the probability model chosen is correct and that  $\theta$  is the true value of the parameter.

**Coverage Probability** To quantify the uncertainty in the interval estimate, we define the coverage probability for the interval estimator  $[L(Y), U(Y)]$  as  $C(\theta) = P(L(Y) \leq \theta \leq U(Y))$

**Confidence Intervals** A 100p% confidence interval for a parameter  $\theta$  is an interval estimate  $[L(y), U(y)]$  for which  $P(L(Y) \leq \theta \leq U(Y)) = p$ . We say that we are 100p% confident that the true parameter is in the interval

**Pivotal Quantity** is a function of the data  $Y$  and the unknown parameter  $\theta$  such that the distribution of the random variable  $Q$  is completely known. We define it as:  $Z = \frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}} \sim N(0, 1)$ .

To compute a 100p% confidence interval for  $\mu$ , determine the threshold  $c$  such that  $P(Z \leq c) = 1 - (\frac{1-p}{2}) = \frac{p+1}{2}$ , where  $Z \sim N(0, 1)$ .

We then construct the interval  $\bar{x} \pm c \frac{\sigma}{\sqrt{n}}$

# Chapter 8

## Confidence Interval II

### 8.1 Definitions

#### 8.1.1 Words

**Treatment** procedure, machine or process being compared by our experiment

**Unit** object being studied by our experiment

**Response variable** a variable of interest whose value we measure during our experiment

**Explanatory Variable** a variable which influences the value of the response variate

#### 8.1.2 Methods

**Independent Sampling** Independent samples are those samples selected from the same population, or different populations, which have no effect on one another. That is, no correlation exists between the samples.

**Matched Pair Samples** Matched samples can arise in the following situations:

Two samples in which the members are clearly paired, or are matched explicitly by the researcher. For example, IQ measurements on pairs of identical twins.

Those samples in which the same attribute, or variable, is measured twice on each subject, under different circumstances. Commonly called repeated measures. Examples include the times of a group of athletes for 1500m before and after a week of special training; or the milk yields of cows before and after being fed a particular diet.

Sometimes, the difference in the value of the measurement of interest for each matched pair is calculated, for example, the difference between before and after measurements, and these figures then form a single sample for an appropriate statistical analysis.

## 8.2 Confidence Intervals

### 8.2.1 Difference of Means for large samples

Suppose  $X_i$  and  $Y_j$  are independent normally distributed samples of size  $n_X$  and  $n_Y$  respectively:

$$X_i \sim N(\mu_X, \sigma_X^2)$$

$$Y_j \sim N(\mu_Y, \sigma_Y^2)$$

For large values of  $n_X$  and  $n_Y$ , by CLT, the random variable:

$$\frac{(\bar{X} - \bar{Y}) - (\mu_X - \mu_Y)}{\sqrt{\frac{\sigma_X^2}{n_X} + \frac{\sigma_Y^2}{n_Y}}} \sim N(0, 1)$$

An approximate 100p% confidence interval for  $\mu_x - \mu_y$  is given by

$$(\bar{x} - \bar{y}) \pm c \sqrt{\frac{s_X^2}{n_X} + \frac{s_Y^2}{n_Y}}$$

where  $s_X$  and  $s_Y$  are the **sample** standard deviations. For  $Z \sim N(0, 1)$ , we choose  $c$  to satisfy

$$P(-c \leq Z \leq c) = p$$

### 8.2.2 Difference of Means for small samples

Keep the basic setup from earlier, except now with smaller sample sizes. Suppose that  $X_i$  and  $Y_j$  are independent samples of sizes  $n_X$  and  $n_Y$  respectively. For **small** values of  $n_X$  and  $n_Y$ , the random variable  $\frac{(\bar{X} - \bar{Y}) - (\mu_X - \mu_Y)}{S_p \sqrt{\frac{1}{n_X} + \frac{1}{n_Y}}} \sim t_{n_X + n_Y - 2}$  Explanation for  $S_p^2$ :

$$S_p^2 = \frac{(n_X - 1)S_X^2 + (n_Y - 1)S_Y^2}{n_X + n_Y - 2}$$

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

$$S_Y^2 = \frac{\sum_{j=1}^{n_Y} (Y_j - \bar{Y})^2}{n_Y - 1}$$

We use the unbiased estimator  $S_p^2$  for  $\sigma^2$  rather than the maximum likelihood estimator  $\tilde{\sigma}^2$ . An approximate 100p% confidence interval for  $\mu_x - \mu_y$  is given by

$$(\bar{x} - \bar{y}) \pm cs_p \sqrt{\frac{1}{n_X} + \frac{1}{n_Y}}$$

where  $c$  is chosen to satisfy  $P(-c \leq t_{n_X + n_Y - 2} \leq c) = p$

### 8.2.3 Difference of means, not independent

Suppose we have paired observations  $(x_i, y_i)$ , where  $X_i \sim N(\mu_X, \sigma^2)$  and  $Y_i \sim N(\mu_Y, \sigma^2)$  are two random variables measured from the same unit of the population

Define new variables  $D_i = X_i - Y_i \sim N(\mu_D, \sigma_D^2)$  with observations

$$d_i = x_i - y_i$$

We have the pivotal quantity  $\frac{\bar{D} - \mu_D}{\frac{s_D}{\sqrt{n}}} \sim t_{n-1}$  where  $P(-c \leq t_{n-1} \leq c) = p$  and  $s_D^2 = \frac{\sum_{i=1}^n (d_i - \bar{d})^2}{n-1}$

The 100p% confidence interval for  $\mu_X - \mu_Y$  is

$$(\bar{x} - \bar{y}) \pm c \frac{s_D}{\sqrt{n}}$$

where  $P(-c \leq T \leq c) = p$  and  $T \sim t_{n-1}$

### 8.2.4 Exact Confidence Intervals

$X_i \sim N(\mu, \sigma^2)$

A 100p% confidence for  $\mu$  where  $\sigma$  is unknown) is given by

$$\bar{x} \pm c \frac{s}{\sqrt{n}}$$

where  $P(-c \leq T \leq c) = p$  and  $T \sim t_{n-1}$

A 100p% confidence interval for  $\sigma$  is given by

$$\left[ \sqrt{\frac{(n-1)s^2}{b}}, \sqrt{\frac{(n-1)s^2}{a}} \right]$$

where  $P(a \leq X \leq b) = p$  and  $X \sim \chi_{n-1}^2$

### 8.2.5 Approximate Confidence Intervals

An approximate 100p% confidence interval for  $\theta$  (Binomial success probability) is given by

$$\hat{\theta} \pm c \sqrt{\frac{\hat{\theta}(1-\hat{\theta})}{n}}$$

where  $P(-c \leq Z \leq c) = p$  and  $Z \sim N(0, 1)$

# Chapter 9

## Hypothesis Testing

### 9.1 Steps

- Define a null hypothesis and alternative hypothesis
- Define a test statistic which is used as evidence against the null hypothesis
- Calculate the  $p$ -value from the test statistic, with a given threshold( $\alpha$ )
- If the  $p$ -value is less than the threshold, then we reject the null hypothesis
- Always state conclusions in the language of the original problem

### 9.2 Difference between One sided and Two sided

For a null hypothesis,  $H_0 : \mu = \mu_0$ , we consider two types of alternative hypothesis

- One sided:  $H_a : \mu > \mu_0$  or  $H_a : \mu < \mu_0$
- Two sided:  $H_a : \mu \neq \mu_0$

Which alternative hypothesis is appropriate depends on the context of the problem

### 9.3 Formulas

For  $\mu \neq \mu_0$ :

$$|T| = \frac{|\bar{D}|}{\frac{s_D}{\sqrt{n}}} \sim |t_{n-1}|$$

$$|Z| \sim |N(0, 1)|$$

For  $\mu > \mu_0$ :

$$T = \frac{\bar{D}}{\frac{s_D}{\sqrt{n}}} \sim t_{n-1}$$

$$Z \sim N(0, 1)$$

For  $\mu < \mu_0$ :

$$-T = -\frac{\bar{D}}{\frac{s_D}{\sqrt{n}}} \sim -t_{n-1}$$

$$-Z \sim -N(0, 1)$$

## 9.4 Contingency tables

- Contingency tables contain the observed counts between two random variables
- Test  $H_0$  : The variables are independent
- If  $r_i$  and  $c_j$  are the row and column totals, then the expected counts are

$$e_{ij} = \frac{r_i c_j}{n}$$

- Test statistic: Under  $H_0$ , when  $n$  is large, the pivotal quantity

$$\Lambda = 2 \sum_{i=1}^a \sum_{j=1}^b Y_{ij} \ln \left( \frac{Y_{ij}}{E_{ij}} \right)$$

has a  $\chi^2_{(a-1)(b-1)}$  distribution