

# Stats E139 Fall 2015

## Study Notes

David Wihl

November 11, 2015

## Contents

<b>1</b>	<b>Logistics</b>	<b>4</b>
1.1	R Demonstration . . . . .	5
<b>2</b>	<b>Unit 1 - Data Collection</b>	<b>6</b>
2.1	Sampling . . . . .	7
2.2	Random Sampling . . . . .	7
<b>3</b>	<b>Unit 2 - Probability</b>	<b>7</b>
3.1	General Probability Review . . . . .	7
3.2	Random Variables . . . . .	8
3.3	Distributions . . . . .	9
3.4	Mean and Variance of $\bar{x}$ . . . . .	9
<b>4</b>	<b>Section 2 - Using R for Statistics and Homework Review</b>	<b>9</b>
<b>5</b>	<b>Unit 3 - Hypothesis Testing</b>	<b>9</b>
5.1	The Hypothesis Testing Framework . . . . .	10
5.1.1	Determine a Test Statistic . . . . .	11
5.1.2	Calculate the $p$ value . . . . .	11
5.1.3	Significance Level of a Test . . . . .	12
5.1.4	Determining the Conclusion . . . . .	12
5.2	Fisher's Randomization Test . . . . .	12
5.2.1	Test Statistic . . . . .	13
5.2.2	Calculating $p$ value . . . . .	14
5.2.3	Conclusion . . . . .	14
5.2.4	Permutation Test . . . . .	14

<b>6</b>	<b>Intro to R</b>	<b>15</b>
6.1	For loops . . . . .	15
<b>7</b>	<b>Unit 4 -T Based Inference</b>	<b>15</b>
7.1	The $\chi^2$ distribution . . . . .	15
7.2	The $t$ distribution . . . . .	17
7.3	One Sample $t$ -based inference . . . . .	18
7.4	The Confidence Interval . . . . .	19
7.5	Caveats of $t$ -based CI and Hypothesis Tests . . . . .	20
7.6	Two Sample $t$ -based Inference . . . . .	20
7.7	Pooled Test . . . . .	22
7.8	Paired $t$ -test . . . . .	22
<b>8</b>	<b>Section 3</b>	<b>23</b>
<b>9</b>	<b>Section 4</b>	<b>23</b>
<b>10</b>	<b>Unit 5 - Assumptions and Robustness of <math>t</math>-based Inference</b>	<b>24</b>
10.1	Assumptions . . . . .	24
10.2	Robustness . . . . .	25
10.3	Breakdown of Independence Assumption . . . . .	26
10.4	Breakdown of Normality Assumption . . . . .	26
10.4.1	Analytical Tests for Normality . . . . .	28
10.4.2	Uniform Distribution . . . . .	28
10.4.3	The Universality of the Uniform Distribution . . . . .	28
10.4.4	Robustness of $t$ -test to Departures from Normality . . . . .	29
10.5	Equal Variance Assumption . . . . .	29
10.5.1	Unpooled (Welch) Two-Sample $t$ -test . . . . .	31
10.5.2	Graphical Check of Equal Variances . . . . .	31
10.6	Formal $F$ -test of Equal Variance . . . . .	31
10.6.1	$F$ -distribution . . . . .	32
<b>11</b>	<b>Unit 6 - Transformations</b>	<b>33</b>
11.1	Non-Linear (Log) Transformations . . . . .	33
11.1.1	Interpretation - Randomized Experiments . . . . .	34
11.1.2	Interpretation - Observational Studies . . . . .	34
11.1.3	Log Transform for Paired $t$ Test . . . . .	35
11.1.4	Graphical Interpretation . . . . .	35
11.1.5	Interpretation of Exponentiated Scale . . . . .	35
11.1.6	Limitations of the Log Transformation . . . . .	35
11.1.7	Other Types of Transformations . . . . .	36
11.2	Choosing a Transformation . . . . .	37

11.3	Summary: Evaluating Validity of $t$ -tools . . . . .	37
11.4	Alternatives to the $t$ -tools . . . . .	37
11.5	Nonparametric Tests . . . . .	38
11.6	Rank-Sum Test . . . . .	38
11.6.1	Rank-Sum Hypotheses . . . . .	38
11.6.2	Exact Sampling of Rank-Sum Distribution . . . . .	39
11.6.3	Rank-Sum: Normal Distribution of Sampling . . . . .	39
11.6.4	Example of Rank-Sum using R . . . . .	40
11.6.5	Rank-Sum Summary . . . . .	42
11.7	Sign Test . . . . .	42
11.7.1	Sign Test Hypotheses . . . . .	43
11.7.2	Exact Calculations and Normal Distributions . . . . .	43
11.8	Wilcoxon Signed Rank Test . . . . .	44
11.8.1	Hypotheses . . . . .	44
11.8.2	Exact Calculations and Normal Approximations . . . . .	44
11.9	Summary of Paired Test Choices . . . . .	45
<b>12</b>	<b>Section 6</b>	<b>45</b>
<b>13</b>	<b>Unit 7 - Power, Sample Size and Error</b>	<b>46</b>
13.1	Type I and II Error and Power . . . . .	46
13.1.1	Statistical Power . . . . .	46
13.2	Sample Size and Power Calculations . . . . .	47
13.3	Find Sample Size for a Desired Power . . . . .	49
13.4	Power and Sample Size Caveats . . . . .	49
13.5	Using R . . . . .	50
13.6	Tradeoff: Assumptions vs. Power . . . . .	50
13.6.1	Pooled vs. Unpooled $t$ -Tests . . . . .	50
13.6.2	Randomization vs. $t$ -Test . . . . .	51
13.6.3	Permutation Test . . . . .	51
13.6.4	When Assumptions are Violated . . . . .	51
13.7	Power Curves . . . . .	52
13.7.1	Multiple Comparisons and Type I Errors . . . . .	52
13.7.2	Bonferroni Correction . . . . .	54
<b>14</b>	<b>Section 7</b>	<b>54</b>
14.1	Rank-Sum . . . . .	54
14.2	Statistical Power . . . . .	55
14.2.1	Power Calculation Example . . . . .	55
<b>15</b>	<b>Unit 8 - Analysis of Variance (ANOVA)</b>	<b>57</b>

15.1	Analysis of Variance (ANOVA)	57
15.1.1	General format and ANOVA's $F$ -test	58
15.1.2	Main Concept of ANOVA	58
15.1.3	One way ANOVA	58
15.1.4	Concept Behind the Test	60
15.1.5	Analysis of Variance Table	61
15.1.6	Critical Value $F^*$	61
15.1.7	Using ANOVA from R	62
15.1.8	Assumptions for ANOVA $F$ -test	62
15.1.9	Contrast Testing	63
15.1.10	Contrast Test Hypothesis	66
15.1.11	Multiple Comparisons	67
15.1.12	Tukey Honest Significant Difference (HSD)	68
15.1.13	Two-way and multi-way ANOVA	68
15.2	Kruskal-Wallis Test	69
<b>16</b>	<b>Section 8 - ANOVA</b>	<b>70</b>
<b>17</b>	<b>Section 9 - Contrast</b>	<b>70</b>
<b>18</b>	<b>Unit 9 - Simple Linear Regression</b>	<b>70</b>

## 1 Logistics

Prerequisites: AP Stats, Stats 101, 101, 102, 104 or 110.

Uses Calculus and includes algebraic derivations. Answers questions not covered by intro Stats class.

Office hours: SC-614, 7:30-8:30pm, also by appointment

Office (not used) 617-495-8711 or preferably via the stats dept 617-495-5496. Best of all use email: [krader@fas.harvard.edu](mailto:krader@fas.harvard.edu)

TA: David Haswell

Combination of undergrad / grad level (master's level) material.

Sections: . Wed 5:30-6:30 Sever 112, followed by Office hours 6:30-7:30, Tory and Alina .  
Thu 7:40-8:40, 1 Story St 307, by David. Available on video, followed by Office hours

Lectures will **not** be broadcast live. There will be delayed video posted within 24 hours.

All material in lecture is tested even if it isn't in the lecture notes. The textbook is useful but not required.

R is used extensively. Course assumes no prior R knowledge. First homework requires some basic R. Start with "Comprehensive Intro to R" on slide 13.

Lectures and classroom is interactive - should bring questions.

Exams: final exam requires using R which is why it is a take-home exam.

Group project is 2-4 people, likely analyzing a data set. Final project grade: 2/3 group, 1/3 take home exam.

Homework is due by class time on Monday. Recommended to use L<sup>A</sup>T<sub>E</sub>X instead of Word. It is ok to scan or photograph hand-written work as long as it is submitted in PDF.

Collaboration is encouraged, but ensure that collaborators are cited appropriately.

There is no curve for exams and homeworks. At the end of the year, all students' grades are summed and a distribution with A cutoff.

The mid-term on Nov 16 is through HW8, Intro to Linear Regression. Two pages, double-sided cheat sheet allowed. **Calculator Required**. Two hours long.

See Math Review sheet on Course Website.

Matrix Algebra will be emphasized later in the class to understand the math behind the techniques. Matrix Algebra will not be on tests.

For graduate students, there will be extra problems. Undergraduate students may do these additional problems for extra credit.

The Harvard College Stats139 class has three hours of lecture whereas this class has only two. So this class has the same breadth, but a little less mathematical rigor and depth.

Statisticians are trained, whereas mathematicians are usually born.

A good statistician is able to communicate their analysis. This is often the most important step of the process.

Units 1, 2, 3, 4 are expected to refresher and will be done quickly.

Quantitative variables can be discrete or continuous.

Categories / Qualitative are not going to be covered much in depth in this class.

## 1.1 R Demonstration

(slides 30 and 31)

Use weather data that is the same source, which is trustworthy. Use Max temperature per day because it is readily available unlike the average. Define what constitutes a “warmer” summer, by choosing an appropriate metric, like comparing means. Include a standard error. If the data is bell shaped, we can use a *t-test*. If the data is not bell shaped, we will explore other approaches.

(demo showing R Studio)

There is a lot more code in the Course site than the lecture notes because of the additional data processing and cleaning.

## 2 Unit 1 - Data Collection

(not many notes in this section as the lecture was primary a dramatic reading of the lecture notes).

Clinical experiments often has as few as 100 participants because the experiments are so expensive.

The balance of confounding factors by having random samples should isolate causation.

Key factors to an experiment: control group, randomized, and replicated. Usually 30 replications is sufficient.

People in a clinical trial may not be representative of the general population and followed regularly by a doctor. They may also be above average health which is why they choose to participate in trials.

One way to avoid confounding is to have *cross-over* groups, so they exchange methods during the course of the class.

Lecture stops at slide 19.

[a: rent or buy textbook

a: learn more R

a: evaluate R Studio]

Lecture 2 by Michael Parzen

Lecture restarts Unit 1, slide 19

Unless you do experimentation, you can't establish causality. There are really only two sources of data: experimental and observational.

Types of Observational studies: (all time related) prospective, retrospective, longitudinal.

## 2.1 Sampling

The sample is a subset of the population of interest. Actual vs. Conceptual: you may not be able to get a representative sample, eg. Alzheimer's patients who are not institutionalized.

Parameters:  $\mu$  population mean,  $\pi$  (or  $p$ ) is true but unknown population proportion.

Statistics is a function of your data - a numerical summary of your sample.

Estimator  $\bar{x}$ , true  $\mu$ .

Estimator  $\bar{X}$ , estimate is  $\bar{x}$ . Ideally a Simple Random Sample (SRS).

Census: sample everyone in the target population.

Sampling Frame: collection of units that are potential members of the sample. Bad sampling (**biased**): voluntary response, convenient sampling, question framing, confusing questions

## 2.2 Random Sampling

Systemic Random Sampling: every  $k$  member. Easy to admin, but population must be well mixed

Variable probability sampling: allow units to have unequal probabilities of being sampled

Interval validity: double check that allocation of groups is appropriate.

# 3 Unit 2 - Probability

## 3.1 General Probability Review

Probability is a measure of uncertainty and cannot be negative. (This section is mostly a recap so it is in summary form.)

Axiom 1:  $P(A) \geq 0$

Axiom 2:  $P(S) = 1$  guaranteed to happen

Axiom 3: If  $A_1, A_2, \dots$  are disjoint events, then

$$P\left(\bigcup_{j=1}^{\infty} A_j\right) = \sum_{j=1}^{\infty} P(A_j)$$

Sampling with replacement:  $n$  objects and making  $k$  choices, results in  $n^k$  possible outcomes.

Sampling without replacement:  $n$  objects and making  $k$  choices:  $n(n-1)(n-2)\dots(n-k+1)$  possible outcomes.

A group of  $k$  people from a population of  $n$  people:

$$\binom{n}{k} = \frac{n!}{k!(n-k)!}$$

Conditional probability:

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

$$P(A \cap B) = P(A|B)P(B) = P(B|A)P(A)$$

Bayes Rule:

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

Events are *independent* if  $P(A|B) = P(A)$ .

Parzen's 2x2 table:

	$B$	$\bar{B}$	
$A$	$P(A \cap B)$	$P(A \cap \bar{B})$	$P(A)$
$\bar{A}$	$P(\bar{A} \cap B)$	$P(\bar{A} \cap \bar{B})$	$P(\bar{A})$
	$P(B)$	$P(\bar{B})$	1.0

Summing in one dimension gives the marginals.

$$\begin{aligned} P(B) &= P(A \cap B) + P(\bar{A} \cap B) \\ &= P(B|A)P(A) + P(B|\bar{A})P(\bar{A}) \end{aligned}$$

Two events  $A$  and  $B$  are *conditionally independent given  $E$*  if and only if:  $P(A \cap B|E) = P(A|E)P(B|E)$ .

## 3.2 Random Variables

Two types: Continuous and discrete random variables. Probability of a specific random variable = 0. Cumulative Distribution Function of a random variable given by  $F_x(x) = P(X \leq x)$ .



(See slides 15-18) Expectation, Properties of Expectation, Variance Definition, Properties of Variance.

### 3.3 Distributions

Normal / Gaussian Distribution (see slides)

In R, use the `pnorm(y, mean, sd)` command for CDF of the Normal.

Sums of Normal are Normal.

The Binomial is the most common discrete distribution.

See slide 26 for an approximation of the binomial distribution (not necessary in R).

### 3.4 Mean and Variance of $\bar{x}$

Because  $\bar{x}$  is random,  $E[\bar{x}] = \mu$ .

Law of large numbers:  $n \rightarrow \infty, \text{var}(\bar{x}) = 0, E(\bar{x}) \rightarrow \mu$ .

Central Limit Theorem: with a large sample size, the distribution of  $\bar{x}$  is normally distributed. Useful only if you are interested in  $\bar{x}$ , not the population.

Example 4: a)  $P(X > 70000)$ . Can't be answered because we don't know the underlying distribution. b)  $P(\bar{x} > 70000)$  can be answered.

## 4 Section 2 - Using R for Statistics and Homework Review

stopped section video at 18:00

Homework review notes: Internal validity - causal relational? Only valid when there is a randomized controlled experiment. External validity - will it generalize? Was the random sampling sufficient to generalize the conclusion?

## 5 Unit 3 - Hypothesis Testing

Section 1.3-1.6 in the text

If we see a relationship, we assume all confounding factors have been balanced so the malaria study, even with a very small sample, does have internal validity.

$$Y_{i,v} = \begin{cases} 1 & \text{if individual gets malaria when vaccinated} \\ 0 & \text{otherwise} \end{cases}$$

$$Y_{i,c} = \begin{cases} 1 & \text{if individual gets malaria when **not** vaccinated} \\ 0 & \text{otherwise} \end{cases}$$

Treat as all independent, so Binomial distribution. Those variables should potentially follow the following distributions:

$$Y_{i,v} \sim \text{Bern}(P_v)$$

$$Y_{i,c} \sim \text{Bern}(P_c)$$

So what are the hypotheses?

$$H_0 : P_v = P_c$$

$$H_0 : P_v \neq P_c (P_v < P_c)$$

Presumably, malaria infections would not be increased by administering the vaccine.

Could use t-test only if distribution was normal or we had a large sample, neither of which applies.

Will use Fisher's Randomization Test.

Will need to determine proportion of infected to non-infected, ie  $\hat{P}_c - \hat{P}_v$  or  $(\bar{Y}_c - \bar{Y}_v)$ , but it will be easier to measure  $\bar{Y}$ . If  $H_0$  is true, the difference will be close to zero. If  $H_a$  is true, the difference will be large. But what will be sufficient difference to have statistical weight?

So divide by standard deviation to determine "how far is far" (like  $Z$  standardize). Since the dataset is so small, we will simulate more data.

## 5.1 The Hypothesis Testing Framework

Four steps in testing Hypotheses:

- Formulate hypotheses  $H_0$  and  $H_A$

- Calculate test statistic
- Calculate p-value based on a reference distribution of the test statistic (assume  $H_0$  is true)
- Determine conclusion of the test by

A **scientific hypothesis** makes a testable statement about the observable universe.

A **statistical hypothesis** is more restricted. It concerns the behavior of a measurable (or observable) random variable. It is often a statement or claim about a parameter of a population or distribution.

These two types of statistical hypotheses for any scientific:

**Null hypothesis** ( $H_0$ ) assumption specifying a possible truth, typically the absence of an effect. Presumes no change, old beliefs. Any discrepancy between the observed data and the hypothesis is due only to chance variation (noise).

**Alternative hypothesis** ( $H_a$ ) assumption describing an alternative truth, typically some effect or some difference. A statement of possible new beliefs. An observed discrepancy between the observed data and the null hypothesis is **not** due to chance variation.

### 5.1.1 Determine a Test Statistic

**Statistic** is a function of the data that summarizes it.  $\hat{\mu} = \hat{\mu}(\mathbf{y})$

**Test statistic** a specific statistic used to weight evidence supporting or contradicting the null hypothesis.

**Reference Distribution** Probability distribution of the test statistic, assuming the null hypothesis is true  $f(\hat{\mu}(Y)|H_0)$ . This is often called the *sampling distribution*.

Fisher's Exact Distribution is a way of doing this.

### 5.1.2 Calculate the $p$ value

**$p$  value** the probability of observing our test statistic or a more extreme one, assuming the null hypothesis to be true. Measure of strength of evidence of the hypotheses.

Calculated by *comparing* the observed test statistic to the reference distribution.

**$p$  value is NOT the probability that the null hypothesis is true.** It is probability of seeing our data on the null hypothesis. Use Bayes Rule to flip this around.

### 5.1.3 Significance Level of a Test

**$p$  value** probability that the test statistic would be at least as extreme as *observed* under the null hypothesis

**significance level ( $\alpha$ )** is the criterion compared against the  $p$  value. The null hypothesis is **rejected** if  $p$  value is lower than  $\alpha$ . (Usually 0.05)

Generally,  $\alpha$  reflects the probability of rejecting the null hypothesis when it is true (a Type I error)

Two sided are considered more conservative because it is usually considered harder to have a conclusive test.

### 5.1.4 Determining the Conclusion

We come to a conclusion about our hypotheses by comparing the  $p$  value to the Type I error rate.

If the  $p$  value  $\leq \alpha$  we “reject the null hypothesis” and say the result is “statistically significant at level  $\alpha$ ”

If the  $p$  value  $> \alpha$ , we are “unable to reject the null hypothesis”. This is different than concluding that the null hypothesis is true.

Based on the study design, we can then generalize internally and/or externally, which the **scope** of the inference procedure.

## 5.2 Fisher’s Randomization Test

Randomization Test is for experiments - for casual inference.

Permutation Tests are for observational studies - for generalization.

In randomized experiments, uncertainty comes from randomness of an assignment. A Fisher Randomization Test is a *distribution free* test for treatment effect in **randomized experiments**. No assumption as to the underlying distribution has to be made. Generally, is  $\mu_1 = \mu_2$ ?

The only assumption is Additive Treatment Effect ( $\delta$ ):

$$Y_{c,i} = Y_{v,i} + \delta$$

In practice, only the control or the experimental can be measured as one study unit cannot be both vaccinated and not vaccinated for example.

$H_0 : \delta = 0$ , ie zero treatment effect for all units. Each unit's outcome is the same regardless of the treatment assigned. So the distribution would be identical in both groups.

$H_a : \delta \neq 0$  non-zero treatment effect for ALL units. (The textbook uses  $Y^* = Y + \delta$ ).

Assumptions:

- Random assignment to groups
- Under  $H_0$ , *independence* of study units. More precisely, there is *interchangability* of study units. The  $Y$  will be the same for a given study unit irrespective of treatment. This is taken advantage of when building the sampling distribution (through simulation).

### 5.2.1 Test Statistic

Difference of outcomes of the two groups:

$$\hat{\delta} = \bar{Y}_c - \bar{Y}_v$$

other test statistics could have been used, such as the difference between the medians.

**Randomization Distribution** is the reference distribution of a test statistic in a randomization test, where variation is due to random assignment of the treatment.

Procedure:  $Y$  is values of malaria.  $X$  is if someone is vaccinated. Simulate a new  $X^*$  with the same values as  $X$  but *in a different order*, leaving the same  $Y$ . Repeat multiple times to simulate and build the histogram. This will show if the distribution is extreme.

The number of simulations could be randomly very large. The maximum number of possible simulations is the binomial distribution:

$$\binom{t}{s}$$

where  $t$  is the total number of study units and  $s$  is the number where the effect was shown.

### 5.2.2 Calculating $p$ value

One sided alternative:

$$H_0 : \delta = 0$$
$$H_a : \delta > 0 \text{ ( or } \delta < 0 \text{ )}$$

Be sure to justify the sides.

The  $p$  value is proportion of values above or below the observed test statistic.

$$\hat{\delta} = 0.5$$
$$p = P(\bar{Y}_c - \bar{Y}_v \geq \bar{y}_c - \bar{y}_v)$$
$$= P(\bar{Y}_c - \bar{Y}_v \geq 0.5) = 0.029$$

so there is sufficient evidence to reject the null.

### 5.2.3 Conclusion

Under the significance level  $\alpha = 0.05$  there is evidence that the vaccine is effective. More generally, if effect is NOT homogeneous, then there is evidence that the vaccine was effective for at least one volunteer.

So  $\hat{\delta}$  means that the vaccine could be effective for 50% of the population.

Scope of inference: Internal Validity maybe not if the same hospital. External validity: possibly not, as these were volunteers.

### 5.2.4 Permutation Test

The Permutation Test is precisely the same although for observation studies.

This example of malaria was used only  $\{1, 0\}$  as  $Y$ . In practice,  $Y$  will usually be a non-discrete number. However,  $X$  will always to be reduced to  $\{1, 0\}$  to determine if study unit received the treatment or not.

See slides.

## 6 Intro to R

### 6.1 For loops

Use for repeated sampling, operating on a vector or matrix, Markov chains, simulation studies.

R has a `for` and `while` loops.

Example:

```
n.iter = 100
# create a vector of 100 NAs (like [float('nan')] * 100 in Python)
variable = rep(NA, n.iter)

for (i in 1:n.iter){
  samp = runif(10)
}
```

## 7 Unit 4 -T Based Inference

Textbook Chapter 2

### 7.1 The $\chi^2$ distribution

Generally right skewed, values from 0 to  $\infty$ .

Let  $Y = Z_1^2 + \dots + Z_k^2$  where  $Z_1, Z_2, \dots, Z_k$  are i.i.d  $\sim N(0, 1)$  (standard normal distribution), then  $Y$  follows a  $\chi^2$  distribution with  $k$  degrees of freedom ( $\sim \chi_{df=k}^2$ ). In English, a sum of standard normal variables squared, written as  $Y \sim \chi_k^2$ . This is related to the sample variance (since this is the square of standard normal distributed elements):

$$S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{x})^2 \sim \frac{\sigma^2}{n-1} \chi_{df=n-1}^2$$

Need the  $\sigma$  in order to standardize.

The PDF is:

$$f(y) = \frac{1}{\Gamma(k/2)} (y/2)^k (1/2) e^{-y/2}$$

It supports  $y > 0$ .

Skew / shape: heavily right skewed.

The mean:

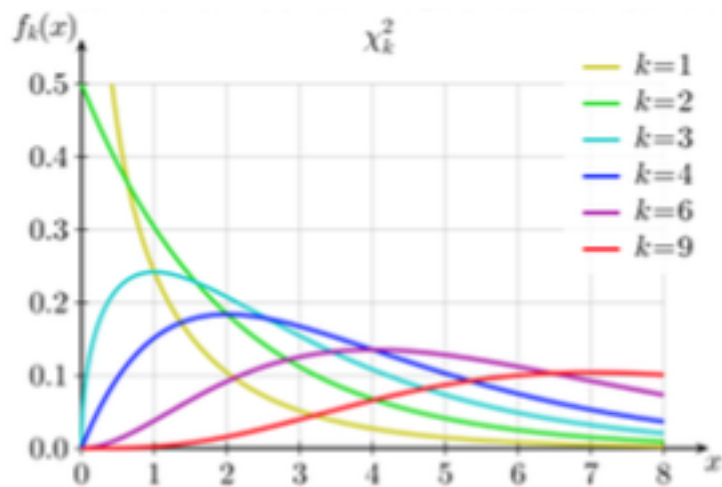
$$Y \sim \chi_{df=k}^2$$

$$\begin{aligned} E(Y) &= \int_0^{\infty} y \cdot f(y) dy \\ &= E(Z_1^2 + \dots + Z_k^2) \\ &= E(Z_1^2) + \dots + E(Z_k^2) \\ &= k \cdot E(Z_i^2) \\ &= k \end{aligned}$$

because each  $Z$  is normally distributed, the following applies:

$$\begin{aligned} Var(X) &= E[(X - E(X))^2] \\ &= E(X^2) - [E(X)]^2 \\ &= E(X^2) - 0^2 = 1 \\ E(X)^2 &= 1 \end{aligned}$$

The expected value of a  $\chi^2$  distribution is the degrees of freedom.





As  $k \rightarrow \infty$ , the  $\chi^2$  distribution looks more like a normal distribution.

We care about  $\chi^2$  distribution because the distribution of the sample variance as a random variable is based on a  $\chi^2$  random variable when underlying observations are Normal.

For i.i.d,  $X_1, \dots, X_n \sim N(\mu, \sigma^2)$ , the sample variance is the r.v.

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

then

$$\frac{(n-1)S^2}{\sigma^2} \sim \chi_{n-1}^2$$

We use  $n-1$  instead of  $n$  because they are not completely independent. They are slightly negatively correlated, canceling out one degree of freedom.

To prove it, we need Gram-Schmidt orthonormalisation. (See Math23a!!! Also in Stat 111).

Use tables to lookup  $\chi^2$  values.

## 7.2 The $t$ distribution

The  $t$  distribution has a slightly wider spread in the tails than Standard Normal. Also known as student-t distribution. This happens when you don't know the real  $\sigma$  and are forced to use  $S$ , which is less reliable.

$$\bar{X} \approx N\left(\mu, \frac{\sigma^2}{n}\right)$$

allows us to Z-score

$$Z = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}}$$

but we don't know  $\sigma$  so in practice we use

$$T = \frac{\bar{X} - \mu}{S/\sqrt{n}}$$

(which has two random variables ( $\bar{X}$  and  $S$ )). Now, take that previous definition of  $T$  and add  $\sigma/\sqrt{n}$  to the top and bottom:

$$\begin{aligned}
T &= \frac{\bar{X} - \mu}{S/\sqrt{n}} \\
&= \frac{(\bar{X} - \mu)/(\sigma/\sqrt{n})}{(S/\sqrt{n})/(\sigma/\sqrt{n})} \\
&= \frac{Z}{\sqrt{S^2/\sigma^2}} \\
&= \frac{Z}{\sqrt{\frac{(n-1)S^2}{(n-1)\sigma^2}}} \\
&= \frac{Z}{\sqrt{\chi^2/df}}, \text{ since } \chi^2 = \frac{(n-1)S^2}{\sigma^2}
\end{aligned}$$

This generates a random variable with  $t$  distribution, which is standard normal, with  $k$  degrees of freedom. Since we are dividing by  $S$  instead of  $\sigma$ , there is greater uncertainty, which is why the distribution is slightly wider.

Recall  $Z \sim N(0, 1)$  and  $\chi^2 \sim \chi^2_{df}$ . It is a  $t$ -test because it is based on a  $t$  distribution. The  $t$  distribution critical values get closer and closer to the normal distribution ( $z$ ) critical values as degrees of freedom increase.

This is important because there are a lot of times we don't know the true variance  $\sigma^2$  but still want to know if the sample mean  $\bar{X}$  matches the population mean.

### 7.3 One Sample $t$ -based inference

One sample  $t$ -based inference follows the same hypothesis testing framework. We won't know the population mean or standard deviation (but assume it is approximately normal), but want to determine from the sample whether  $\mu$  is particular value or not.

$$H_0 : \mu = 70 (\mu_0)$$

$$H_a : \mu \neq 70$$

$$t = \frac{\bar{X} - \mu_0}{s/\sqrt{n}}$$

Obviously if you have the real  $\sigma$  or  $\mu$ , use that! But it will be a  $Z$  score, not a  $t$  test.

$$z = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}}$$

Getting back to the  $t$ -test, with an example:

$$t = \frac{24923 - 24000}{22517/\sqrt{100}} = 0.410$$

$$\begin{aligned} p\text{-value} &= P(|t_{df=99}| > 0.410) \\ &= 2P(t_{df=99} > 0.41) = 2 * 0.3413 = 0.6826 \end{aligned}$$

0.3413 comes from R: `1 - pt(0.410, df=99)` (probability from a  $t$  distribution)

Since our  $p$ -value (0.68) is  $> \alpha = 0.05$ , we fail to reject the null. There is insufficient evidence to suggest financial aid has changed.

**Note:**  $t$  values are measures of *distance*.  $p$  values are measure of *probability*.

Even though the population may not be normally distributed, if there are sufficient samples, the samples will be normally distributed so the  $t$  distribution can be trusted.

Could also use R's `t.test()` to do this in one step, and also provides a confidence interval, which is an estimate for the true population mean  $\mu$ .

## 7.4 The Confidence Interval

We want to make an *inference* about the population (e.g.  $\mu$  or  $p$ ) using information from the observed sample data. Instead of just a point estimate, it is more useful to have a *margin of error*.

CI is Estimate  $\pm$  margin of error or Estimate  $\pm (z \text{ value}) \times (\text{SD of estimate})$ .

We figure it out, using the formula from above:

$$t = \frac{\bar{X} - \mu_0}{s/\sqrt{n}}$$

and rearrange it to solve for  $\mu_0$ :

$$\mu_0 = \bar{X} \pm t \cdot s/\sqrt{n}$$

(we use  $\pm$  since we don't care if it is above or below  $\bar{X}$ ).

At this point  $t$  is no longer a statistic - it is a critical value coming from the  $t$  distribution.

To find the 95% Confidence Interval for the **true mean**  $\mu$ :

$$\bar{x} \pm t^*(s/\sqrt{n})$$

$t^*$  has to be looked up in R. Use `qt()` (quantile coming from a  $t$  distribution) instead of `pt()`. `qnorm()` is the same as `qt()` with infinite degrees of freedom.

If the confidence interval and  $\alpha$  add up to one (e.g. 95% confidence interval  $+\alpha(0.05) = 1.0$ ), then there is a direct relationship between CI and two-sided tests of hypotheses. When we reject the null, the calculated mean should be in the confidence interval. This may not hold for one-sided tests. (No one uses one-sided confidence intervals which is a value between  $x$  and infinity).

## 7.5 Caveats of $t$ -based CI and Hypothesis Tests

Assumptions:

1. Observations are **independent**. This is true if the sample is randomly selected, but false if there is bias in the sampling.
2. **Normal Distributions** of the observations. This happens when the sample is large enough (via the Central Limit Theorem) or it is known that the population is normally distributed.

If these assumptions are not correct, none of the  $t$ -tests will work correctly.

## 7.6 Two Sample $t$ -based Inference

This is the more common way to compare two different groups. Two different formulas for unpooled and pooled. Could also be paired.

Birthweight sample and smoking study. There was a lot of data, so there is likely correlations. The question is whether causality has been established. In this case, it was very specific, namely baby boys, survived at least 28 days and were single births. Good for confounding factors (internal consistency) but not good for generalization (external consistency).

In R:

```
cbind # combines columns
by # standard deviation, smoking
```

Anyone who did not answer the smoking question will be ignored, which is ok because there are only 10 out of 1200, so it isn't significant. If it was larger, we would have to worry about non-response bias. If we cared, we could distribute them in the two groups. Assign them to smoking status's that makes them most and least and significant. If that does not materially change the results, then non-response bias has no effect in this study.

Let  $X_1$  birthweight of a baby born from a smoking mother.

Let  $X_2$  birthweight of a baby born from a non-smoking mother.

$$X_1 \sim N(\mu_1, \sigma_1^2) X_2 \sim N(\mu_2, \sigma_2^2)$$

$$H_0 : \mu_1 = \mu_2$$

$$H_a : \mu_1 \neq \mu_2$$

alternatively, to generate the statistic we want to use:

$$H_0 : \mu_1 - \mu_2 = 0$$

$$H_a : \mu_1 - \mu_2 \neq 0$$

$$\left( \frac{\bar{X}_1 - \bar{X}_2}{\sqrt{\sigma_1^2/n_1 + \sigma_2^2/n_2}} \right) = 0$$

however we don't have the true population variances so we have to use  $S$  **WARNING: there be dragons here. Some of these equations are incorrect. TODO: Fix**

$$\left( \frac{\bar{X}_1 - \bar{X}_2}{\sqrt{S_1^2/n_1 + S_2^2/n_2}} \right) = 0$$

because of properties of Expectations:  $Var(X - Y) = Var(X) + Var(Y)$ . We don't have to worry about covariances because these are supposed to be independent.

This will be an approximate  $t$  distribution because they are approximately normally distributed. The degrees of freedom will be based on both  $S_1$  and  $S_2$ . To be conservative, we'll use the minimum of the two sample sizes  $df = \min(n_1, n_2) - 1$ .

There isn't a causation here - there is only an association because it is an observational study. It also can't be generalized that well.

Confounding factors: mothers who choose to smoke might be unhealthy in other ways.

Confidence intervals work the same (see slide 33).

Unpooled test because we used two different variances. Pooled means you use a single variance.

R does this for you, however it has positive values. Make sure the order of variables match. When data is grouped, use: `t.test(Y ~ X, var.equal=T)`. If data is not grouped, use `t.test(Y1, Y2, var.equal=T)`.

## 7.7 Pooled Test

Assume a single variance:

$$t = \frac{(\bar{x}_1 - \bar{x}_2) - (\mu_1 - \mu_2|H_0)}{s_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$$

where:

$$s_p^2 = \frac{(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2}{n_1 + n_2 - 2}$$

which is effectively combining the two sum of squares and then dividing by the combined variances. In this case,  $df = n_1 + n_2 - 2$

Combining variance gives more statistical power but could make a significant mistake. Possible if the two variances are very close. In practice, everyone does the unpooled test. Technicality: the pooled test is truly normally distributed.

If the ratio of variances is  $> 2$ , use unpooled (always put larger variance in the numerator). If the ratio is  $< 2$  you might be able to use the pooled test.

## 7.8 Paired $t$ -test

There are pairs of observations for each subject in the study, e.g. take a reading on one day and another a week later. The two measurements should be correlated, in theory. By looking at the difference, we control for person-to-person variability. The column of differences is all we need in order to perform a one-sample  $t$  test to see if the differences are  $= 0$ .

In R, `t.test(x1,x2, paired=T)`. If the sample sizes in the two groups, you might be able to use the Paired  $t$ -test. If the sample sizes are different, you cannot use the Paired  $t$  test.

## 8 Section 3

Steps of Hypothesis testing:

1. State the null and alternative hypothesis
2. Choose a relevant test statistic
3. Calculate the observed test statistic, e.g.  $t = |\bar{x}_0 - \bar{x}_1|$
4. Find the probability of observing the test statistic under the null (i.e. the p-value), e.g.  $p(t \geq 0.5 | H_0) = 0.03$
5. Conclusion: either fail to reject or reject, don't "accept." Mention scope of inference and internal validity. How general are the conclusions?

## 9 Section 4

One sample  $t$ -test, used when we want to test the mean of a single population.  $T$  stat has a  $df=n-1$

$$t = \frac{\bar{X} - \mu_o}{s/\sqrt{n}}$$
$$t \sim T_{df=n-1}$$

test the mean of a population for a single parameter. "We think the mean of the population is this? Can we accept / reject?" The Student T is slightly more spread out than the standard distribution.

Paired t-test:

Same mechanics as 1-sample test (where single sample is the difference in the pairs). Used when there is some natural pairing between subjects in each sample (e.g. twins).  $df = n - 1$  (where each sample has  $n$  observations).

$$t = \frac{\bar{x}_{Diff} - \mu_{Diff}}{\frac{s}{\sqrt{n_{Diff}}}}$$

Use paired over pooled test when possible if you expect a correlation between the two samples.

Unpooled test:

Unpooled works irrespective how close the two variances. Unless you are really sure the variances are the same, don't use the pooled test.

$$df = \min(n - 1, m - 1)$$

R would have a different  $df$ . Why is it a  $t$  distribution ?

$$\frac{Z}{\sqrt{\chi^2/n}} \sim T_n$$

Pooled Test:

only when you are really sure  $\sigma_x^2 = \sigma_y^2$ ,  $df = n + m - 2$ . Larger  $df$  means more statistical power. But there is an additional assumption to worry about. Note different  $S_p^2$  for pooled test and test statistic.

Assumptions:

in life, truly nothing is truly normally distributed. We are still making this assumption of normality. We are also assuming independence of observations. The Central Limit Theorem makes this more palatable. This is more about how the data was collected than any specific analysis. A good, randomized sampling method, means is more important than the math.

## 10 Unit 5 - Assumptions and Robustness of $t$ -based Inference

Chapter 3 in the text.

### 10.1 Assumptions

We've seen four  $t$ -based assumptions so far:

1. One sample  $t = \frac{\bar{X} - \mu}{s/\sqrt{n}}$

If the assumptions are true, the sample will be normally distributed.

2. Two Sample (independent groups)

- Unpooled  $t = \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$ .



If the assumptions are true, the samples will be *approximately* normally distributed. Also assumes that the two groups are independent.

- Pooled  $t = \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{S_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$

If the assumptions are true, the sample will be normally distributed.

3. Paired  $t = \frac{\bar{D} - 0}{S_d / \sqrt{n}}$

If the assumptions are true, the sample will be normally distributed.

Recap of  $t$ -based assumptions:

1. Independence of Observations  $X_i \stackrel{i.i.d.}{\sim} N(\mu, \sigma^2)$
2. Observations come from a normal distribution
3. Two sample  $t$ -based Assumptions also assumes the two groups are independent. There is no covariance in either equation. We assume that the observations are independent both within the groups and with each other.  $X_{1,i} \stackrel{i.i.d.}{\sim} N(\mu_1, \sigma_1^2)$  and independently  $X_{2,i} \stackrel{i.i.d.}{\sim} N(\mu_2, \sigma_2^2)$
4. Pooled also assumes that the variances are the same.  $\sigma_1 = \sigma_2$  and that the samples are independent.

## 10.2 Robustness

The performance of an inferential procedure when the assumptions fail is called **robustness**.

This is measured by:

**Type I error** probability of incorrectly rejecting the null hypothesis when it is actually true. False positive rate. Fixed typically  $\alpha = 0.05$ . If the assumptions are wrong, especially with small sample sizes, the probability of Type I errors increases dramatically.

**Power** probability of correctly rejecting the null hypothesis when the alternative is actually true. Usually, the alternative has a range of values so a value has to be picked. Defined as  $(1 - \text{Type II error})$ .

Both types of error concern rejecting the null, but under different conditions.

### 10.3 Breakdown of Independence Assumption

The assumption requires independence within the group and between the groups. If between-group independence is violated, we can do a paired  $t$ -test.

Based on collecting data, if correlations are related within groups, then there is a breakdown of independence assumption.

For example collecting information about gender, but both samples of gender are from the same peer group. **Clustering** is when subgroups of units are similar to each other.

If they are centered around different averages, then regression will allow us to “control for” it.

There are extensions of linear regression to control for this grouping based on the type of correlation:

**Serial effect** Dependence over time

**Spatial effect** Interference across space

You can only take these clusters into account if you can measure how these clusters are close to each other either in time or space. Think carefully how was the data was collected. Did it all come from the same neighborhood?

**When independence assumption is violated then the Sampling Variance calculations will be incorrect.** When two results are always the same, the sample size is artificially inflated that may lead to incorrect statistically significant results. More advanced calculations are needed or redefine units. (See textbook, chapter 15). Example: selecting NFL team stats over multiple years and treating as independent data, when it isn't.

Use of graphics and plots can help identify these clusters. Split the data into batches manually and see if there is correlation (e.g. boxplot, periodicity). Check across space and time to see if there are patterns.

### 10.4 Breakdown of Normality Assumption

Use visual examination of shapes (left or right skews, symmetry (long or short tailed), unimodal / multimodal) of the sample distribution. Does the sample distribution follow a normal distribution?

Overlay the density plot with the **kernel density** and the normal curve (for example, fitted by the method of moments, MOM). If the two curves fit reasonably well, it is likely a normal distribution.

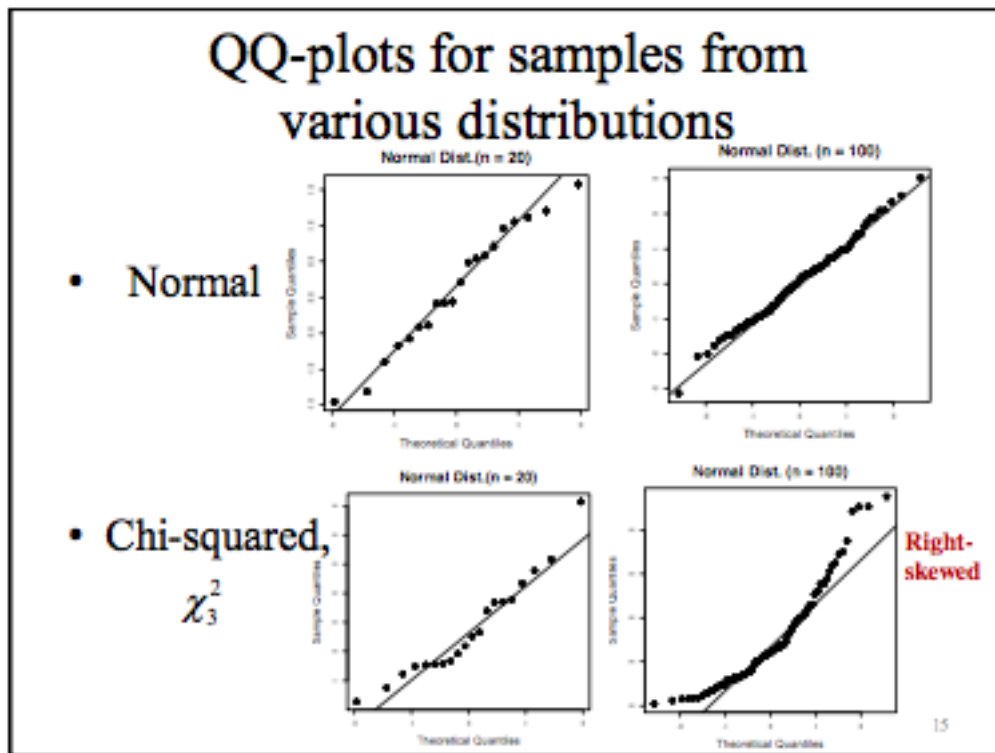


Figure 1: QQ Plots

In R,

```
x <- data$tuition[data$public==1]
hist(x, col="light gray", main="Public", xlab="Tuition, in $1000",
     prob=TRUE, breaks = 10, cex.main=2,cex.lab=2)
# Kernel density
lines(density(x, adjust = 2), col = "blue", lwd=2)
# Approx. normal curve, fitted using MOM
points(seq(min(x), max(x), length.out=500),
       dnorm(seq(min(x), max(x), length.out=500), mean(x), sd(x)),
       type="l", col="red", lwd=2, lty=2)
```

Use a boxplot show outliers.

Best of all, use a quantile-quantile (QQ) plot which should generate a straight line when both distributions are normal as well as chunking of points around 0 and less towards the extremities. Good for showing distributions away from normality that a histogram wouldn't

necessarily show. When the sample size is small (e.g. 20), even a normal distribution will not look normal on a QQ plot. R: `qnorm()`, `qqline()`. QQ plots can contrast any distribution like  $\chi^2$  (see homework 3). Samples should follow a  $\chi^2$  distribution. Even when sampling from a normal distribution, the plot might be skewed, especially when the sample size is small.

#### 10.4.1 Analytical Tests for Normality

There exists analytical tests for normality such as Anderson-Darling or Shapiro-Wilk. However, use **caution!** When there are only small sample sizes (e.g.  $< 20$ ), these tests have low power and may not even reject for a non-normal sample.

For large samples, these tests will always reject, but we know from the Central Limit Theorem, that the larger the sample, the more likely it is to be normally distributed, irrespective of the population distribution.

The non-normality should be assessed relative to the problem, which implies using plots.

#### 10.4.2 Uniform Distribution

$X \sim Unif(a, b)$ ,  $a < b$  means that the values are evenly distributed in the interval  $[a, b]$  so all values are equally probable.

Drawing the PDF, shows a rectangle with length  $b - a$  and height  $\frac{1}{b-a}$  (since all PDFs have to have an area = 1.) So  $y = f(x) = 1/(b - a)$ .

Standard Uniform Distribution:  $a = 0, b = 1$ .

#### 10.4.3 The Universality of the Uniform Distribution

If you take random samples from a uniform distribution, the samples should be normally distributed  $X_i \sim N(\mu, \sigma^2)$ . Then if you calculate one sample  $t$ -test from that sample population  $T = \frac{\bar{X} - \mu}{S/\sqrt{n}}$ , the resulting  $T$  values should follow a  $t$  distribution  $T \sim t_{df=n-1}$ .

If you then generate  $p$ -values based on the  $t$ -test, the resulting  $p$ -values should have a **uniform distribution**. If the  $t$ -values do not follow the  $t$  distribution, then the resulting  $p$ -values will not be uniformly distributed, and our assumptions are incorrect.

See sample code for Unit 5 to see how to sample from a distribution to verify this.

Since  $\alpha = 0.05$ , we will reject 5% of the time. If 5% of the p-values are less than 0.05, then we can conclude that the test is valid. Changing sample sizes (to even small sample sizes) will not change the distribution of the pooled case.

#### 10.4.4 Robustness of $t$ -test to Departures from Normality

The two-sampled  $t$ -test is robust under moderate skewness as long as the sample size is large. (See Example 3 in Unit 5 code).

Sample goes for Exponential distribution. As long as the sample size is large (e.g.  $\geq 50$ ), the  $t$ -test will still be robust. (See Example 5, in Unit 5 code). When the sample size is small (e.g. 10), we get 6% of rejection, increasing our Type I errors. (See Example 6).

Note that mean and standard deviation had to be coerced into 0, 1 which isn't typical for other distributions like exponential.

In summary,  $t$ -tests are fairly robust to departures from normality, especially in large samples (CLT). When the sample sizes are not equal,  $t$ -tests are more sensitive to skewedness and long-tailedness. For small samples,  $t$ -tests are somewhat sensitive to markedly different skewedness in two groups. Watch out for outliers.

When normality assumption is violated:

- $t$ -test is usually still valid, or
- Use data transformation, or
- Use non-parametric test (Unit 6).

#### 10.5 Equal Variance Assumption

For the Pooled  $t$ -test, we need to verify the Equal Variance Assumption.

When sample sizes are equal, the pooled  $t$ -test is fairly robust to unequal variances.

When sample sizes are unequal, the pooled  $t$ -test is typically not valid for unequal variances; the unpooled  $t$ -test is a robust alternative.

When the equal variance assumption is violated:

- Use unpooled  $t$ -test; or
- Use data transformation; or
- Maybe the populations are not comparable?

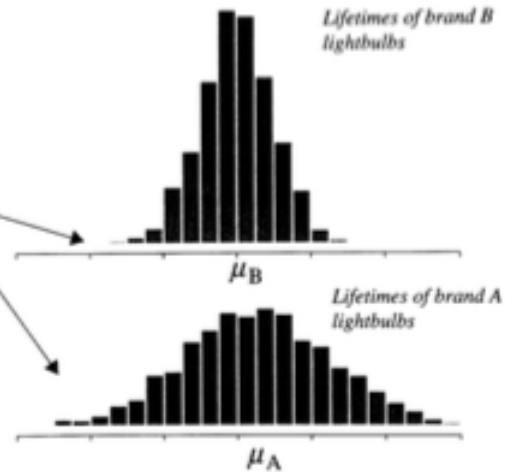
## Maybe the populations are not comparable?

**DISPLAY 4.11**

The conceptual difficulty with comparing population means when population spreads are not the same

*On average, brand A bulbs last longer; but there is also a greater chance of early burnout with brand A.*

*The question of which brand is better may be more complex than simply "Which mean is larger?"*



30

Figure 2: Maybe the populations are not comparable?

It isn't enough to compare the means of the two populations - it is more important to compare the lower tail of the population having equal distribution. E.g. minimizing the number of light bulbs that burn out.

(See Unit 5 code, example 7). Recall that the ratio of the variances should be  $< 2$  (always put larger on top) in order to use the pooled test.

The larger group, with the larger variance give p-values closer to 1, whereas the smaller group with the smaller variance brings p-values closer to 0. See the formula for pooled  $t$  test:

$$t = \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{S_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$$

Depending which sample we decide is  $X_1$  or  $X_2$  will affect the skewness of the  $p$ -value distribution.

Having a very high Type I error rate, like 37% (See Example 9) is very bad for the quality of our model. So, either make the group with the larger variance be  $X_1$  or simply and preferably use an unpooled test.

### 10.5.1 Unpooled (Welch) Two-Sample $t$ -test

This is what R uses to calculate degrees of freedom. Not needed for this class, but the formula is on Unit 5, slide 35.

### 10.5.2 Graphical Check of Equal Variances

Use box plots. Look out for outliers! (since variance is sensitive to outliers).

In practice, Statisticians use the unpooled test most often as it is more robust.

## 10.6 Formal $F$ -test of Equal Variance

Start with graphical check and look at the *ratio of sample variances*.

Alternatively, use the  $F$ -test for comparison of two variances. In R, `var.test()`.

Same assumptions as for  $t$ -tests (i.i.d, Normality)

$$H_0 : \sigma_x^2 = \sigma_y^2$$

$$H_a : \sigma_x^2 \neq \sigma_y^2, \sigma_x^2/\sigma_y^2 \neq 1$$

Test Statistic:

$$F = \frac{S_x^2}{S_y^2} \sim F_{n_x-1, n_y-1}$$

**Caution:** will reject for large samples, even when the ratio is very close to 1.

### 10.6.1 $F$ -distribution

Let  $X \sim \chi_{n_x}^2$  and  $Y \sim \chi_{n_y}^2$ , independent of each other.

$F$  is distributed as the ratio of the two variances.

$$F_{n_x-1, n_y-1} = \frac{\sigma_x^2 \cdot \frac{\chi_{n_x-1}^2}{n_x-1}}{\sigma_y^2 \cdot \frac{\chi_{n_y-1}^2}{n_y-1}}$$

Under the null hypothesis,  $\sigma_x^2/\sigma_y^2 = 1$ , so this simplifies to:

$$F_{n_x-1, n_y-1} \sim \frac{\frac{\chi_{n_x-1}^2}{n_x-1}}{\frac{\chi_{n_y-1}^2}{n_y-1}}$$

If this is true, then the following ratio has an  $F$ -distribution with  $n_x$  and  $n_y$  degrees of freedom.

$$R = \frac{X/n_x}{Y/n_y} \sim F_{n_x, n_y}$$

The above assumes a population. Generally, this is written as degrees of freedom for a sample, ie.

$$R = \frac{X/(n_x - 1)}{Y/(n_y - 1)} \sim F_{n_x-1, n_y-1}$$

As a ratio of variances, it will be positive, usually heavily right skewed.



(Later: very handy for ANOVA and model selection in Linear Regression)

(Don't worry about Two-sided  $p$ -values for  $F$ -test, Unit 5, slide 40).

## 11 Unit 6 - Transformations

Section 3.5 and Chapter 4 in the text.

We use transformations when assumptions of whatever test we are using (e.g.  $t$ -test) are violated.

### 11.1 Non-Linear (Log) Transformations

Linear transformations are in the form  $ax + b$ . Nonlinear transformations are square root, natural log, exponent, etc.

The log transformation is the most common transformation because it flattens exponential scales. It turns a multiplicative effect (which is very right skewed) and into an additive effect (which is more symmetrically skewed). It may also fix very large variance differences. This does not work for square root transformations.

Right skewed data cannot be used for a  $t$  test, whereas symmetric data can be subject to a  $t$ -test.

Most common choice is natural log transformation (for ease of derivation):

$$Z_i = \ln(Y_i), i = 1 : n$$

Used when data is very skewed to the right ("positively skewed") or spread is larger in a group with a larger center.

The ideal result after the transformation is two symmetric samples with similar spreads but possibly different centers.

### 11.1.1 Interpretation - Randomized Experiments

Let  $T$  be the treatment group, and  $C$  be the control group.

$$\begin{aligned} Z_{T,i} &= \ln(Y_{T,i}), i = 1 : n_T \\ Z_{C,i} &= \ln(Y_{C,i}), i = 1 : n_C \\ Z_{T,i} &= Z_{C,i} + \delta, \text{ which is equivalent to} \\ \frac{Y_{T,i}}{Y_{C,i}} &= e^\delta \end{aligned}$$

For unit  $i$  responses, you have to exponentiate the difference:

$$\exp(\bar{Z}_T - \bar{Z}_C) \text{ estimates } \frac{Y_{T,i}}{Y_{C,i}}$$

### 11.1.2 Interpretation - Observational Studies

$$\begin{aligned} Z_{1,i} &= \ln(Y_{1,i}), i = 1 : n_1 \\ Z_{2,i} &= \ln(Y_{2,i}), i = 1 : n_2 \end{aligned}$$

For symmetric distributions,  $E(Z_{j,i}) = \text{Median}(Z_{j,i})$ . If the distribution is symmetric on the log transform scale, the means and medians are the same. When transforming back to the exponentiated scale, the mean will be messed up, but the medians will not be affected.

This is called *Monotonicity of logs*:  $\text{Median}(\ln(Y_{j,i})) = \ln(\text{Median}(Y_{j,i}))$

Interpretation in terms of *a ratio of population medians*:

Let  $m_1 = \text{Median}(Y_{1,j})$  and  $m_2 = \text{Median}(Y_{2,j})$ , then

$$\exp(\bar{Z}_2 - \bar{Z}_1) \text{ estimates } \frac{m_2}{m_1}$$

So the median of the second population is  $\exp(\bar{Z}_2 - \bar{Z}_1)$  times as large as the median of the first population.

### 11.1.3 Log Transform for Paired $t$ Test

$$Z_i = \ln(Y_{2,i}) - \ln(Y_{1,i}), i = 1 : n$$

take differences in logged outcomes within pairs **before** averaging!

In Randomized Experiments:  $Y_{2,i}/Y_{1,i}$ .  $\exp(\bar{Z})$  estimates a multiplicative effect treatment.

In Observational Studies: Let the median of ratios be  $m$  ( $\neq$  ration of medians!),  $\text{Median}(Y_{2,i}/Y_{1,i}) = m$ , then  $\exp(\bar{Z})$  estimates  $m$ .

### 11.1.4 Graphical Interpretation

Look at distribution of log scale, QQplots. Outliers are not gone, but they are minimized.

### 11.1.5 Interpretation of Exponentiated Scale

Recapping the hypothesis test, we want to see the result on the exponentiated scale. Starting with log scale, we use a hypothesis test to check if the true (log) means are equal:

$$\begin{aligned} H_0 : \mu_1^* &= \mu_2^* \\ H_a : \mu_1^* &\neq \mu_2^* \\ t &= \frac{(\bar{Z}_1 - \bar{Z}_2)}{\sqrt{\frac{S_1^{2*}}{n_1} + \frac{S_2^{2*}}{n_2}}} \end{aligned}$$

( $\bar{Z}$  is mean of data on log scale)

Example: The above hypothesis test produces a very low  $p$ -value and a (logarithmic) 95% confidence interval of (0.24, 0.72) for ratio of  $\frac{m_1}{m_2}$ . So the exponentiated confidence interval is  $(e^{0.24}, e^{0.72})$  or (1.27, 2.05). This result makes sense because the interval does not 1. A multiplicative factor of 1 on the log scale is equivalent to a 0 in CI of non-log scale.

### 11.1.6 Limitations of the Log Transformation

Log scale does not work well for data with many small values [why?], especially zeros - best if all values are  $> 1$ .

$$Y_{j,i}^* = Y_{j,i} + \epsilon, Z_{j,i} = \ln(Y_{j,i}^*)$$

If there are not many zeros, one can shift all observations by a small number (say, 0.01, or 0.1).

Cannot be used with negative data.

It doesn't matter what base of log is used - the shape of the transformation will always be the same. Sometimes  $\log_2$  or  $\log_{10}$  are used to make interpretation better: allows for use of a doubling or a ten-fold increase for a one-unit change in the log-transformed variable.

### 11.1.7 Other Types of Transformations

If the log transformation is too strong (turns right skewed into left skewed), other transformations can be applied.

**Square root** or other polynomial

- Good for moderately right-skewed measurement data (e.g. counts, area sizes, etc.)
- Difficult to interpret results
- More appropriate for data values  $< 1$  and  $> 1$
- Using  $Y^2$  if there is left-skewed data

**Reciprocal** taking  $1/Y$

- For severely skewed data (waiting or failure times)
- Good for left skewed data
- Large values become small values.
- Often use negative reciprocal
- Can be used with negative data

**Logit**  $\log(\frac{Y}{1-Y})$ , or  $\arcsin(2Y - 1)$

- Good proportions or percentages to transform to real line
- Used with logistic regression

Could also flip axes to turn left skewed into right skewed data (e.g turn literacy rate into illiteracy rate).

## 11.2 Choosing a Transformation

Transformations, other than  $\log(Y)$  are difficult to interpret.

Choices are made by considering the nature of the data (e.g. counts, waiting times, proportions, etc.)

Plot the *transformed* data and assess where the result conforms with the assumptions of a chosen test (equal spread, normality, outliers, etc.).

Financial data often benefits from log transformations since the underlying data is often distributed in log-normal form (it is normal after taking the log).

## 11.3 Summary: Evaluating Validity of $t$ -tools for a Problem

Evaluate independence assumption by considering the data collection method.

Use graphs to evaluate normality, similarity of shapes equality of variances and outliers

If needed

- transform data
- consider possible justifications for removing outliers
- or use robust  $t$ -tools (e.g. unpooled test)

Be cautious about removing outliers. Do not remove outliers unless there is justification to remove the outlier.

## 11.4 Alternatives to the $t$ -tools

aka Nonparametric Tests based on Ranks.

Data is ordered and ranked from smallest to largest.

Use if sample sizes are small (e.g.  $< 30$  for bell-shaped data or ever larger if skewed), or there are still outliers (and transformations did not help).

Transform all observations based on **ranks**. For two independent samples, use the **Rank-Sum Test** (aka Mann-Whitney or Wilcoxon Test).

For paired-samples, use the **Sign Test** or **Wilcoxon Signed-Rank test**.

## 11.5 Nonparametric Tests

**Parametric** procedure makes an assumption about underlying distribution of the observations (e.g.  $t$ -tests assume Normal Distribution).

**Nonparametric** procedure makes no assumption of the data generating process and therefore does not assume the observed data follows a specific distribution. Example: randomization and permutation tests

## 11.6 Rank-Sum Test

Suppose we have two samples  $X$  and  $Y$  with different sizes drawn independently from two populations where  $n_x \leq n_y$ .

Ranking could be thought of as a transformation of the underlying data.

First the samples are transformed by **Rank** (e.g. sorted) into a combined set,  $Z$ . Ties (or duplicate values) are handled by averaging the corresponding Ranks (e.g. Rank 1, 2 becomes Rank 1.5)

Second, a permutation test is performed on ranks using the following test statistic:

$$T = \sum_{i_x=1}^{n_x} Z_{1,i}$$

which sums up the ranks for each group. Assumes  $n_x \leq n_y$ .

Then compare  $T$  for each group to see if they are close enough. If there is a difference in sample size, the rank sum has to be proportioned by sample size.

Important: See Allocation of Units to Groups, Unit 6, slide 21.

### 11.6.1 Rank-Sum Hypotheses

$H_0$  and  $H_a$  are essentially the same as those for a permutation (or randomization test). Rank-sum is preferred if there are *censored observations* (not available due to some value being exceeded). Permutation test is preferred if there are many ties.

If the shapes and spreads of the two populations are similar:

**Permutation test** (using  $\bar{Y}_2 - \bar{Y}_1$ ),  $H_a$ : there is a difference in the *averages* between the two populations.

**Rank-Sum test** (using  $T$ ),  $H_a$ : there is a difference in the *medians* between the two populations.

### 11.6.2 Exact Sampling of Rank-Sum Distribution

Because  $n_x$  and  $n_y$  are known in theory, we can calculate  $T = \sum_{i=1}^{n_x} Z_{1,i}$  for *all possible regroupings* and obtain the *exact p-value*.

Under  $H_0$ , and no ties, the ranks are Discretely Uniformly distributed,

$$Z_{ji} \sim \text{DUnif}(1, 2, \dots, n_x + n_y)$$

### 11.6.3 Rank-Sum: Normal Distribution of Sampling

Since a sum of independent values has a normal distribution, and the ranking is independent [why?], then rank-sum can be approximated by a normal distribution.

For moderate or larger sample size (e.g.  $n_i > 10$ ), under  $H_0$  and no ties,

$$T \sim N\left(\frac{n_x(n_x + n_y + 1)}{2}, \frac{n_x n_y (n_x + n_y + 1)}{12}\right)$$

If there are ties, let  $\bar{R}$  be the sample mean and  $S_R^2$  be the sample variance for a combined set of  $n_x + n_y$  ranks, then

$$T \sim N\left(n_x \bar{R}, S_R^2 \frac{n_x n_y}{n_x + n_y}\right)$$

then the following  $Z$  statistic could be used (irrespective of ties):

$$Z = \frac{T - n_x \bar{R}}{S_R \sqrt{\frac{n_x n_y}{n_x + n_y}}}$$

So where do these means and variances come from?

Ranks =  $Z = 1, 2, 3, \dots, n$  where  $n = n_x + n_y$ . So what is the total sum of Ranks?

It is a simple sum of the first 10 integers, ie.  $\sum_{i=1}^n i = \frac{n(n+1)}{2}$

If  $H_0$  is true, what proportion of those total sum of Ranks, should be attributed to Group X? It should be in proportion to the number of samples in that group ( $n_x/n$ ).

Given that  $T = \sum_{i=1}^n Z_i$ , so

$$\begin{aligned} E(T) &= \frac{n_x}{n} \times \frac{n(n+1)}{2} \\ &= \frac{n_x(n+1)}{2} \end{aligned}$$

So what is the Variance? Recall the basic formula for variance:

$$Var(T) = E(T^2) - [E(T)]^2$$

Using basic algebra, sum of first  $n$  integers squared:

$$\sum_{i=1}^n i^2 = \frac{n(n+1)(2n+1)}{6}$$

so taking this proportion to the total Rank:

$$\begin{aligned} E(T^2) &= \frac{n_x}{n} \times \frac{n(n+1)(2n+1)}{6} \\ &= \frac{n_x(n+1)(2n+1)}{6} \end{aligned}$$

Bringing this back to the variance:

$$\begin{aligned} Var(T) &= \frac{n_x(n+1)(2n+1)}{6} - \left( \frac{n_x(n+1)}{2} \right)^2 \\ &= \frac{2n_x n^2 + 6n_x n + 2n_x - n_x^2 n^2 - 2n_x^2 - n_x^2}{12} \end{aligned}$$

[TODO: how to simplify?]

If there are ties, we need to change the mean and variances appropriately.

Why do we care? We want to build a reference distribution for the test statistic  $T$  in order to generate a  $p$ -value.

#### 11.6.4 Example of Rank-Sum using R

```
library(coin)
chocolate = c (0,5,9,2,1,3,3,4,4,2,0,5)
female = c(0,0,1,1,1,0,0,0,0,0,1,1)
rank(chocolate)
sum(rank(chocolate)[female==0]) # = 46.5
sum(rank(chocolate)[female==1]) # = 31.5
```



highest value was 9, but after rank it is now  $n$  ( $= 12$ ).

No matter shapes and distributions, the hypotheses are  $H_0$  : distributions of chocolate eaten is the same for men and women.  $H_a$  : distributions are different.

If shapes and distributions are similar:  $H_0 : Med_f = Med_m, H_a : Med_f \neq Med_m$

$$H_0 : \sum_{men} = \text{Ranks} = 46.5$$

$$H_a :$$

How do we know if this is significant? Use the normal distribution approximation. From R's calculation,  $\bar{R} = 6.5, S_R = 3.574$

Recall:

$$T \sim N\left(n_x \bar{R}, S_R^2 \frac{n_x n_y}{n_x + n_y}\right)$$

Then  $Z$ -score the  $H_0$  sum.

$$\begin{aligned} n_x &= 7 \\ \bar{R} &= 6.5 \\ Z &= \frac{\bar{X} - n_x \bar{R}}{\sqrt{S_R^2 \left(\frac{n_x n_y}{n_x + n_y}\right)}} \\ &= \frac{46.5 - 7(6.5)}{\sqrt{3.574^2 \frac{7 \cdot 5}{7+5}}} \\ &= 0.164 \end{aligned}$$

In R

```
p = 2*(1-pnorm(0.164))
```

Conclusion:  $p \approx 0.87$

Alternatively, you could run a permutation test. See sample code for both this and the built-in Wilcox test.

### 11.6.5 Rank-Sum Summary

Rank-Sum is useful when

- There are no distribution assumptions.
- The only assumption is independence of units.
- There is censoring of values. (e.g. someone who survived more than  $n$  years is just considered  $n$ )
- There are outliers.

See text 4.2.4 for explanation on how to obtain confidence intervals for a treatment effect  $\delta$  in randomized experiments.

### 11.7 Sign Test

This is useful for Paired Data. Aka binomial test.

The Rank Sum test is a non-parametric test for a quantitative response variable for two independent groups. A similar approach can be taken when the data are paired, i.e., when there are two observations taken on the same study unit or two units are paired together naturally.

A common type of paired data used in scientific research is from twin studies. If you are a twin, you have likely been recruited for many medical studies. Twins are good study units because they control for genetics and possibly environment very easily.

P: Probability of schizophrenia when unaffected twin has a larger hippocampus than affected twin.

$$\begin{aligned}H_0 : p &= 0.5 \\H_a : p &\neq 0.5 \\X &= 14 \\X &\sim \text{Binom}(p, n = 15) \\p\text{-value} : P(X \geq 14) + P(X \leq 1)\end{aligned}$$

Recall: in a binomial distribution,  $E(X) = np$  and binomials are discrete. Be sure to include the end-points when using the Sign Test.

In R:

```
dbinom(14, p=0.5, size=15)
dbinom(14, p=0.5, size=15) + dbinom(15,p=0, size=15 + dbinom(1,p=0.5,size=15) +
    dbinom(0,p=0,size=15)
```

Suppose we have  $n$  independent pairs of Random Variables  $(X_1, Y_1) \dots (X_n, Y_n)$ .

$D_i = Y_i - X_i$  and use the following test statistic:

$$K = \sum_{i=1}^n I(D_i > 0)$$

ie. a number of pairs, when  $Y_i > X_i$ .

Zeroes are dropped, which may be a problem.

$K$  is a binomial distribution.

### 11.7.1 Sign Test Hypotheses

Most general formulation:

$$\begin{aligned} H_0 : P(D_i > 0) &= P(Y_i > X_i) = 0.5 \\ H_a : P(D_i > 0) &\neq 0.5 (or > 0.5, or < 0.5) \end{aligned}$$

There are two other potential hypotheses. See Unit 6, slide 38 for continuous and symmetric random variables. Also see textbook 4.4.1.

### 11.7.2 Exact Calculations and Normal Distributions

Under  $H_0$ , the exact distribution of  $K$  is binomial,  $K \sim Bin(m, 0.5)$  where  $m$  is the final number of pairs with nonzero  $D_i$ .

If samples are small, use the real distribution based on the binomial distribution, which is more precise than using the Normal approximation.

For larger samples, we can use a Normal approximation:

$$Z = \frac{K - m/2}{\sqrt{m/4}} \sim N(0, 1)$$

In R, use `sign.test()` in the BSDA package. [Example](#).

## 11.8 Wilcoxon Signed Rank Test

Like the Rank-Sum test for paired data.

First calculate the differences, then calculate the magnitudes of the differences, then rank the magnitudes.

In other words, after calculating the differences within pairs, we rank their absolute values  $|D_i|$

$$Z_i = \text{Rank}(|D_i|), i = 1, \dots, n$$

The test statistic is the sum of ranks for positive differences:

$$S = \sum_{i=1}^n Z_{i|D_i>0}$$

In R, use `wilcox.test(..., pair=TRUE)` and maybe `conf.int=TRUE`.

### 11.8.1 Hypotheses

Most general formulation of hypotheses:

$H_0$  : the rank of the magnitude of within-pair difference is *unrelated* to the sign of the difference

$H_a$  : the rank of the magnitude of within-pair difference is *related* to the sign of the difference

See Unit 6, slide 43 for continuous and symmetric hypotheses.

### 11.8.2 Exact Calculations and Normal Approximations

Exact sampling distribution of  $S$  may be generated by

- randomly switching the group status but keeping the same observations within each pair;
- ranking absolute differences for new group assignments;
- calculating a new test statistic  $S^*$ .

Under combinatorics, it can be shown that under  $H_0$ :

$$Z = \frac{S - m(m+1)/4}{\sqrt{m(m+1)(2m+1)/24}} \sim N(0, 1)$$

for  $m \geq 20$ , where  $m$  is a number of non-zero differences.

If the sample size is smaller than 20, do a permutation or re-sampling test instead. Assign the sign by randomly generating the sign.

In R:

```
ranks=1:15
ranks[9] = ranks[10] = 9.5
rbinom(15,0.5,size=1) # simulate coin flips
sum(ranks[rbinom(15,0.5,size=1)==1]) # repeat for multiple simulations
# and then count numbers of simulations exceeding the critical values in the two-sided test
```

## 11.9 Summary of Paired Test Choices

There are four options: Paired  $t$  test, Transformed  $t$  test, Sign Test, Wilcoxon Signed Rank Test

Choose the most valid test. Choose  $t$  or Transformed  $t$ -test **unless the assumptions are not met**. If the assumptions are not met, usually use Wilcoxon Sign Rank over Sign Test as it has higher statistical power. Alternatively, run a permutation or re-sampling test (Lecture 6, part 2, 9:00)

## 12 Section 6

HW5Q3, simulations. Frequentist approach.

```
# set up problem
n.sim = 10000

#initialize vectors with correct size
x1 = rep(NA, 50)
x2 = rep(NA, 50)
```

put in cheatsheet appropriateness of test - see section notes

Use R Markdown to embed R output into L<sup>A</sup>T<sub>E</sub>X

## 13 Unit 7 - Power, Sample Size and Error

Textbook 4.5 and 6.3

### 13.1 Type I and II Error and Power

Recall:

Type I: reject null when we shouldn't have. Fixed typically at  $\alpha = 0.05$ .

Type II: fail to reject null when it is false. Not typically fixed.

		Truth about the population	
		$H_0$ true	$H_a$ true
Decision based on sample	Reject $H_0$	Type I error	Correct decision
	Accept $H_0$	Correct decision	Type II error

Figure 3: Error Type Terminology

#### 13.1.1 Statistical Power

$P(\text{type II error}) = P(\text{Conclude } H_0 \text{ when } H_a \text{ is true})$  is labeled  $\beta$ .

In decision theoretic approach to hypothesis testing:

- the type I error rate  $\alpha$  is fixed (usually 0.05)
- $1 - \beta = P(\text{Conclude } H_a | H_a) = P(\text{correct decision when } H_a \text{ is true})$

$1 - \beta$  is called the **statistical power** of the test. Computing this probability can be subtle, and depends closely on particular problem. It is done *before* you collect the data. It is intimately tied with determining sample size. If this is true, what is the probability we can show it?

Note: Power is always between 0 and 1 because it's a probability. Closer to one is better.

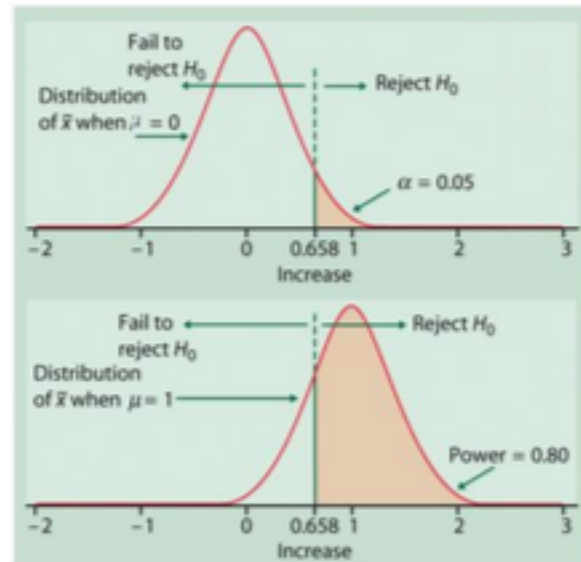


Figure 4: Statistical Power - Null on top, Alternative on bottom

In order to calculate you need one specific value, e.g.  $\mu = 1$  (arbitrarily coming from the domain expert) to create a rejection region. Either the null is true and we use the distribution on the top, or the alternative is true and we use the distribution on the bottom.

Power is affected by:

- effect size – how different is the alternative from the null hypothesis?
- the variance. If the variance increases, which decreases power
- sample size. Increasing sample size which increases power.

## 13.2 Sample Size and Power Calculations

Example: we know the standard deviation of financial aid  $\sigma = 20000$ .

Suppose we wish to design the study so that the margin of error in a 95% confidence interval for the mean is 2000, ie sample mean  $\pm 2000$ . How large does the sample size need to be?

Recall: the formula for margin of error ( $m$ ) in a confidence interval for a mean: If we know the population variance,

$$m = z^* \left( \frac{\sigma}{\sqrt{n}} \right)$$

Otherwise, we use the sample variance:

$$m = t^* \left( \frac{s}{\sqrt{n}} \right)$$

We assume we know the population variance because we don't have the data yet. We are picking an acceptable value based on the problem constraints.

$$n = \left( \frac{(z^*)(\sigma)}{m} \right)^2 = \left( \frac{1.96(20000)}{2000} \right)^2 = 384.16 \approx 385$$

Calculating Power is a two step process:

1. Determine the rejection region assuming is true (what values of  $\bar{X}$  are needed to reject the test).
2. Calculate the probability of finding a sample statistic (in this case the sample mean) that falls in the rejection region if the alternative hypothesis is actually true. We will need a specific estimate of  $\mu$  under  $H_a$ .

Example, continued:

$$H_0 : \mu = 24000$$

$$H_a : \mu \neq 24000$$

1. Rejection region is:

$$\bar{X} > \mu_0 + 1.96 \frac{\sigma}{\sqrt{n}} = 24000 + 1.96 \frac{20000}{\sqrt{385}} = 25998$$

$$\bar{X} < \mu_0 - 1.96 \frac{\sigma}{\sqrt{n}} = 24000 - 1.96 \frac{20000}{\sqrt{385}} = 22002$$

See Unit 7 slide 13.



2. Choosing a  $\mu$  within the rejection region ( $H_a$ ), say 27000.

$$\begin{aligned} &P(\bar{X} > 25998 | \mu_a = 27000) + P(\bar{X} < 22002 | \mu_a = 27000) \\ &\approx P(\bar{X} > 259998 | \mu_a = 27000) + 0 \\ &= P\left(\frac{\bar{X} - \mu_a}{\sigma/\sqrt{n}} > \frac{25998 - 27000}{20000/\sqrt{385}}\right) \\ &= P(z > -0.98) \\ &= 1 - 0.1635 = 0.8365 \end{aligned}$$

In R use

```
z=(25998 - 27000) / (20000 / sqrt(385))
1-pnorm(z)
```

### 13.3 Find Sample Size for a Desired Power

Assuming you have a desired power and then want to find the appropriate sample size:

Standardize  $\bar{X}$  assuming  $H_a$  is true. See formulas in Unit 7, slide 15. You need to know (or determine)  $\sigma$  as well as the effect size. Both are based on the domain expert. In practice, these values are bumped up 10% because 1) this is a  $t$  not normal distribution, 2) the experts assume they have more precision than they really do.

### 13.4 Power and Sample Size Caveats

The calculation of Power and Sample Size for a study are tied hand-in-hand.

Oftentimes in practice, if the planning is actually done, power is chosen to be around 0.80 or 0.90 for a specific value ('effect size') in the alternative hypothesis, and the sample size is then determined based on that calculation.

Warning: this calculation has its merits and its flaws: 1) It's automatically ad hoc: we have to make an educated guess of the effect size and of the true variance,  $\sigma$ .

2) Should be done using the analysis that is planned to be carried out...this gets very complicated very quickly. So often a table (or sets of tables) is produced for various levels of power (or sample size) and various effect sizes in the alternative world, and the scientific expert then chooses based on these properties. When the calculation gets complicated, often Monte Carlo simulations are used to do the calculation.

## 13.5 Using R

These calculations are based on a z-test (assuming we know the true standard deviation), but in practice we'll use the t-test.

Why does this make the calculation more difficult? Need to incorporate the randomness of using  $s$ , plus the  $t$  critical values will change depending on  $df = n - 1$

In R use in the package `samplesize`:

```
power.t.test(n, delta, sd, sig.level, power,
  type = c("two.sample", "one.sample",
    "paired"), alternative = c("two.sided",
    "one.sided")
n.ttest(...)
```

end of lecture 6.

## 13.6 Tradeoff: Assumptions vs. Power

lecture 7

How to increase power of a study (if all assumptions are met):

1. Increase  $\alpha$ . Not useful if  $\alpha$  is too high
2. (not applicable in practice) - conduct a one-tailed test
3. Increase the anticipated effect size of interest
4. Increase sample size
5. Improve precision

Refer to Power instead of type I error because we are examining the likelihood of rejecting the **alternative** hypothesis.

When variances are the same, use the same sample size.

### 13.6.1 Pooled vs. Unpooled $t$ -Tests

From most to least powerful, if population variances are the same, use Pooled, then Unpooled. (Pooled assumes same normal distribution and same variance). It is more conservative to use unpooled, since it is safer and not much loss in power.

Power gain is especially obvious when sample sizes are different.

However, the difference vanishes as  $n_x$  and  $n_y$  grow.

### 13.6.2 Randomization vs. $t$ -Test

Permutation tests are preferred to  $t$ -tests since there are no assumptions (non parametric). However, the alternative may not be informative [why?]

Choose  $t$ -tests over permutation tests if there are clearer null and alternative hypotheses. The  $t$ -test will have higher power *if assumptions are met*.

So in general, nonparametric tests are somewhat less powerful than  $t$ -tests, when normality and equal variance assumptions are met.

If assumptions of  $t$ -tests are met, then a “general” order from most to least power is:

1.  $t$ -tools
2. Permutaton test
3. Rank-sum test (or signed-rank test for paired data)
4. Sign test, or equivalent like Fisher’s Exact Test

For large samples, rejection rates and power are practically equivalent.

### 13.6.3 Permutation Test

The permuation test is sensitive to differences in variances.

Permutation test’s typical null hypothesis is  $H_0$  : distributions are equal in two groups. But this assumptions could be violated by change in spread or symmetry.

If the two distributions has similar shapes and spreads, then the null is  $H_0 : \mu_1 = \mu_2$

### 13.6.4 When Assumptions are Violated

Example: comparing uniform vs. exponential distribution. Every time there is an outlier (which is frequent for exponential distributions), then  $H_0$  will be rejected every time.

(See Unit 7, slide 40 for a diagram)

Use the Rank-Sum test in such a case.

## 13.7 Power Curves

Power is a function of:

1. Choice of test
2.  $\alpha$
3. Sample size
4. Effect size

We can compare the different tests for varying levels of sample size (sometimes done) or effect size (often done) to see how they compare (since  $\alpha = 0.05$  is usually fixed).

Plotting power as a function of  $n$  or, more commonly, the effect size is called the **power curve**.

The  $x$ -axis is  $\mu_y - \mu_x$ , centered at 0, which is the null hypothesis. The  $y$ -axis the probability of rejecting  $H_0$ . The lowest  $y$ -value is  $\alpha = 0.05$ . The highest  $y$ -value is 1.

If it is one-sided, ie.  $H_a : \mu_y > \mu_x$ , will be ‘S’ shaped, approaching 0 to the left (depending on sign) and 1 on the right, at a slightly faster rate than the two sample case, because the threshold is lower for one-sided tests.

### 13.7.1 Multiple Comparisons and Type I Errors

Imagine there are  $I$  groups to see if all the means are the same. If you do enough comparisons, something will show as significant even if that isn’t true. (Related to publication bias)

Example:  $I = 9$

$$\binom{I}{2} = \binom{9}{2} = \frac{9!}{2! \cdot 7!} = \frac{9 \cdot 8}{2} = 36$$

Probability of rejecting a true  $H_0$  at least once:

$$\begin{aligned} P(\text{type I error}) &= P(\text{at least one rejection} | H_0 \text{ is always true}) \\ &= 1 - P(\text{no rejections} | H_0 \text{ is always true}) \\ &= 1 - (0.95^{36}) = 0.84 \end{aligned}$$

For comparing multiple groups means, use the overall  $F$ -test in ANOVA (unit 8).

Multiple comparisons occur frequently:

## Power Curves: Paired Data

- Below are the power curves for various effect sizes for:
  - 1) Paired  $t$ -test (black)
  - 2) 2-sample unpooled  $t$ -test (red)
  - 3) Signed Rank Test (green)
  - 4) Sign Test (blue)
- Assume:  $X_i \sim N(0, 1)$ ,  $Y_i \sim N(\mu_Y, 1)$ ,  $n_X = n_Y = 10$ ,  $\rho_{XY} = 0.577$

nsims = 2000

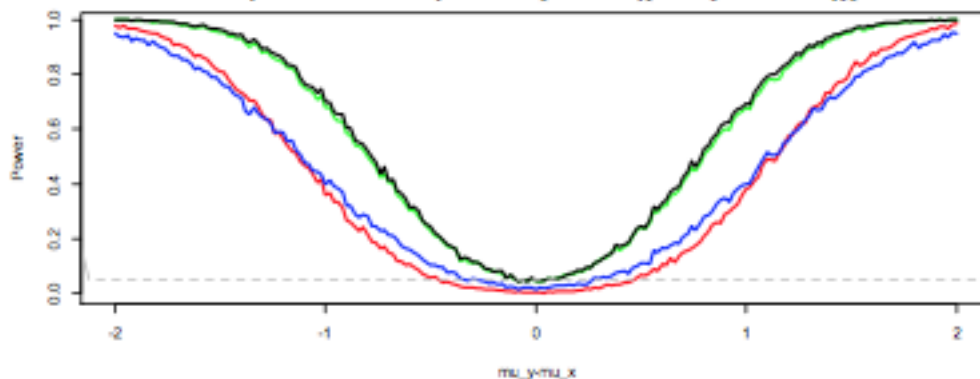


Figure 5: Sample Power Curve

1. performing many independent hypothesis tests
2. looking at many predictors in a regression model
3. multiple investigators from many universities attempting to come up with new relationships in science (aka Publication bias)

### 13.7.2 Bonferroni Correction

One solution is to adjust the  $\alpha$  level using the Bonferroni Correction. It is the most conservative approach, but often overcorrects.

Perform each individual test as  $\alpha / \text{number of test performed}$ , ie.

$$\alpha^* = \frac{\alpha}{\binom{I}{2}}$$

Likewise, each confidence interval should be widened and put  $1 - \alpha^*/2$  error in each tail.

Example: comparing three groups:

$$\binom{3}{2} = \frac{3!}{2!(1)!} = 3 \text{ possible comparisons}$$

Perform each of the three tests at  $\alpha^* = 0.05/3 = 0.0167$  or Confidence Interval at 98.33% level ( $1 - \alpha^*$ ).

## 14 Section 7

### 14.1 Rank-Sum

Basics of Rank-Sum:

**Parametric** not assuming any distribution.

**Robust**

**Independence of Observations** is the only assumption.

The  $t$ -tests are used in a lot of places where it shouldn't be. Rank-Sum could be used instead.

The null hypothesis is typically permutation / randomization:

$$H_0 : \forall_i \in [1, \dots, N], Y_i = X_i$$

The test statistic for  $X, Y$  :  $T_{stat} = \sum_{i=1}^n \text{Rank}(X_i), n_x \leq n_y$ . Sum over whichever observation has fewer samples.

Example:  $X = \{17, 65, 2, 8\}, Y = \{10, 1, 6, 99\}$ .  $\text{Rank}(X) = \{6, 7, 2, 4\}$ ,  $\text{Rank}(Y) = \{5, 3, 1, 8\}$

$$p = P(T_{obs} \leq T_{stat} | H_0)$$

$\text{Rank}(X_i) \sim \text{DUnif}(1, \dots, n_x + n_y)$ , assuming no ties.

If  $n$  is large, it is possible to do a normal approximation.

If the distributions of  $X, Y$  are similar, the null hypothesis could be that the medians (or a specific percentile) are the same.

## 14.2 Statistical Power

$P(\text{Rejection} | H_0) = \alpha = 0.05$  (want this to be small - can be set for our tests)

The lower we make  $\alpha$ , the more likely we are to reject:  $P(\text{Rejection} | H_a) = 1 - \beta$ . A smaller  $\alpha$  leads to a smaller  $\beta$ , so the probability  $(1 - \beta)$  gets larger. We might be too conservative and fail to reject too frequently.

We could select **any**  $\alpha = 5\%$  region of the CDF. It is more powerful if we reject at the extremes of the CDF. “You don’t just want to be wrong 5% of the time under the null, you also want to right the highest percentage of time under the alternative.”

Power is defined in a specific alternative hypothesis, not generally. For example  $H_a : \mu \neq \mu_0$ , would not work since we don’t actually know what  $\mu$  is.

### 14.2.1 Power Calculation Example

$H_0 : \mu = \mu_0, H_a : \mu = \mu_a$ . Power will depend specifically on what  $\mu_a$  is.

Power =  $P(R | H_a)$

We reject when we are in the outer 2.5%.

Example:

$$T = \frac{\bar{X} - \mu_0}{\sigma/\sqrt{n}} \sim N(0, 1)$$

Reject when  $|T| \geq Z_{1-\alpha/2}$

Assumptions:  $\mu_a > \mu_0, 1 - \beta \geq 0.5, 1 - \alpha/2 \geq 0.5$ .

We then won't have to the case when  $T < Z_{\alpha/2}$ . Simplifying  $T \geq Z_{1-\alpha/2}$ . So what is the probability of this is true under the alternative?

$$P\left(\frac{\bar{X} - \mu_o}{\sigma/\sqrt{n}} \geq Z_{1-\alpha/2} | H_a\right)$$

We know the distribution of this under the alternative:  $\bar{X} \sim N(\mu_a, \sigma^2)$

Under the alternative, the distribution is

$$T|H_a \sim N\left(\frac{\mu_a - \mu_0}{\sigma/\sqrt{n}}, \sigma^2\right)$$

So taking this to standardize the probability above by subtracting  $T|H_a$  from both sides of the inequality:

$$P\left(Z \geq Z_{1-\alpha/2} - \frac{\mu_a - \mu}{\sigma/\sqrt{n}}\right)$$

Let's say we want this probability to be above a threshold:

$$P\left(Z \geq Z_{1-\alpha/2} - \frac{\mu_a - \mu}{\sigma/\sqrt{n}}\right) \geq 1 - \beta$$

The probability of 'greater than' is the same as saying it is 1 - probability of 'less than' (ie.  $P(X > Y) = 1 - P(X \leq Y)$ ).

$$1 - P\left(Z < Z_{1-\alpha/2} - \frac{\mu_a - \mu}{\sigma/\sqrt{n}}\right) \geq 1 - \beta$$

$$P\left(Z < Z_{1-\alpha/2} - \frac{\mu_a - \mu}{\sigma/\sqrt{n}}\right) \leq \beta$$

meaning

$$Z_{1-\alpha/2} - \frac{\mu_a - \mu}{\sigma/\sqrt{n}} \leq Z_\beta$$



continuing the derivation - subtracting from both sides:

$$-\frac{\mu_a - \mu}{\sigma/\sqrt{n}} \leq Z_\beta - Z_{1-\alpha/2}$$

The standard normal is symmetric, ie.  $Z_\beta = -Z_{1-\beta}$

$$-\frac{\mu_a - \mu}{\sigma/\sqrt{n}} \leq -Z_{1-\beta} - Z_{1-\alpha/2}$$

Multiplying by -1:

$$\frac{\mu_a - \mu}{\sigma/\sqrt{n}} \geq Z_{1-\beta} + Z_{1-\alpha/2}$$

Rearranging terms. Also, we know  $\mu_a - \mu_0 > 0$ ,

$$\frac{\sqrt{n}(\mu_a - \mu)}{\sigma} \geq Z_{1-\beta} + Z_{1-\alpha/2}$$

$$\sqrt{n} \geq \frac{\sigma(Z_{1-\beta} + Z_{1-\alpha/2})}{\mu_a - \mu}$$

Squaring both sides (assuming RHS is positive):

$$n \geq \left( \frac{\sigma(Z_{1-\beta} + Z_{1-\alpha/2})}{\mu_a - \mu} \right)^2$$

## 15 Unit 8 - Analysis of Variance (ANOVA)

Chapter 5, Section 13.1-13.2 in the text

### 15.1 Analysis of Variance (ANOVA)

Requires lot of algebra, but analysis of variance can also be modeled by regression.

Example: inference for 3+ means, bone density. Drop rats from 1-3 ft to increase bone density. Three groups (no jump, 30 cm jump, 60 cm jump). Can't use  $t$ -test because there is no formula for 3 groups.

To know if there is a difference, find the difference in means between the groups and the individual spread in the different groups. Which leads to the analysis of variance.

Variance between groups, and variance within group is what we are going to measure.

### 15.1.1 General format and ANOVA's $F$ -test

For  $I = 3, H_0 : \mu_1 = \mu_2 = \dots = \mu_i$

The alternative is that at least one pair of means is not equal.

$$Y_{i,j} \sim N(\mu_i, \sigma^2)$$

where  $i$  is the group number ( $1, 2, \dots, I$ ), and  $j$  is the individual within group  $i$ , ( $1, 2, 3, \dots, n_i$ ). Note that there is a single variance, so this is an extension of the pooled  $t$ -test.

The total sample size:

$$\sum n_i = n$$

### 15.1.2 Main Concept of ANOVA

Within  $I$  populations / groups, there are two types of variability:

**A - variability within groups** variation of individuals around their group means

**B - variability between groups** variation of group means around the overall mean

If (A) is small relative to (B), then the group means are different.

ANOVA determines whether variability in data is mainly due to (A) or (B).

### 15.1.3 One way ANOVA

$$Y_{1,j} \sim N(\mu_1, \sigma^2), Y_{2,j} \sim N(\mu_2, \sigma^2), \dots, Y_{i,j} \sim N(\mu_i, \sigma^2),$$

where  $i = 1 : I$  groups (aka subpopulations) for a specific *factor* and  $j = 1 : J$  for individuals sampled in each group.

Can be written as:

$$Y_{ij} = \mu_i + \epsilon_{ij}$$

ie. DATA = MODEL + RESIDUAL / ERROR, where  $\epsilon$  is the unexplained residuals and assumed to be  $N(0, \sigma^2)$  (still one common variance).

The Sample means are  $\bar{Y}_i$ , the sample variances are  $S_i^2$ . The grand mean is  $\bar{Y}$ .

To check the assumption that all variances are the same, take the ratio of the largest and smallest group variances  $S^2$  and ensure it is  $\leq 2$ , ie.

$$\frac{S_{\text{largest}}^2}{S_{\text{smallest}}^2} \leq 2$$

**Variance within groups** Combine pooled estimates:

$$\begin{aligned}
 S_p^2 = S_w^2 &= \frac{(n_1 - 1)S_1^2 + \dots + (n_I - 1)S_I^2}{n - I} \\
 &= \frac{\sum_{i=1}^I (n_i - 1)S_i^2}{n - I} \\
 &= \frac{SSE_{Error}}{df_{Error}} \\
 &= MSE
 \end{aligned}$$

$S_w^2$  is the within groups estimate. A.k.a. “mean square error within groups (MSW)” and “mean square error (MSE)”

$$SSE = SS_{Within} = \sum_{i=1}^I \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_i)^2$$

which is the sample variance without dividing by  $(n_i - 1)$ . so, SSE can also be written as:

$$SSE = \sum_{i=1}^I (n_i - 1)S_i^2$$

[why? - reconcile slides]  $df = n - I$

**Variance between groups** a.k.a.  $SS_{Model} = SS_{Between}$

$$\begin{aligned}
 SSM &= \sum_{i=1}^I \sum_{j=1}^{n_i} (\bar{Y}_i - \bar{Y})^2 \\
 &= \sum_{i=1}^I n_i (\bar{Y}_i - \bar{Y})^2
 \end{aligned}$$

If the null hypothesis  $H_0 : \mu_1 = \mu_2 = \dots = \mu_I$  is then examining the individual sample means is as if we are sampling  $I$  times from the same population with mean  $\mu$  and variance  $\sigma^2$ .

Recall: the sampling distribution of sample means:  $\bar{Y} \sim N(\mu, \frac{\sigma^2}{n})$

So under  $H_0$ ,  $\sigma^2$  can be estimated by:

$$\begin{aligned}
S_B^2 &= \frac{n_1(\bar{Y}_1 - \bar{Y})^2 + \dots + n_j(\bar{Y}_j - \bar{Y})^2}{I - 1} \\
&= \frac{\sum_{i=1}^I n_i(\bar{Y}_i - \bar{Y})^2}{I - 1} \\
&= \frac{SS_{Groups}}{df_{Groups}} \\
&= MSG = MSB
\end{aligned}$$

$S_B^2$  refers to the between groups estimate.

Mean Square between groups (MSB or MSG) or Mean Square of the model (MSM).

Only a valid (bias) estimate of  $\sigma^2$  if  $H_0$  is true, otherwise it is inflated.

$$df = I - 1$$

Sums of squares / variance gives an estimate of variance.

#### 15.1.4 Concept Behind the Test

We now take these two estimates of variances and put them in a ratio.

If  $H_0 : \mu_1 = \mu_2 = \dots = \mu_I$  is true, then  $S_w^2$  and  $S_B^2$  both estimate  $\sigma^2$  and should be of similar magnitude.

If  $H_0$  is not true, then  $S_B^2$ , the between groups estimate, will in general be larger than  $S_w^2$ , the within groups estimate.

This ratio of variances is the  $F$ -test, which is test statistic in this hypothesis:

$$F = \frac{S_B^2}{S_w^2} = \frac{MSB}{MSW}$$

This test statistic has an  $F$  distribution with  $I - 1$  and  $n - I$  degrees of freedom when  $H_0$  is true.

(Recall, the  $F$  distribution is a ratio of  $\chi^2/df$  distributions. If  $X_1 \sim \chi_k^2$ ,  $X_2 \sim \chi_l^2$  and independent, then  $\frac{\chi_k^2/k}{\chi_l^2/l} \sim F_{k,l}$ , since both SSM and SSW are the sums of normalized random variables, ie. a  $\chi^2$  distribution.)

Reject  $H_0$  for large values of  $F$ . If  $F > F_{\alpha, I-1, n-I}^*$  or the corresponding  $p$ -value is the less than  $\alpha$ . See slide 17. It is a one-sided test - look up in a table (or use R) based on the degrees of freedom.

The total sum of squares  $SST = SSB + SSW$ .

$$SST = \sum_{i=1}^I \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y})^2$$

The degrees of freedom are now partitioned  $DFT(n-1) = DFB(I-1) + DFW(n-I)$ . It is  $n-1$  degrees of freedom because we had to calculate the overall mean  $\bar{Y}$ . It  $n-I$  degrees of freedom for DFW since we had to calculate the mean for each group.

The mean squares (MS) are formed the same way in every partition:

$$MS = \frac{SS}{df} = \frac{\text{Sum of squares}}{\text{degrees of freedom}}$$

TODO: Add ANOVA table to cheat sheet

### 15.1.5 Analysis of Variance Table

The is what an ANOVA table looks like in R:

Source	SS	DF	MS	F
Groups (between)	SSB	$I-1$	$S_B^2 = SSB/DFB$	$F = MSB/MSW$
Error (within)	SSW	$n-I$	$S_w^2 = SSW/DFW$	
Total	SST	$n-1$	SST/DFT	

The  $F$  statistic tests if there is a difference among any of the  $I$  population means. MSW is still our estimate of  $\sigma^2$ : the variance of the residuals or the average variance of the observations from their group means (also called the pooled variance estimate).

### 15.1.6 Critical Value $F^*$

To find the overall average:

$$\bar{Y} = \frac{n_1 \bar{Y}_1 + n_2 \bar{Y}_2 + \dots + n_I \bar{Y}_I}{n}$$

See full example Unit 8, slide 17

Find the critical value  $F^*$  using  $F$  distribution and the two degrees of freedom of the problem. In R: `qf(0.95, 2, 27)` where 0.95 is  $1 - \alpha$ , 2 = degrees of freedom for between, 27 = degrees of freedom of within.

To find the  $p$ -value, `pf(7.98, 2, 27)`, where  $F = 7.98$  and the degrees of freedom are as above.

### 15.1.7 Using ANOVA from R

Creating an ANOVA table is a two-step procedure:

```
> model1 = aov(data$bone.density~data$group)
> anova(model1)
Analysis of Variance Table

Response: data$bone.density
          Df Sum Sq Mean Sq F value    Pr(>F)
data$group 2  7433.9   3716.9    7.9778 0.001895 **
Residuals 27 12579.5    465.9
```

where AOV is the analysis of variance, which predicts a response based on many different groups. The ANOVA table summarizes the results.

To find SST, just add the sums of squares =  $7433.9 + 12579.5$

### 15.1.8 Assumptions for ANOVA $F$ -test

The  $F$ -test has the following assumptions:

1. a) Independence between groups b) Independence of observations within groups
2. Equal variances within group
3. Each observation is normally distributed around the group's mean

The distribution of SSB is:

$$SSB = \sum_{i=1}^I n_i (\bar{Y}_i - \bar{Y})^2 \sim \sigma^2 \chi_{I-1}^2$$

The distribution of SSW is:

$$SSW = \sum_{i=1}^I (n_i - 1) S_i^2 \sim \sigma^2 \chi_{n_i-I}^2 + \dots + \sigma^2 \chi_{n_j-I}^2 = \sigma^2 \chi_{N-I}^2$$

These two random variables SSB and SSW are independent from each other.

Recall that

$$F = \frac{X/df_1}{Y/df_2} \sim F_{df_1, df_2}$$

if  $X \sim \chi_{df_1}^2$  and  $Y \sim \chi_{df_2}^2$  are independent of each other.

If assumptions are true, the  $F$  stat will follow the  $F$  distribution, and the resulting  $p$ -values will be uniformly distributed. The type I error rate will be the usual  $\alpha$ .

If assumptions are violated, e.g. the correlation between groups is positive, the  $F$  statistic will be small and the  $p$ -values will be too small.

If constant variance is violated, there will be an inflated type I error rate. This could worsened by having a smaller sample size.

If the variance of one group is too small, there will be an inflated  $F$  stat.

If all are skewed, there isn't much of a problem. If one is skewed, then there will be a problem.

Summary of Assumption Violations:

- The  $F$ -test is robust (conservative at least) to positive correlation between groups (though, there may be a better approach) but not to positive correlation within groups (we did not show).
- The  $F$ -test is sensitive (not robust) to violations of constant variance, especially if sample sizes are different in the groups
- The  $F$ -test is somewhat sensitive to skewedness, especially if the groups have different skewedness and/or different sample sizes.
- If all groups have the same skew, perform a transform on all groups and then run the ANOVA test

### 15.1.9 Contrast Testing

Lecture 8

Recap:

Use ANOVA to compare a variable across multiple groups. Since there are multiple groups, we can no longer use a  $t$ -test.

Individuals  $j$  within group  $i$ ,  $Y_{ij} \sim N(\mu_i, \sigma^2)$

Three assumptions:

1. Independence within and between groups. Defined by study design.
2. Normal distribution. Look at histograms, box plots, QQplot, etc.
3. Constant variance. Look at box plots, histograms. Check that the largest to smallest variance ratio is  $< 2$ .

The ANOVA hypothesis:

$$H_0 : \mu_1 = \mu_2 = \dots = \mu_I$$
$$H_a : \text{there is one difference}$$

The test statistic is

$$F = \frac{SSM/df_m}{SSE/df_e}$$

in other words,  $F$  is the ratio of variance explained by the model divided the variance not explained by the model. If  $F$  is large, there is evidence against the null hypothesis.

Errors between groups: how far away is the variance of one group from all the other groups?

$$SSM = SSB = \sum_{i=1}^I \sum_{j=1}^{n_i} (\bar{Y}_i - \bar{Y})^2 = \sum_{i=1}^I n_i (\bar{Y}_i - \bar{Y})^2$$

with degrees of freedom  $df_M = I - 1$ .

Errors within groups: how far away is the variance within the group?

$$SSE = SSW = \sum_{i=1}^I \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_i)^2 = \sum_{i=1}^I (n_i - 1) S_i^2$$

with degrees of freedom:  $df_E = n - I$

There are two ways of determining which group differ. There are a number of techniques.

1) measure all pair-wise combinations, 2) use contrasts.

Steps in a Complete ANOVA Procedure:

1. Examine the data, checking assumptions



2. If possible, formulate in advance some working (alternative) hypotheses about how the population group means might differ. For example, if we are examining a control group vs two study groups, use

$$H_0 : \mu_1 = \frac{1}{2}(\mu_2 + \mu_3)$$

or, equivalently:

$$H_0 : \mu_1 - \frac{1}{2}(\mu_2 + \mu_3) = 0$$

This is called a **contrast**, which is a linear combination of  $\mu$ 's. The coefficients sum up to 0. If the coefficients did not sum up to 0, then the null hypothesis would not hold. If the groups had different sample sizes, you could make the coefficient a weighting factor to adjust for this.

3. Check the evidence against the global null hypothesis of no differences among the groups by calculating the F-test, or by testing all pairwise combinations.
4. If the F-test is significant (that is, if it leads to a rejection of the null hypothesis of equal population group means):
  - Test the individual hypotheses specified at the second step
  - If there was not enough information to formulate working hypotheses, test all pairwise comparisons of means, adjusting for multiple comparisons (example also coming)

After the omnibus  $F$  test has shown overall significance investigate other comparisons of groups using *contrasts* can be performed. Two ways to do this: contrast tests, and pair-wise  $t$ -tests comparing all group permutations.

A contrast is a linear combination of  $\mu_i$ 's.

$$\psi = \sum a_i(\mu_i)$$

where  $\sum a_i = 0$ . The corresponding contrast of sample means is  $c = \sum a_i(\bar{Y}_i)$

To compare group 2 and 3, ignoring the first group:

$$\psi = \mu_2 - \mu_3 = (0)\mu_1 + (1)\mu_2 + (-1)\mu_3$$

### 15.1.10 Contrast Test Hypothesis

This is similar to a pooled test (we are assuming the variance among each group is the same):

$$\begin{aligned} H_0 : \psi &= 0 \\ H_a : \psi &\neq 0 \\ T &= \frac{\sum a_i \bar{Y}_i}{S_p \sqrt{\sum \frac{a_i^2}{n_i}}} \end{aligned}$$

This  $t$ -test has a  $t$  distribution with degrees of freedom associated with  $S_p$  under  $H_0$ . The test can be 1- or 2-sided.

For this  $t$ -test, what is the variance? If the separate groups were not independent, it would be  $Var(X + Y + Z) = Var(X) + Var(Y) + Var(Z) + 2 \cdot Cov(X, Y) + 2 \cdot Cov(X, Z) + 2 \cdot Cov(Y, Z)$ . However, we are assuming that the variance of each of the groups is independent, so  $Var(X + Y + Z) = Var(X) + Var(Y) + Var(Z)$ .

Recall  $\psi = \mu_1 + \frac{1}{2}(\mu_2 + \mu_3)$ . The variance would be

$$\begin{aligned} Var &= Var(\bar{Y}_1) + Var(-\frac{1}{2}\bar{Y}_2) + Var(-\frac{1}{2}\bar{Y}_3) \\ &= \frac{\sigma^2}{n_1} + (-\frac{1}{2})^2 \frac{\sigma^2}{n_2} + (-\frac{1}{2})^2 \frac{\sigma^2}{n_3} \\ &= \sigma^2 \left( \frac{1}{n_1} + \frac{(-1/2)^2}{n_2} + \frac{(-1/2)^2}{n_3} \right) \\ sd &= \sigma \sqrt{\left( \frac{1}{n_1} + \frac{(-1/2)^2}{n_2} + \frac{(-1/2)^2}{n_3} \right)} \\ E(S_p^2) &= \sigma^2 \\ sd &= S_p \sqrt{\left( \frac{1}{n_1} + \frac{(-1/2)^2}{n_2} + \frac{(-1/2)^2}{n_3} \right)} \\ T &= \frac{\bar{Y}_1 - \frac{1}{2}(\bar{Y}_2 + \bar{Y}_3)}{S_p \sqrt{\left( \frac{1}{n_1} + \frac{(-1/2)^2}{n_2} + \frac{(-1/2)^2}{n_3} \right)}} \\ S_p^2 &= SSE/df_E \text{ (denom of } F \text{ test)} \end{aligned}$$

See worked example, unit 8, slide 35, 36.

R calls  $S_p^2$  the Residuals of Mean Squares.

### 15.1.11 Multiple Comparisons

Since contrast tests involves multiple comparisons, we have an inflated risk of type I errors. This can be corrected using the Bonferroni correction, where  $\alpha^* = \alpha/(\text{number of tests})$ .

For  $N$  simultaneous tests that are independent, the overall type I error (ie. the probability of at least one type I error among  $N$  tests) is  $1 - (1 - \alpha)^N$ .

For  $N$  simultaneous tests that are perfectly positively dependent (if one rejects, then the other combinations are likely to reject as well, even if the groups are independent), the overall type I error is just  $\alpha$ .

For  $N$  simultaneous tests that have unknown dependence, the maximum overall type I error is  $\min(N\alpha, 1)$  (since it is a probability, it can't be  $> 1$ ), where

$$\alpha \leq 1 - (1 - \alpha)^N \leq N\alpha$$

*Individual confidence level* is a probability that a *single* confidence interval covers the true value.

*Overall (familywise) confidence level* is a probability that *all* confidence cover the corresponding true values at the same time. Perform the confidence levels at the adjusted rate. ( $\alpha$  and confidence levels add up to 1).

Analogously, if the success rate of a  $(1-\alpha)$  100% confidence interval is  $(1-\alpha)$ , the simultaneous success rate of several  $(1-\alpha)$ 100% confidence intervals is **less than  $(1-\alpha)$** .

Methods for Multiple Comparison:

Differ by types of multipliers for CIs or modifications to reference distribution:

$$\text{Estimate} \pm \text{multiplier}SE$$

The *multiplier* comes from a  $t$  distribution such as Bonferroni, or Tukey's HSD. [We do not cover Fisher's Protected least significant difference and Scheffe's procedure]

Bonferroni is, as before, most conservative. For pairwise mean comparisons: Margin of Error for  $(1-\alpha)$  100% CI:  $t_{n-I, 1-\alpha/(2N)} \cdot SE$ .

### 15.1.12 Tukey Honest Significant Difference (HSD)

This is the correct adjustment under ANOVA. Basic idea: consider the *largest difference* between any two sample means for for  $I$  groups.

Intuition: There are  $\bar{Y}_1, \bar{Y}_2, \dots, \bar{Y}_I$  sample means. To do a pairwise comparison would generate  $\binom{I}{2}$  2 sample  $t$ -tests. We are going to reject one of these with 5% comparison rates. So it reduces to only have to check the ratio of max and min sample means.

It is appropriate when interested in differences between all pairs of group means. We assume normality, equal variances, equal sample sizes( $\bar{n}$ ), and equal means.

$$Q = \frac{\bar{Y}_{max} - \bar{Y}_{min}}{S_p \sqrt{\bar{n}}} \sim q(1 - \alpha, n - I), \text{ where } n = \bar{n}I$$

The Tukey distribution  $q$  is called the **Studentized Range Distribution**.

In R, use `qtukey()` to obtain quantiles.

The margin of for the Confidence Interval:

$$\left( \frac{q(1 - \alpha, I, n - I)}{\sqrt{2}} \right) SE, \text{ where } SE = S_p \sqrt{\frac{1}{\bar{n}} + \frac{1}{\bar{n}}}$$

This is also used for the critical value as part of the  $t$ -test.

In R,

```
model <- aov(...)
TukeyHSD(0.95, 2, 1000)/sqrt(2)
qtukey(0.95, 3, 1000)/sqrt(2)
pt(qtukey(0.95, 3, 1000)/sqrt(2), df=1000) # one tail
1 - pt(qtukey(0.95, 3, 1000)/sqrt(2), df=1000) # other tail
2 * (1 - pt(qtukey(0.95, 3, 1000)/sqrt(2), df=1000)) # two tail
```

### 15.1.13 Two-way and multi-way ANOVA

One way ANOVA - only one factor affecting the groups.

Two way ANOVA - multiple grouping variables to compare groups.

In R:

```

model1 = aov(logtexts ~ female) # One way ANOVA
model3 = aov(logtexts ~ classyear + female) # Two way ANOVA
model4 = aov(logtexts ~ classyear * female) # Interaction variable

```

This can be reduced to boxplot

**Interactive effect:** when one variable changes the other variable. E.g. by going from sophomore to junior, the number of texts would change.

degrees of freedom =  $(I_1 - 1) * (I_2 - 1) = (4 - 1) * (2 - 1) = 3$

## 15.2 Kruskal-Wallis Test

This is a form of ANOVA on Ranks rather than  $t$ -test. The KW test is just an extension of the Wilcoxon Rank Sum test to 3 or more groups.

The same procedure as other rank tests:

1. Rank all the combined data ignoring groups from 1 to  $N$  (treating them all like one sample). For any ties, average those ranks.
2. Calculate an  $F$  like  $\chi^2$  test statistic:

$$K = (N - 1) \frac{\sum_{i=1}^I n_i (\bar{R}_i - \bar{R})^2}{\sum_{i=1}^I \sum_{j=1}^{n_i} (R_{ij} - \bar{R})^2} \sim \chi_{I-1}^2$$

The hypothesis is like the Rank-Sum test: the quantiles in the group are all the same.

It is  $\chi^2$  because the numerator and denominator are tightly related. Only the numerator is a random variable.

In R:

```

ranks = rank(logtexts)
ni=as.vector(table(classyear))
ribar = as.vector(by(ranks, classyear, mean))
rbar = mean(ranks)
N = length(ranks)
I = length(ni)
K = (N-1) * sum(ni * (ribar - rbar)^2) / sum((ranks-rbar)^2)
K
1-pchisq(K,df=I-1)
# Or....
kruskal.test()

```

## 16 Section 8 - ANOVA

One way ANOVA hypotheses  $H_0 : \mu_1 = \mu_2 \dots = \mu_I$ , then the alternative is one is not equal  $H_a$  : not all  $\mu_i$  are equal.

Is the variance within the groups less than the variance between groups.

Assumptions:

1. Independence between and within groups. Use plots, by time and / or space. Can also judge by study design, ie. data comes from the same family would not be independent.
2. Equal variances within and between groups. Example, one of the groups might be bimodal.

HW7 Q1:

$$\begin{aligned}\sum e_{ij} &= 0 \\ \sum (Y_{ij} - \bar{Y}_i) &= 0 \\ \sum (Y_{ij} - \frac{\sum Y_{ij}}{n}) &= 0 \\ n\bar{Y} - n\bar{Y} &= 0\end{aligned}$$

$$\sum_{i=1}^n Y_i = n \cdot \bar{Y}$$

## 17 Section 9 - Contrast

$$3. \text{ a } P(\text{at least one type I error}) = 1 - (1 - \alpha)^I = 1 - 0.95^I$$

## 18 Unit 9 - Simple Linear Regression

Chapters 7 and 8 in the text.

Simple definition: fitting a line to a scattering of points.

Simple linear regression: one predictor variable.

The trash model: ignore  $x$  to predict  $y$ .

$$\mu_{y|x} = \beta_0 + \beta_1 x$$

## 19 Section 9 - Exam Prep

1. transformations on a distribution
2. ANOVA  $\alpha^*$  confidence interval
3. Items for the cheat sheet
4. Spring Exam

not the same means, so more likely to reject so the histogram of the p-values will heavily right tailed.

go through assumptions being violated. Any time the sample size is different, the distortions in the p-values is amplified.

be able to generate some sample distributions using the TI nspire