

## Contents

1	Terminology	1
2	Notation	1
3	Mean and Variance of RVs	1
3.1	Linear Combinations of RVs . . . . .	1
3.2	Summation Identities . . . . .	1
4	Normal RVs and $\chi^2$ RV	2
4.1	Normal RVs . . . . .	2
4.2	$\chi^2$ Distribution . . . . .	2
5	Hypothesis Testing and Confidence Intervals	2
5.1	Testing for difference in means . . . . .	2
5.2	Confidence Interval for difference in means .	2
5.3	Assumptions . . . . .	3
6	Experimental Design	3
6.1	Sampling . . . . .	3
6.2	Factors . . . . .	3
6.3	Crossed vs. Nested . . . . .	3
6.4	Blocking . . . . .	3
6.5	ANOVA Designs . . . . .	3

## 1 Terminology

**Subject** : A person, place or thing from which we measure data.

**Population** : The collection of all subjects of interest.

**Sample** : A subset of the population, from which we collect data.

**Response/Dependent Variable** : The variable which is of primary interest.

**Explanatory/Independent Variable** : The variable which we believe helps explain some of the variance in the response variable.

For this class, the response variable is **numerical**, and the explanatory variable/s are **categorical**. This is often phrased as “how does this numerical variable differ by group?”

**Random Variable** : A variable whose outcome is random. These are typically denoted by capital letters.

## 2 Notation

1 Let  $Y$  = the random variable denoting all possible values of the response variable.

1 Let  $y$  = an observed value of  $Y$  (in other words, measured observations).

2 Let  $Y_{ij}$  = the rv denoting all possible values of the  $j$ th observation in group  $i$ .

2 Let  $y_{ij}$  = the  $i$ th observed value of the  $j$ th group in  $Y$ . As an example, let  $i = 1$  denote sex M, and  $i = 2$  denote sex F.

- $Y_{13}$  = All possible values of height for the 3rd male. The outcome is random and unknown.

- $y_{13}$  = 72 inches. The specific observed value of the 3rd male once measured.

$\mu_i$  = The population mean for group  $i$  (a single value).

$\bar{Y}_i$  = All possible values of the sample mean for group  $i$ .

$\bar{y}_i$  = A specific, observed value of  $\bar{Y}_i$ . In other words, the mean of one given sample.

$\sigma_i$  = The population standard deviation for group  $i$ .

$S_i$  = All possible values of the sample standard deviation.  
 $s_i$  = A specific, observed value of  $S_i$ . In other words, the standard deviation of one given sample.

The book does not make this distinction  $\implies Y = y$  and  $\bar{Y} = \bar{y}$

**Parameter** : The unknown population value of some statistic. For example  $\mu_i, \sigma_i$ . These values are constant (non-random) if we could measure the population we would know the true value.

The goal of statistics 106 is to estimate parameters with sample values, and use the assumed distribution of those sample values to form **Hypothesis Tests (HTs)** and

**Confidence Intervals (CIs)**.

## 3 Mean and Variance of RVs

Let  $Y_i$  be drawn from a distribution with population mean  $\mu_Y$  and population standard deviation  $\sigma_Y$ .

Let the **mean** of  $Y_i = \mu_{Y_i} = E\{Y_i\} = \mu_Y$ .

Let the **standard deviation** of  $Y_i = \sigma_{Y_i} = \sigma\{Y_i\} = \sigma_Y$ .

### 3.1 Linear Combinations of RVs

Let  $Y^*$  be a linear combination of  $Y$  where  $a, b \in \mathbb{R}$ , then

Combination	Mean	Variance
$Y^* = a + Y$	$\mu_{Y^*} = a + \mu_Y$	$\sigma_{Y^*}^2 = \sigma_Y^2$
$Y^* = bY$	$\mu_{Y^*} = b\mu_Y$	$\sigma_{Y^*}^2 = b^2\sigma_Y^2$
$Y^* = a + bY$	$\mu_{Y^*} = a + b\mu_Y$	$\sigma_{Y^*}^2 = b^2\sigma_Y^2$

### 3.2 Summation Identities

Let  $Y_1, Y_2, \dots, Y_n$  be RVs

$$\begin{aligned}
 E\left\{\sum_{i=1}^n Y_i\right\} &= E\{Y_1 + Y_2 + \dots + Y_n\} \\
 &= E\{Y_1\} + E\{Y_2\} + \dots + E\{Y_n\} \\
 &= \sum_{i=1}^n E\{Y_i\}
 \end{aligned}$$

If  $Y_1, Y_2, \dots, Y_n$  are independent RVs,

$$\sigma^2\left\{\sum_{i=1}^n Y_i\right\} = \sum_{i=1}^n \sigma^2\{Y_i\}$$

Let  $\bar{Y} = \frac{1}{n} \sum_{i=1}^n Y_i$ ,  $Y_i$  is independent with mean  $\mu_Y$  and std.dev.  $\sigma_Y$ .

$$\begin{aligned} E\{\bar{Y}\} &= E\left\{\frac{1}{n} \sum_{i=1}^n Y_i\right\} = E\left\{\frac{1}{n}(Y_1 + Y_2 + \dots + Y_n)\right\} \\ &= \frac{1}{n} E\{Y_1 + Y_2 + \dots + Y_n\} \\ &= \frac{1}{n} (E\{Y_1\} + E\{Y_2\} + \dots + E\{Y_n\}) \\ &= \frac{1}{n} \sum_{i=1}^n E\{Y_i\} = \frac{1}{n} \sum_{i=1}^n \mu_Y \\ &= \frac{1}{n} (\mu_Y + \mu_Y + \dots + \mu_Y) = \frac{1}{n} (n * \mu_Y) \\ E\{\bar{Y}\} &= \mu_Y \\ \sigma^2\{\bar{Y}\} &= \sigma^2 \left\{ \frac{1}{n} \sum_{i=1}^n Y_i \right\} = \left( \frac{1}{n} \right)^2 \sigma^2 \left\{ \sum_{i=1}^n Y_i \right\} \\ &= \left( \frac{1}{n} \right)^2 \sum_{i=1}^n \sigma^2\{Y_i\} \end{aligned}$$

## 4 Normal RVs and $\chi^2$ RV

### 4.1 Normal RVs

A normal RV follows a bell curve created by a probability density function (pdf).

If  $Y$  is normally distributed with mean  $\mu_Y$  and std dev  $\sigma_Y$ , we say that  $Y \sim N(\mu_Y, \sigma_Y)$

$Y \sim N(\mu_Y, \sigma_Y) \implies Y^* = a + bY \sim N(a + b\mu_Y, b\sigma_Y)$

From this we can get two more results,

If  $Y_1, \dots, Y_n$  independent and  $Y_i \sim N(\mu_Y, \sigma_Y)$ , then

1.  $\bar{Y} \sim N(\mu_Y, \sigma_Y/\sqrt{n})$
2.  $\sum Y_i \sim N(n\mu_Y, \sqrt{n}\sigma_Y)$

The standard normal distribution is a specific linear combination of a general normal distribution, denoted  $Z$ .

Let  $Y \sim N(\mu_Y, \sigma_Y)$

$$\begin{aligned} Z &= \frac{Y - \mu_Y}{\sigma_Y} = \frac{Y}{\sigma_Y} - \frac{\mu_Y}{\sigma_Y} \\ E\{Z\} &= \frac{-\mu_Y}{\sigma_Y} + \mu_Y \left( \frac{1}{\sigma_Y} \right) = 0 \\ \sigma_Z^2 &= \left( \frac{1}{\sigma_Y} \right)^2 \sigma_Y^2 = 1 \end{aligned}$$

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

### 4.2 $\chi^2$ Distribution

The  $\chi^2$  distribution (chi-squared) is a sum of independent squared  $Z$  distributions.

Let  $Z_1, Z_2, \dots, Z_n$  be independent RVs where  $Z_i \sim N(0, 1)$

$X = Z_1^2 + \dots + Z_n^2 \sim \chi_r^2$  with degrees of freedom

$r =$  The number of summed and squared  $Z_i^2$

$$E\{\chi_r^2\} = r$$

## 5 Hypothesis Testing and Confidence Intervals

### 5.1 Testing for difference in means

Step 1: Declare Hypothesis

$$H_0 : \mu_1 = \mu_2 \text{ or } \mu_1 \leq \mu_2 \text{ or } \mu_1 \geq \mu_2$$

$$H_A : \mu_1 \neq \mu_2 \text{ or } \mu_1 > \mu_2 \text{ or } \mu_1 < \mu_2$$

Step 2: Calculate test-statistic

$$t_s = \frac{(\bar{y}_1 - \bar{y}_2) - \Delta_0}{\sqrt{\left( \frac{s_1^2}{n_1} + \frac{s_2^2}{n_2} \right)}}$$

If equal variances are assumed, the following test statistic formula can be used.

$$\begin{aligned} t_s &= \frac{(\bar{y}_1 - \bar{y}_2) - \Delta_0}{\sqrt{s_p^2 \left( \frac{1}{n_1} + \frac{1}{n_2} \right)}} \sim t, df = n_1 + n_2 - 2 \\ s_p^2 &= \frac{s_1^2(n_1 - 1) + s_2^2(n_2 - 1)}{n_1 + n_2 - 2} \end{aligned}$$

$t_s$  = The number of estimated standard deviations our sample difference in means is from the null.

Step 3: Calculate the p-value

If  $H_A \implies$  p-value

$$\begin{array}{ll} H_A \implies & p \\ \mu_1 \neq \mu_2 & 2P\{t > |t_s|\} \\ \mu_1 < \mu_2 & P\{t < t_s\} \\ \mu_1 > \mu_2 & P\{t > t_s\} \end{array}$$

p-value =  $P\{\text{our data or more extreme} | H_0 \text{ TRUE}\}$

p-value = probability of observing our sample data or more extreme, if in reality the null hypothesis were true.

Step 4: State decision rule and conclusion

If  $p\text{-value} < \alpha$ , reject  $H_0$

If  $p\text{-value} \geq \alpha$ , fail to reject  $H_0$

Recall that

$$\alpha = P\{\text{Type I Error}\} = P\{\text{reject } H_0 | H_0 \text{ true}\}$$

### 5.2 Confidence Interval for difference in means

The corresponding  $(1 - \alpha)100\%$  CI for  $(\mu_1 - \mu_2)$  is

$$(\bar{y}_1 - \bar{y}_2) \pm t_{1-\alpha/2; n_1+n_2-2} \sqrt{s_p^2 \left( \frac{1}{n_1} + \frac{1}{n_2} \right)}$$

$t_{1-\alpha/2; n_1+n_2-2}$  is the  $(1-\alpha)100th$  percentile of a t distribution with  $df = n_1 + n_2 - 2$

## 5.3 Assumptions

1. Random samples from both groups.
2. Groups are independent
3.  $\sigma_1 = \sigma_2$  if using  $s_p$  formula.
4.  $\bar{Y}_1 - \bar{Y}_2$  is distributed normally, either because
  - (a) Both populations are normal
  - (b)  $n_1$  and  $n_2 \geq 30$  (Central limit theorem)

# 6 Experimental Design

## 6.1 Sampling

In ANOVA studies, the sampling scheme is very important. Typically, the **categorical variable** is seen as a **treatment**, and the goal is to see if it had an effect on the numerical variable.

In an **experiment**, subjects are randomly assigned a **treatment**, and the results are assessed to find a causal relationship between variables.

In an **observational study**, subjects are randomly sampled and may fall into natural **treatment groups**, but are not assigned one. The data is assessed to find **correlations** between variables.

## 6.2 Factors

**Factors** are the variables that experimenters control during an experiment in order to determine their effect on the response variable. A factor can take on only a small number of values, which are known as factor levels. Examples of factors are brand of equipment, where the factor levels are brand A, B, and C.

A **treatment** is a combination of factors that has been applied to a subject. Ex: A study with two factors - control vs drug group, and patient blood type.

bt/drug	A	B	AB	O
C	C,A	C,B	C, AB	C, O
D	D,A	D,B	D, AB	D, O

C,A and D,A are two possible treatments.

## 6.3 Crossed vs. Nested

When we have two factors, the design can be either **crossed or nested**.

A **crossed design** is where every possible treatment (combinations of factor levels) is present in the study.

A **nested design** is where not all possible treatments are present. For example, if we have 8 schools and two teaching methods, but not all schools teach both types.

	A	A	B	B	C	C	C	C
1	1,A	1,A	1,B	1,B				
2					2,C	2,C	2,C	2,C

Here, we would say that schools are nested within class format.

## 6.4 Blocking

Consider an experiment that is trying to determine if a new supplement increases vitamin C absorption.

Let Y response variable = vitamin C absorption

Let Factor A = group with levels "control" and "new". There is a **total variance** in how subjects absorb vitamin C. If we can explain more of that variance, we are more likely to be able to tell if factor A had an effect.

**Blocking** is using another explanatory variable to further split the subjects. For example, perhaps gender affects how subjects absorb vitamin C. Then we could first block (separate) subjects by gender, then randomly assign them to factor A. This may reduce unexplained variance in Y.

## 6.5 ANOVA Designs

Most ANOVA models assume an underlying structure to the data,

$$Y = [\text{overall constant}] + [\text{same things}] + [\text{individual error}]$$

For example, we may say that the height of a tree has some **overall value** which could be affected by the **same things**, and then also **individual variance (error)**. Depending on the design of the study, we use different models.

**Completely Randomized Designs** are where treatments are assigned to subjects randomly.

For example, say we assign a sample of trees randomly to 4 different fertilizers (A,B,C,D).

Then our model would be

$$\text{height} = [\text{some constant}] + [\text{fertilizer effect}] + [\text{ind error}]$$