

# **MATH 326**

Probability & Statistics II

Tom Carty – Spring 2022

TeX'ed by Matthew Wilder (BU '23)

# Contents

Abbreviations and Notations . . . . .	3
<b>7 Sampling Distributions and the Central Limit Theorem</b>	<b>4</b>
7.2 Sample Means . . . . .	4
Theorem 7.1 . . . . .	4
Theorem 7.2 . . . . .	5
Theorem 7.3 (Fisher's Theorem) . . . . .	5
7.4 The Central Limit Theorem . . . . .	6
Table 4: Normal Curve Areas . . . . .	8
<b>8 Estimation</b>	<b>9</b>
8.1 An Estimator . . . . .	9
8.2 The Bias and Mean Square Error of Point Estimators . . . . .	10
8.3 Some Common Unbiased Point Estimators . . . . .	14
8.4 Evaluating the Goodness of a Point Estimator . . . . .	17
8.5 Confidence Intervals . . . . .	18
8.6 Large Sample Confidence Intervals . . . . .	21
8.7 Selecting the Sample Size . . . . .	23
8.8 Small-Sample Confidence Intervals for $\mu$ and $\mu_1 - \mu_2$ . . . . .	25
Table 5: <i>t</i> -Distributions . . . . .	26
8.9 Confidence Intervals for $\sigma^2$ . . . . .	30
Table 6: $\chi^2$ Lower Tail . . . . .	31
Table 6: $\chi^2$ Upper Tail . . . . .	32
<b>9 Properties of Point Estimators and Methods of Estimation</b>	<b>34</b>
9.2 Relative Efficiency . . . . .	34
9.3 Consistency . . . . .	37
Consistency Theorem . . . . .	38
Probabalistic Convergence Limit Laws . . . . .	40
9.4 Sufficiency . . . . .	42
The Factorization Theorem . . . . .	44
9.5 Rao-Blackwell Theorem and Minimum-Variance Unbiased Estimation . . . . .	48
9.6 The Method of Moments . . . . .	52
9.7 The Method of Maximum Likelihood . . . . .	55
<b>10 Hypothesis Testing</b>	<b>62</b>
10.3 Z-tests (large samples) . . . . .	65
10.4 More about errors and sample size . . . . .	68
10.8 T tests . . . . .	71

<b>11 Linear Models and Estimation by Least Squares</b>	<b>77</b>
11.1 Introduction . . . . .	77
11.2 Linear Statistical Models . . . . .	77
11.3 The Method of Least Squares . . . . .	78
11.4 Properties of the Least-Squares Estimators: Simple Linear Regression . . . . .	82
11.5 Inferences concerning the point estimators . . . . .	87
11.6 Predictions via least squares regression line . . . . .	89
11.7 Predictions on $y$ . . . . .	91
11.10 Multiple Linear Regression . . . . .	93
11.11 A big ol theorem . . . . .	97
11.12 Hypothesis Testing C.I. . . . .	98
<b>13 One-way Analysis of Variance</b>	<b>102</b>

## Common Abbreviations and Notations

dist'n = Distribution  
 r.v. = Random variable  
 DOF = degrees of freedom  
 d.f. = degrees of freedom  
 MGF = moment generating function  
 mgf = moment generating function  
 pdf = probability distribution function  
 CLT = Central Limit Theorem  
 iid = Independent Identically Distributed  
 MSE = Mean Square Error  
 MoM = Method of Moments  
 MVUE = Minimum Variance Unbiased Estimator  
 MLE = Maximum Likelihood Estimator  
 $\mu$  = mean (average)  
 $\sigma$  = standard deviation  
 $\sigma^2$  = variance  
 $\bar{X}$  = Sample mean  
 $S$  = Sample standard deviation  
 $S^2$  = Sample variance  
 $N(\mu, \sigma^2)$  = Normal distribution with: mean =  $\mu$ , variance =  $\sigma^2$   
 $Z$  = Z-score

# Chapter 7

## Sampling Distributions and the Central Limit Theorem

### Chapter 7 Review

Our first goal in Math 326 is to learn how to estimate global parameters of “population” like  $\mu$  and  $\sigma$ . To do this, we need to understand how the random variable we use to estimate them are distributed. For example:

$$\text{Parameters } \begin{cases} \mu \\ \sigma^2 \end{cases} \quad \text{Statistics } \begin{cases} \bar{X} = \frac{1}{n} \sum X_i \\ S^2 = \frac{1}{n-1} \sum (X_i - \bar{X})^2 \end{cases}$$

### 7.2 Sample Means

#### ① Sample means:

**Theorem 7.1:** Let  $Y_1, \dots, Y_n$  be a random sample from  $N(\mu, \sigma^2)$  then  $\bar{Y}$  is distributed by  $N\left(\mu, \frac{\sigma^2}{n}\right)$ . The distribution of sample means,  $\bar{Y}$ , is also Normal.

**Reason:** Linearity of Expectation and Variance properties

**Discussion:** Recall in working with  $N(\mu, \sigma^2)$ , we learned it was easier to standardize everything via  $Z$ -scores:

$$Z = \frac{x - \mu}{\sigma}$$

and  $Z$  is distributed  $N(0, 1)$ , the **Standard Normal Distribution**.

#### ② Sample variance:

Recall standard deviation  $\sigma$  is a measure of the spread of the random variable and it's derived from  $(Y_i - \bar{Y})^2$  terms. In all math, we normalize to take the “units” out of things.

$$U_i = \underbrace{\frac{Y_i - \bar{Y}}{S}}_{\text{Data Driven = Stat}} \approx \underbrace{\frac{Y_i - \mu}{\sigma}}_{\text{Not a stat}} = Z_i$$

**Reason:**  $Z_i$  depends on unknown population parameters, nonetheless, it makes it easier to pretend that we start here.

**Theorem 7.2** If  $Y_1, \dots, Y_n$  are a random sample of  $N(\mu, \sigma^2)$ . Then,

$$U = \sum Z_i^2 = \sum \left( \frac{Y_i - \mu}{\sigma} \right)^2$$

has a  $\chi^2$  distribution with  $n$  degrees of freedom (df).

Recall:  $\chi^2$  distribution is a Gamma  $\left(\frac{\nu}{2}, 2\right)$  where  $\nu = \text{df}$  (degrees of freedom).

**Reason:** In old homework (325),  $Z_i^2$  is  $\chi^2$  with  $\text{df} = 1$ . By product of mgf,  $\sum Z_i^2$  is  $\chi^2$  with  $\text{df} = n$ . Now to get sample variance, we do some algebra:

$$\begin{aligned} S^2 &= \frac{1}{n-1} \sum (Y_i - \bar{Y})^2 \\ &\approx \frac{1}{n-1} \sum (Y_i - \mu)^2 \\ \iff &\frac{S^2}{\sigma^2} \approx \frac{1}{n-1} \sum \left( \frac{Y_i - \mu}{\sigma} \right)^2 \\ \iff &\frac{S^2(n-1)}{\sigma^2} \approx \sum Z_i^2 \end{aligned}$$

Showing it is okay to replace  $\bar{Y}$  with  $\mu$  is the point of the proof of the following theorem:

**Theorem 7.3 (Fisher's Theorem)** The distribution of sample variance  $S^2$

If  $Y_1, \dots, Y_n$  is a random sample from  $N(\mu, \sigma^2)$ . Then,

- ①  $\frac{S^2(n-1)}{\sigma^2}$  has  $\chi^2$  distribution with  $(n-1)$  degrees of freedom.
- ②  $\bar{Y}$  and  $S^2$  are independent random variables.
- ③  $t$ -distribution and  $F$ -distribution

$$\underbrace{\frac{Y_i - \mu}{\sigma}}_{\text{Normal}} \approx \underbrace{\frac{Y_i - \mu}{S}}_{\text{t-dist'n}} \approx \underbrace{\frac{Y_i - \bar{Y}}{S}}_{\text{Statistic}}$$

Address when to use what is the point of this class

Skip definitions (such as moments) for the time being

Also, The Law of Large Numbers. Used to prove the Central Limit Theorem.

## 7.4 The Central Limit Theorem

The reason Normal distributions play an outsized role in applied statistics is that the distribution of  $\bar{Y}$  can be made *nearly normal*, no matter the underlying distribution of  $Y$  (Does not need to start “life” Normal). So *nearly* normal, that we just pretend it is.

### Theorem 7.4: The Central Limit Theorem (CLT)

Let  $Y_1, \dots, Y_n$  be independent and identically distributed (iid) random variables with

$$E(Y_i) = \mu \quad \text{and} \quad \text{Var}(Y_i) = \sigma^2 < \infty$$

(Need finite variance for CLT to hold)

Define

$$U_n = \frac{\bar{Y} - \mu}{\sqrt{\frac{\sigma^2}{n}}} = \frac{\sum(Y_i) - n\mu}{\sigma\sqrt{n}}.$$

The CLT says the distribution function of  $U_n$  converges to  $N(0, 1)$  as  $n \rightarrow \infty$ .

Big idea: The distribution  $\bar{Y}$  can be thought of as  $N\left(\mu, \frac{\sigma^2}{n}\right)$ .

**Definition:** The support of  $f$  is the domain on which  $f$  is non-zero.

**Example:** Let  $\bar{X}$  denote the mean of a random sample of size  $n = 15$ , from the distribution whose pdf is

$$f(x) = \frac{3}{2}x^2, \quad x \in [-1, 1]$$

Can be shown that

$$\mu = E[X] = 0 \quad \text{and} \quad \sigma^2 = E[(X - \mu)^2] = \frac{3}{5}$$

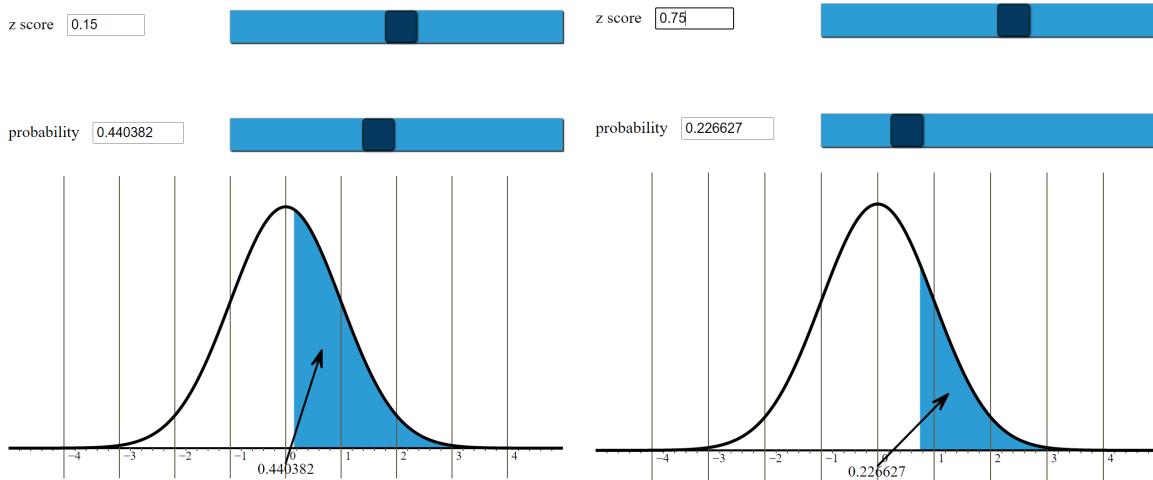
To compute  $\Pr(0.03 \leq \bar{X} \leq 0.15)$ , we use the CLT and assume  $\bar{X}$  is distributed by

$$N\left(\mu, \frac{\sigma^2}{n}\right) = N\left(0, \frac{3/5}{15}\right) = N\left(0, \frac{1}{25}\right).$$

Then

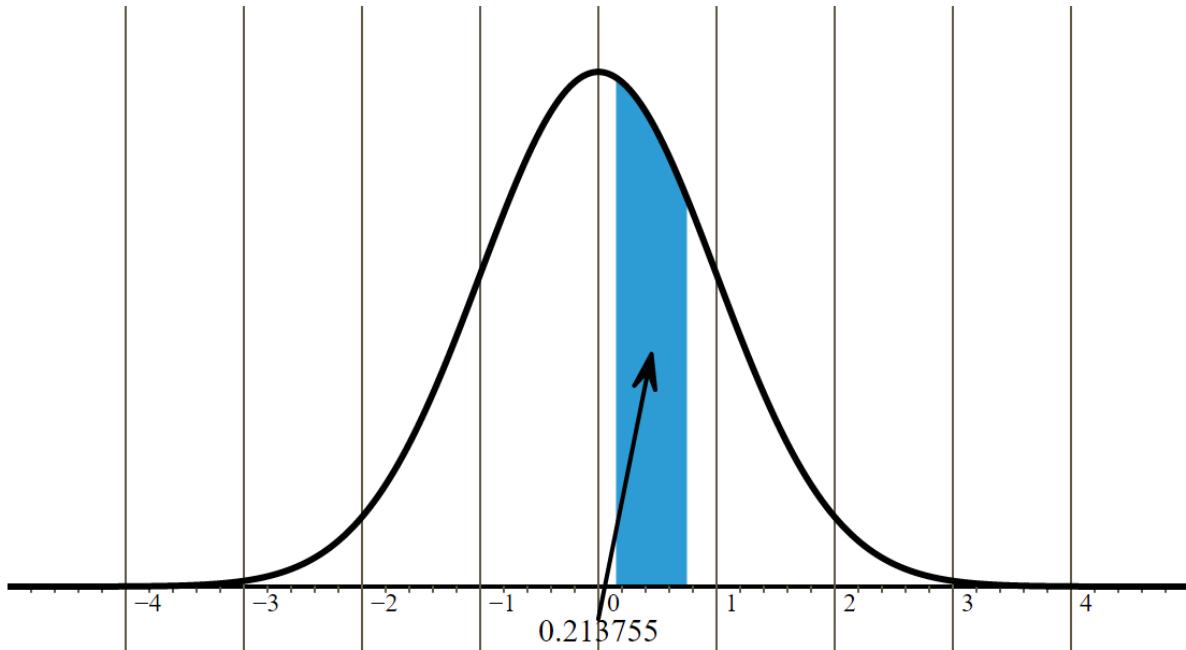
$$\begin{aligned} \Pr(0.03 \leq \bar{X} \leq 0.15) &= \Pr\left(\frac{0.03 - \mu}{\sigma} \leq \frac{\bar{X} - \mu}{\sigma} \leq \frac{0.15 - \mu}{\sigma}\right) \\ &= \Pr\left(\frac{0.03 - 0}{\sqrt{1/25}} \leq Z \leq \frac{0.15 - 0}{\sqrt{1/25}}\right) \\ &= \Pr\left(\frac{0.03}{\sqrt{1/25}} \leq Z \leq \frac{0.15}{\sqrt{1/25}}\right) \\ &= \Pr(0.15 \leq Z \leq 0.75) \\ &= \text{Table4}(0.15) - \text{Table4}(0.75) \\ &= 0.4404 - 0.2266 \\ &= 0.2138 \end{aligned}$$

Recall: Table 4 gives upper tail probabilities.



Therefore,

$$\begin{aligned} \text{Table4}(0.15) - \text{Table4}(0.75) &= \Pr(Z > 0.15) - \Pr(Z > 0.75) \\ &= 0.4404 - 0.2266 \\ &= 0.2138 \end{aligned}$$



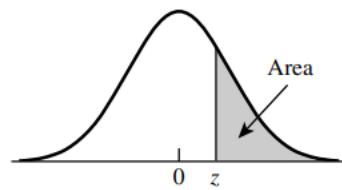
**Discussion:** About  $n$ . Surprisingly, “large”  $n$  doesn’t need to be that large. Usually for any random variable  $X$  and any distribution, then

$$n \geq 30 \quad \text{is enough.}$$

That is, when  $n \geq 30$ , the CLT says  $N\left(\mu, \frac{\sigma^2}{n}\right)$  yields a good approximation of  $\bar{X}$ .

When  $X$  is symmetric, unimodal, and continuous,  $n = 4$  or  $n = 5$  is often enough.

**Table 4 Normal Curve Areas**  
**Standard normal probability in right-hand tail**  
**(for negative values of  $z$ , areas are found by symmetry)**



$z$	Second decimal place of $z$									
	.00	.01	.02	.03	.04	.05	.06	.07	.08	.09
0.0	.5000	.4960	.4920	.4880	.4840	.4801	.4761	.4721	.4681	.4641
0.1	.4602	.4562	.4522	.4483	.4443	.4404	.4364	.4325	.4286	.4247
0.2	.4207	.4168	.4129	.4090	.4052	.4013	.3974	.3936	.3897	.3859
0.3	.3821	.3783	.3745	.3707	.3669	.3632	.3594	.3557	.3520	.3483
0.4	.3446	.3409	.3372	.3336	.3300	.3264	.3228	.3192	.3156	.3121
0.5	.3085	.3050	.3015	.2981	.2946	.2912	.2877	.2843	.2810	.2776
0.6	.2743	.2709	.2676	.2643	.2611	.2578	.2546	.2514	.2483	.2451
0.7	.2420	.2389	.2358	.2327	.2296	.2266	.2236	.2206	.2177	.2148
0.8	.2119	.2090	.2061	.2033	.2005	.1977	.1949	.1922	.1894	.1867
0.9	.1841	.1814	.1788	.1762	.1736	.1711	.1685	.1660	.1635	.1611
1.0	.1587	.1562	.1539	.1515	.1492	.1469	.1446	.1423	.1401	.1379
1.1	.1357	.1335	.1314	.1292	.1271	.1251	.1230	.1210	.1190	.1170
1.2	.1151	.1131	.1112	.1093	.1075	.1056	.1038	.1020	.1003	.0985
1.3	.0968	.0951	.0934	.0918	.0901	.0885	.0869	.0853	.0838	.0823
1.4	.0808	.0793	.0778	.0764	.0749	.0735	.0722	.0708	.0694	.0681
1.5	.0668	.0655	.0643	.0630	.0618	.0606	.0594	.0582	.0571	.0559
1.6	.0548	.0537	.0526	.0516	.0505	.0495	.0485	.0475	.0465	.0455
1.7	.0446	.0436	.0427	.0418	.0409	.0401	.0392	.0384	.0375	.0367
1.8	.0359	.0352	.0344	.0336	.0329	.0322	.0314	.0307	.0301	.0294
1.9	.0287	.0281	.0274	.0268	.0262	.0256	.0250	.0244	.0239	.0233
2.0	.0228	.0222	.0217	.0212	.0207	.0202	.0197	.0192	.0188	.0183
2.1	.0179	.0174	.0170	.0166	.0162	.0158	.0154	.0150	.0146	.0143
2.2	.0139	.0136	.0132	.0129	.0125	.0122	.0119	.0116	.0113	.0110
2.3	.0107	.0104	.0102	.0099	.0096	.0094	.0091	.0089	.0087	.0084
2.4	.0082	.0080	.0078	.0075	.0073	.0071	.0069	.0068	.0066	.0064
2.5	.0062	.0060	.0059	.0057	.0055	.0054	.0052	.0051	.0049	.0048
2.6	.0047	.0045	.0044	.0043	.0041	.0040	.0039	.0038	.0037	.0036
2.7	.0035	.0034	.0033	.0032	.0031	.0030	.0029	.0028	.0027	.0026
2.8	.0026	.0025	.0024	.0023	.0023	.0022	.0021	.0021	.0020	.0019
2.9	.0019	.0018	.0017	.0017	.0016	.0016	.0015	.0015	.0014	.0014
3.0	.00135									
3.5	.000 233									
4.0	.000 031 7									
4.5	.000 003 40									
5.0	.000 000 287									

# Chapter 8

## Estimation

### 8.1 An Estimator

We are seeking an unknown population parameter. General name of  $\theta$  (e.g.  $\mu, \sigma^2, \rho, \dots$ ).

The rule used to approximate or guess  $\theta$  is called an estimator. Any guess based on observations is called an estimate.

**Example:** Want  $\theta = \mu$ .

Estimator:  $\bar{X} = \underbrace{\frac{1}{n} \sum x_i}_{\text{a function}}$

Estimate: Given 7 observations,  $Y_1, \dots, Y_7$ , then  $\bar{Y} = \frac{1}{7} \sum_1^7 Y_i$  is an estimate.

This is an example of a point estimator.

These are also interval estimators.

**Example:** Confidence interval,  $\theta = \mu$ .

$$\bar{X} - SZ_x \leq \mu \leq \bar{X} + SZ_x \quad \iff \quad |\bar{X} - \mu| \leq SZ_x$$

## 8.2 The Bias and Mean Square Error of Point Estimators

A start at trying to determine if a point estimator is any good.

Let  $\hat{\theta}$  be a point estimator for a parameter  $\theta$ . (e.g.  $\theta = \mu$ ,  $\hat{\theta} = \bar{X}$ )

- ① One goal of a good estimator is that  $E(\hat{\theta}) = \theta$  (e.g. by CLT,  $E(\bar{X}) = \mu$ )  
But this is not always the case.

**Definition:** If  $E[\hat{\theta}] = \theta$ , we say that  $\hat{\theta}$  is an unbiased point estimator.  
(e.g.  $\bar{X}$  is unbiased).

If  $\hat{\theta}$  is a biased point estimator, we define the **bias** to be

$$\text{bias}(\hat{\theta}) = E[\hat{\theta}] - \theta$$

- ② Another goal of a good estimator might be that its “spread” of observations is tightly packed (talking about variance), and hopefully near  $\theta$ .

spread  $\implies$  variance

**Definition:** The mean square error of  $\hat{\theta}$  is

$$MSE(\hat{\theta}) = E[(\hat{\theta} - \theta)^2]$$

Recall:  $V[X] = E[X^2] - (E[X])^2$

**Corollary:**

$$MSE(\hat{\theta}) = \text{Var}(\hat{\theta}) + [\text{bias}(\hat{\theta})]^2$$

*Proof.*

$$\begin{aligned} MSE(\hat{\theta}) &= E[\hat{\theta}^2 - 2\hat{\theta}\theta + \theta^2] \\ &= E[\hat{\theta}^2] - 2E[\hat{\theta}\theta] + E[\theta^2] \\ &= E[\hat{\theta}^2] - 2E[\hat{\theta}]\theta + \theta^2 \\ &= E[\hat{\theta}^2] - (E[\hat{\theta}])^2 + (E[\hat{\theta}])^2 - 2E[\hat{\theta}]\theta + \theta^2 \\ &= \text{Var}(\hat{\theta}) + (E[\hat{\theta}] - \theta)^2 \\ &= \text{Var}(\hat{\theta}) + (\text{bias}(\hat{\theta}))^2 \end{aligned}$$

□

*Remark:* We have decomposed MSE by:

$$MSE(\hat{\theta}) = \underbrace{\text{Var}(\hat{\theta})}_{\text{Precision}} + \underbrace{\text{bias}(\hat{\theta})^2}_{\text{Accuracy}}$$

Idea of a “best” estimator is tricky. We’d like MSE to be as small as possible, but this is an impossible problem because MSE depends on the unknown  $\hat{\theta}$ .

**Example:** (Population proportions). Want  $p$ .

To estimate, let

$$Y = \begin{cases} 1 & \text{“success”} \\ 0 & \text{“failure”} \end{cases}$$

and  $\{Y_i\}$  be an iid random sample.

**Estimator 1:**  $\hat{p} = \frac{1}{n} \sum Y_i = \bar{Y}$

This sums the  $\frac{\text{people said yes}}{\text{total people}}$

Recall:  $\sum Y_i$  is  $\text{binom}(n, p)$ .

And for a Binomial Distribution, we know

$$\mathbb{E}[X] = np \quad \text{and} \quad \text{Var}(X) = npq = np(1-p)$$

Bias? Then what is the  $\mathbb{E}[\hat{p}]$ ?

$$\mathbb{E}[\hat{p}] = \mathbb{E}\left(\frac{1}{n} \sum Y_i\right) = \frac{1}{n} \mathbb{E}\left(\sum Y_i\right) = \frac{1}{n} \cdot np = p$$

Because  $\mathbb{E}[\hat{p}] = p$ , then  $\hat{p}$  is an unbiased estimator. For  $\text{MSE}(\hat{p})$ , we need  $\text{Var}[\hat{p}]$ :

$$\text{Var}[\hat{p}] = \text{Var}\left(\frac{1}{n} \sum Y_i\right) = \left(\frac{1}{n}\right)^2 \text{Var}\left(\sum Y_i\right) = \left(\frac{1}{n}\right)^2 np(1-p) = \frac{p(1-p)}{n}$$

Therefore,

$$\begin{aligned} \text{MSE}(\hat{p}) &= \text{Var}(\hat{p}) - \text{bias}(\hat{p})^2 \\ &= \frac{p(1-p)}{n} + \underbrace{0^2}_{\text{“unbiased”}} \end{aligned}$$

**Estimator 2:**  $\tilde{p} = \frac{\sum_1^n Y_i + 1}{n + 2}$

$$n = 1, \quad y = \begin{cases} 1 \\ 0 \end{cases} \quad \tilde{p} = \begin{cases} 2/3 \\ 1/3 \end{cases}$$

$$n = 2, \quad y = \begin{cases} 2 \\ 1 \\ 0 \end{cases} \quad \tilde{p} = \begin{cases} 3/4 \\ 2/4 \\ 1/4 \end{cases}$$

$$\mathbb{E}[\tilde{p}] = \frac{\mathbb{E}(Y_i) + \mathbb{E}(1)}{n+2} = \frac{np+1}{n+2} = p + \frac{1-2p}{n+2} \neq p$$

Hence,  $\tilde{p}$  is biased. And,

$$\text{bias}(\tilde{p}) = \text{E}[\tilde{p}] - p = p + \frac{1-2p}{n+2} - p = \frac{1-2p}{n+2}$$

**Calc question:** What happens to  $\text{E}[\tilde{p}]$  as sample size grows? ( $\text{E}[\tilde{p}] \rightarrow p$ )

Need

$$\begin{aligned}\text{Var}(\tilde{p}) &= \text{Var}\left(\frac{\sum_1^n Y_i + 1}{n+2}\right) \\ &= \text{Var}\left(\frac{\sum_1^n Y_i}{n+2} + \frac{1}{n+2}\right) \\ &= \text{Var}\left(\frac{\sum_1^n Y_i}{n+2}\right) \\ &= \frac{1}{(n+2)^2} \text{Var}\left(\sum_1^n Y_i\right) \\ &= \frac{np(1-p)}{(n+2)^2}\end{aligned}$$

Then

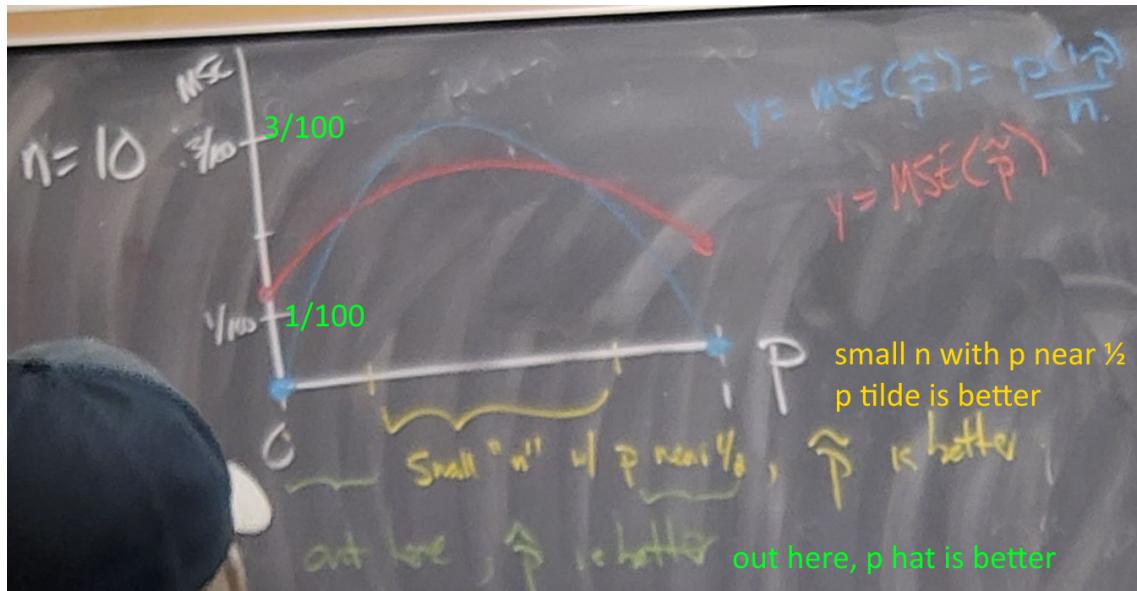
$$\begin{aligned}\text{MSE}(\tilde{p}) &= \text{Var}(\tilde{p}) + \text{bias}(\tilde{p})^2 \\ &= \frac{np(1-p)}{(n+2)^2} + \left(\frac{1-2p}{n+2}\right)^2 \\ &= \frac{np(1-p) + (1-2p)^2}{(n+2)^2}\end{aligned}$$

**Q:** Which estimator is better?      **A:** *it depends!*

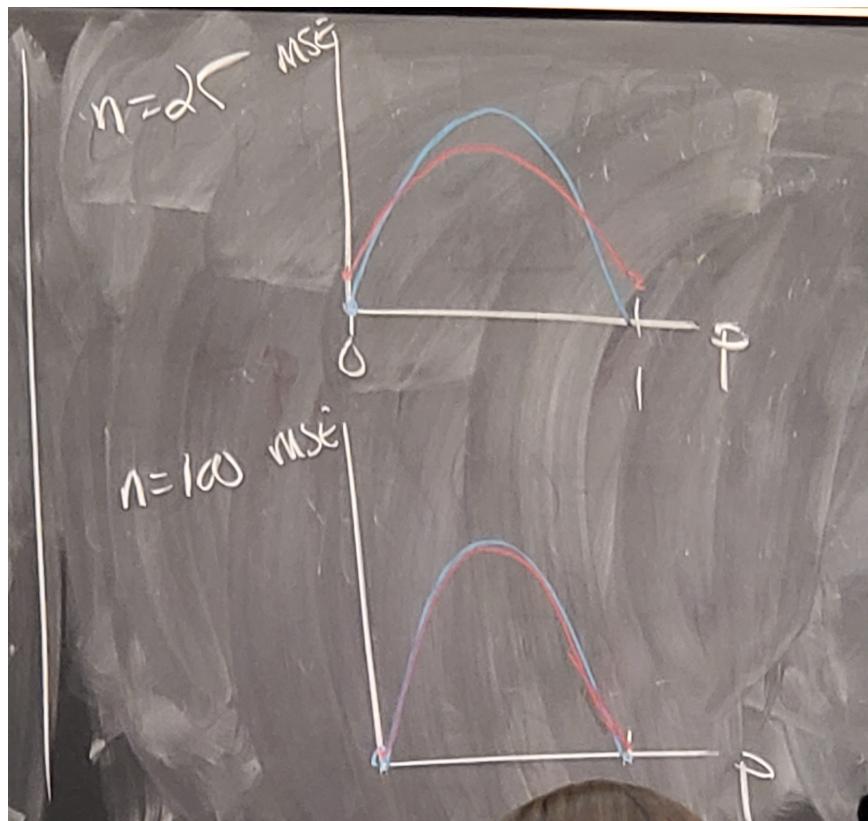
Note that for both,  $\text{MSE}(\tilde{p})(n, p)$  are functions of  $n$  and  $p$ . Let

$$y = \text{MSE}(\tilde{p})(n, p) \quad \text{and} \quad \textcolor{blue}{y} = \text{MSE}(\hat{p})(n, p)$$

fix  $n = 10$ ,  $p \in [0, 1]$



$n = 25$  and  $100$   $p \in [0, 1]$



For larger  $n$ , they become indistinguishable.

### 8.3 Some Common Unbiased Point Estimators

Table 8.1, page 397: Common unbiased estimator for  $\mu$ ,  $\rho$ ,  $\mu_1 - \mu_2$ ,  $\rho_1 - \rho_0$ .

**Example:** Variance of a data set and the  $n - 1$ .

Natural definition is  $S^2 = \frac{1}{n} \sum (X_i - \bar{X})^2$

Why do we use  $n - 1$ ?

Check bias ( $S^2$ ). Is  $E(S^2) = \sigma^2$ ? Let

$$E(X) = \mu \quad \text{and} \quad \text{Var}(X) = \sigma^2$$

and  $X_1, \dots, X_n$  iid random sample. Then,

$$\begin{aligned} E(S^2) &= E\left[\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2\right] \\ &= \frac{1}{n} \underbrace{\left[\sum (X_i^2 - 2X_i\bar{X} + \bar{X}^2)\right]}_{\text{focus on this}} \\ &= \frac{1}{n} \left(E\left[\sum X_i^2 - 2 \sum X_i\bar{X} + \sum \bar{X}^2\right]\right) \\ &= \frac{1}{n} \left(\underbrace{E\left[\sum X_i^2\right]}_{\textcircled{1}} - \underbrace{2E\left[\sum X_i\bar{X}\right]}_{\textcircled{2}} + \underbrace{E\left[\sum \bar{X}\right]^2}_{\textcircled{3}}\right) \end{aligned}$$

① and ③ uses the same “trick” of  $\text{Var}(X) = E[X^2] - (E[X])^2$ . So,

$$\begin{aligned} \text{Part } \textcircled{1} &= E\left[\sum_{i=1}^n X_i^2\right] \\ &= \sum E[X_i^2] \\ &= \sum (\text{Var}(X_i) + (E[X_i])^2) \\ &= \sum (\sigma^2 + \mu^2) \\ &= n(\sigma^2 + \mu^2) \end{aligned}$$

$$\begin{aligned} \text{Part } \textcircled{3} &= E\left[\sum_{i=1}^n \bar{X}^2\right] \\ &= \sum (\text{Var}(\bar{X}) + (E[\bar{X}])^2) \\ &= \sum \left(\frac{\sigma^2}{n} + \mu^2\right) \\ &= n \left(\frac{\sigma^2}{n} + \mu^2\right) \\ &= \sigma^2 + n\mu^2 \end{aligned}$$

$$\begin{aligned}
\text{Part ②} &= -2 \mathbb{E} \left[ \sum_i X_i \left( \frac{1}{n} \sum_j X_j \right) \right] \\
&= -\frac{2}{n} \mathbb{E} \left[ \sum_i \sum_j X_i X_j \right] \\
&= -\frac{2}{n} \mathbb{E} \left[ \underbrace{\sum_i X_i X_j}_{n \text{ terms}} + \underbrace{\sum_{i \neq j} X_i X_j}_{n^2 - n \text{ terms}} \right] \\
&= -\frac{2}{n} \left( \underbrace{\sum_i \mathbb{E}[X_i^2]}_{\text{this is part ① again}} + \underbrace{\sum_{i \neq j} \mathbb{E}[X_i X_j]}_{\text{covariance like}} \right)
\end{aligned}$$

Recall:

$$\begin{aligned}
\text{Cov}(X_i, X_j) &= \mathbb{E} [(X_i - \bar{X})(X_j - \bar{X})] \\
&= \mathbb{E}[X_i X_j] - \mathbb{E}[X_i] \mathbb{E}[X_j]
\end{aligned}$$

But iid, “i” for independent  $\implies \text{Cov}(X_i, X_j) = 0$ . Therefore,

$$\mathbb{E}[X_i X_j] = \mathbb{E}[X_i] \mathbb{E}[X_j] = \mu \cdot \mu = \mu^2$$

$$\begin{aligned}
\text{Part ②} &= -\frac{2}{n} \left( n(\sigma^2 + \mu^2) + (n^2 - n)\mu^2 \right) \\
&= -2 [(\sigma^2 + \mu^2) + (n - 1)\mu^2]
\end{aligned}$$

Substituting back into the original parts,

$$\begin{aligned}
E(S^2) &= \frac{1}{n} \left( \underbrace{E \left[ \sum X_i^2 \right]}_{\textcircled{1}} - 2 \underbrace{E \left[ \sum X_i \bar{X} \right]}_{\textcircled{2}} + \underbrace{E \left[ \sum \bar{X} \right]^2}_{\textcircled{3}} \right) \\
&= \frac{1}{n} \left( \underbrace{n(\sigma^2 + \mu^2)}_{\textcircled{1}} - 2 \underbrace{\left( (\sigma^2 + \mu^2) + (n-1)\mu^2 \right)}_{\textcircled{2}} + \underbrace{\sigma^2 + n\mu^2}_{\textcircled{3}} \right) \\
&= \frac{1}{n} (n\sigma^2 + n\mu^2 - 2\sigma^2 - 2\mu^2 - 2(n-1)\mu^2 + \sigma^2 + n\mu^2) \\
&= \frac{1}{n} (n\sigma^2 + n\mu^2 - 2\sigma^2 - 2\mu^2 - 2n\mu^2 + 2\mu^2 + \sigma^2 + n\mu^2) \\
&= \frac{1}{n} (n\sigma^2 - 2\sigma^2 + \sigma^2) \\
&= \frac{1}{n} (n\sigma^2 - \sigma^2) \\
&= \frac{\sigma^2(n-1)}{n} \\
&\neq \sigma^2
\end{aligned}$$

Therefore, it is a biased point estimate.

To make an unbiased estimator, rescale the summation of the natural definition...

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

## 8.4 Evaluating the Goodness of a Point Estimator

**Definition:** Given any  $\hat{\theta}$ , the exact (theoretical) error is defined as

$$\varepsilon = |\hat{\theta} - \theta|$$

Note that  $\varepsilon$  is itself a random variable, and thus we can make probability statements about it. Obviously we want  $\varepsilon$  to be small. Consider  $|\hat{\theta} - \theta| < b$ . We can make conclusions about

$$P(|\hat{\theta} - \theta| < b)$$

Chebyshev:

$$P(|\hat{\theta} - \theta| < k\sigma) > 1 - \frac{1}{k^2}$$

*Remark:* We call  $\sigma = \sqrt{\sigma^2}$  the **standard error** in stats.

In practice,  $k = 2$  is a common choice.

$$\Pr(|\hat{\theta} - \theta| < 2\sigma) > 1 - \frac{1}{2^2} = \frac{3}{4} = 0.75$$

In the real world,  $2\sigma$  is usually much better than this (see Table 8.2) when the underlying distribution  $\hat{\theta}$  is symmetric and unimodal.

Last working assumption for a few lectures: in real life we use  $S^2 \approx \sigma^2$  (more later).

**Example:** (Difference in means)

2 iid random samples:

$$n_1 = 100 \quad \bar{Y}_1 = 26,400 \quad S_1^2 = 1,440,000$$

$$n_2 = 200 \quad \bar{Y}_2 = 25,100 \quad S_2^2 = 1,960,000$$

Estimating difference in means  $\mu_1 - \mu_2$ . Here,

$$\bar{Y}_1 - \bar{Y}_2 = \underbrace{26,400 - 25,100}_{\text{a point estimate}} = 1300$$

*Recall:*  $\sigma_{\bar{Y}_1 - \bar{Y}_2}^2 = \frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}$ . Use  $S_{\bar{Y}_1 - \bar{Y}_2}^2$  to find a  $2\sigma$  interval estimate:

$$\sigma_{\bar{Y}_1 - \bar{Y}_2}^2 \approx S_{\bar{Y}_1 - \bar{Y}_2}^2 = \frac{S_1^2}{n_1} + \frac{S_2^2}{n_2} \approx \frac{1,440,000}{100} + \frac{1,960,000}{200} = 22,800$$

$$\sigma_{\bar{Y}_1 - \bar{Y}_2} \approx S_{\bar{Y}_1 - \bar{Y}_2} = \sqrt{S_{\bar{Y}_1 - \bar{Y}_2}^2} = \sqrt{22,800} \approx 151$$

Therefore, a  $2\sigma$  interval estimate for  $\mu_1 - \mu_2$  is  $1300 \pm 2(151) \implies 1300 \pm 302$  or  $(998, 1602)$

## 8.5 Confidence Intervals

Goal: Take interval estimates from §8.4 but make more precise comments about the probability  $\Pr(|\hat{\theta} - \theta| < k\sigma)$

$$-k\sigma < \hat{\theta} - \theta < k\sigma \iff \hat{\theta} - k\sigma < \theta < \hat{\theta} + k\sigma$$

Consider

$$\Pr(\hat{\theta}_L \leq \theta \leq \hat{\theta}_U) = 1 - \alpha$$

$\hat{\theta}_L$  and  $\hat{\theta}_U$  are the lower and upper confidence limits, respectively.

$1 - \alpha$  is the confidence coefficient.

### Example:

Want 9% CI, chose  $\alpha = 5\%$  when we know how  $\hat{\theta}$  is distributed. We can use standardization methods to find the limits.

Two sided confidence interval:

$$\Pr(\hat{\theta}_L \leq \theta \leq \hat{\theta}_U) = 1 - \alpha$$

One sided:

$$\Pr(\hat{\theta}_L \leq \theta) = 1 - \alpha \quad \Pr(\hat{\theta}_U \geq \theta) = 1 - \alpha$$

**Discussion:** The pivotal method for finding confidence intervals.

- ① We know how a r.v.  $Y$  is distributed, but not some underlying parameter  $\theta$ .
- ② Using a distribution of an estimator  $\hat{\theta}$ , we convert to a probability distribution that does not depend on  $\theta$  (standardizing).

### Example:

$\bar{X}$  distributed  $N(\mu, \frac{\sigma^2}{n})$ , via

$$Z = \frac{\bar{X} - \mu}{\sqrt{\sigma^2/n}} \implies \underbrace{N(0, 1)}_{\text{Independent of } \mu}$$

We use  $N(0, 1)$  to rescale the limits of  $\hat{\theta}_L$  and  $\hat{\theta}_U$ .

### Example:

$Y_1, \dots, Y_n$  an iid random sample of size  $n$  from a uniform distribution on the interval  $(0, \theta)$ .

We want an estimate for  $\theta$ , use

$$\hat{\theta} = \max\{Y_1, \dots, Y_n\}$$

We know how  $\hat{\theta}$  is distributed, but it is clearly dependent on  $\theta$ .

$$\text{For } Y_i, \quad f(y) = \frac{1}{\theta}$$

$$\implies F(y) = \Pr(Y \leq y) = \int_0^y \frac{1}{\theta} = \frac{y}{\theta}, \quad y \in [0, \theta]$$

$$F(y) = \begin{cases} 0 & y < 0 \\ y/\theta & y \in [0, \theta] \\ 1 & y > \theta \end{cases}$$

The max order stat  $\hat{\theta}$  has CDF:

$$\Pr(\hat{\theta} < w) = \Pr(Y_1 \leq w, Y_2 \leq w, \dots, Y_n \leq w)$$

$$= [\Pr(Y \leq w)]^n$$

$$= \begin{cases} 0 & w < 0 \\ (w/\theta)^n & w \in [0, \theta] \\ 1 & w > \theta \end{cases}$$

Use a change of variables to find an associated pivotal distribution:  $U = \frac{\hat{\theta}}{\theta}$

The CDF of U,

$$\Pr(U \leq u) = \Pr\left(\frac{\hat{\theta}}{\theta} \leq u\right)$$

$$= \Pr(\hat{\theta} \leq \theta u)$$

$$= \left(\frac{\theta u}{\theta}\right)^n$$

$$= u^n \text{ for } u \in [0, 1]$$

$$F(u) = \begin{cases} 0 & u < 0 \\ u^n & u \in [0, 1] \\ 1 & u > 1 \end{cases}$$

Pivotal CDF of  $u$ , no longer depends upon  $\theta$ .

We use  $U$ 's CDF to construct a confidence interval.

Goal: Find a 95% lower confidence interval for  $\theta$ . Want:

$$\Pr\left(\underbrace{\hat{\theta}_L \leq \theta}_{\text{One-sided C.I.}}\right) = 0.95$$

Using  $U$ 's CDF:  $\Pr(U \leq u) = 0.95 \iff u^n = 0.95$

$$u = \sqrt[n]{0.95} = (0.95)^{1/n}$$

Then,

$$\begin{aligned}\Pr\left(\frac{\hat{\theta}}{\theta} \leq (0.95)^{1/n}\right) &= 0.95 \\ \Pr\left(\frac{\hat{\theta}}{(0.95)^{1/n}} \leq \theta\right) &= 0.95\end{aligned}$$

But  $\hat{\theta} = \max(Y_1, \dots, Y_n) = Y_{(n)}$

So, our 95% confidence interval is

$$\frac{Y_{(n)}}{(0.95)^{1/n}} \leq \theta.$$

### Example:

For a random sample

$$\underbrace{0.76, 0.88, 1.68, 1.74, 1.78}_{n = 5} \\ Y_{(5)} = 1.78$$

A 95% C.I. for  $\theta$  is given by  $n = 5$ ,  $Y_{(5)} = 1.78$

$$\frac{1.78}{(0.95)^{1/5}} \leq \theta$$

$$\begin{aligned}\frac{1.78}{0.98979} &\leq \theta \\ 1.798 &\leq \theta\end{aligned}$$

## 8.6 Large Sample Confidence Intervals

The unbiased point estimates for  $\mu$ ,  $\rho$ ,  $\mu_1 - \mu_0$ ,  $\rho_1 - \rho_2$  all have near Normal distributions by the Central Limit Theorem.

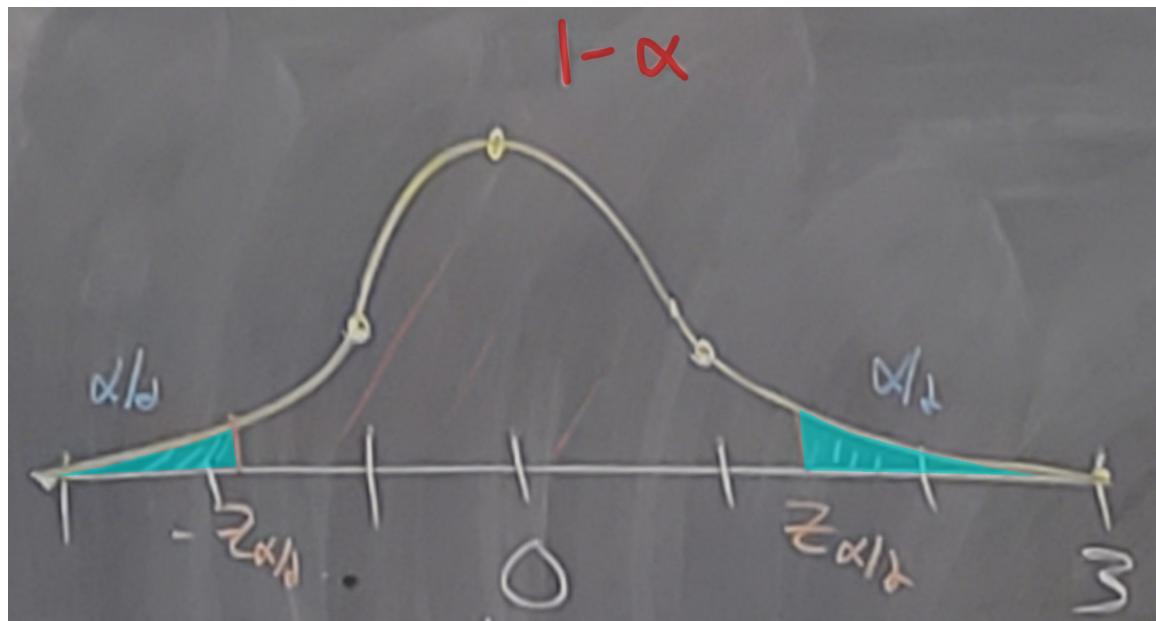
Moreover, using

$$Z = \frac{\hat{\theta} - \theta}{\sigma_{\hat{\theta}}}$$

$Z$  is a pivotal quantity  $N(0, 1)$ .

For two-sided confidence intervals,

$$\Pr(-Z_{\alpha/2} \leq Z \leq Z_{\alpha/2}) = 1 - \alpha$$



Some common standard errors:

$$90\% \text{ C.I.} \implies Z_{0.05} = 1.645$$

$$95\% \text{ C.I.} \implies Z_{0.025} = 1.960$$

$$99\% \text{ C.I.} \implies Z_{0.005} = 2.576$$

Then

$$-Z_{\alpha/2} \leq Z \leq Z_{\alpha/2}$$

$$-Z_{\alpha/2} \leq \frac{\hat{\theta} - \theta}{\sigma_{\hat{\theta}}} \leq Z_{\alpha/2}$$

$$-Z_{\alpha/2}\sigma_{\hat{\theta}} \leq \hat{\theta} - \theta \leq Z_{\alpha/2}\sigma_{\hat{\theta}}$$

$$-Z_{\alpha/2}\sigma_{\hat{\theta}} - \hat{\theta} \leq -\theta \leq Z_{\alpha/2}\sigma_{\hat{\theta}} - \hat{\theta}$$

$$\underbrace{\hat{\theta} - Z_{\alpha/2}\sigma_{\hat{\theta}}}_{\hat{\theta}_L} \leq \theta \leq \underbrace{\hat{\theta} + Z_{\alpha/2}\sigma_{\hat{\theta}}}_{\hat{\theta}_U}$$

Of course, we don't know  $\sigma_{\hat{\theta}}$  exactly. For "large" samples, we can use  $\sigma_{\hat{\theta}} \approx S_{\hat{\theta}}$ .

**Example:**

$$\bar{X} = 19.07 \quad S^2 = 10.60 \quad \text{with } n = 32$$

Recall: For "large" samples ( $n \geq 30$ ),  $\bar{X}$  nearly distributed by  $N(\mu, \sigma_{\bar{X}}^2) \approx N(\mu, \frac{10.60}{0.32})$ . For a 95% CI for  $\mu$ ,

$$\sigma_{\bar{X}} \approx \sqrt{\frac{10.60}{0.32}} \approx 0.576$$

Here,  $\alpha = 0.05$  and  $\frac{\alpha}{2} = 0.025$ . So  $\bar{X} \pm Z_{0.025}\sigma_{\bar{X}}$ , thus

$$\begin{aligned} & 19.07 \pm (1.96)(0.576) \\ & \Rightarrow 19.07 \pm 1.128 \quad \text{or} \quad (17.94, 20.20) \end{aligned}$$

**Discussion:** Differences in population proportions.

Estimating  $p_1 - p_2$  by  $\hat{p}_1 - \hat{p}_2$  for samples of size  $n_1$  and  $n_2$  respectively.

$$\begin{aligned} \hat{\theta} \pm Z_{\alpha/2}\sigma_{\hat{\theta}} & \implies \hat{p}_1 - \hat{p}_2 \pm Z_{\alpha/2}\sqrt{\sigma_{\hat{p}_1}^2 + \sigma_{\hat{p}_2}^2} \\ & \implies \hat{p}_1 - \hat{p}_2 \pm Z_{\alpha/2}\sqrt{\underbrace{\frac{p_1 q_1}{n_1} + \frac{p_2 q_2}{n_2}}_{\text{Depends on the unknowns } p_1, p_2}} \end{aligned}$$

Two standard fixes:

$$\textcircled{1} \quad pq = p(1-p) \leq 1/4$$

Easy to show by calculus. Yields a "max" error via

$$\sqrt{\frac{1}{4n_1} + \frac{1}{4n_2}}.$$

However, in practice we use \textcircled{2} as smaller confidence intervals are desirable.

$$\textcircled{2} \quad \text{When } n_i \text{ is large enough, we can use } \hat{p}_i \text{ for } p_i.$$

The  $1 - \alpha$  confidence interval is

$$\hat{p}_1 - \hat{p}_2 \pm Z_{\alpha/2}\sqrt{\frac{\hat{p}_1(1-\hat{p}_1)}{n_1} + \frac{\hat{p}_2(1-\hat{p}_2)}{n_2}}$$

## 8.7 Selecting the Sample Size

For our C.I. we have  $\hat{\theta} \pm Z_{\alpha/2}\sigma_{\hat{\theta}}$

- $\sigma_{\hat{\theta}}$  is the standard error
- $E = Z_{\alpha/2}\sigma_{\hat{\theta}}$  is the error bounds of our C.I.

In real world, we choose our confidence level  $1 - \alpha$  and/or the error  $E$  ahead of time. This leads to the problem of choosing an appropriate sample size  $n$ .

**Example:** For  $\bar{X}, \dots$  distributions is  $N(\mu, \sigma^2)$ .

$$\sigma_{\bar{X}} = \sqrt{\sigma^2/n} \quad \text{and} \quad E = Z_{\alpha/2}\sqrt{\sigma^2/n}$$

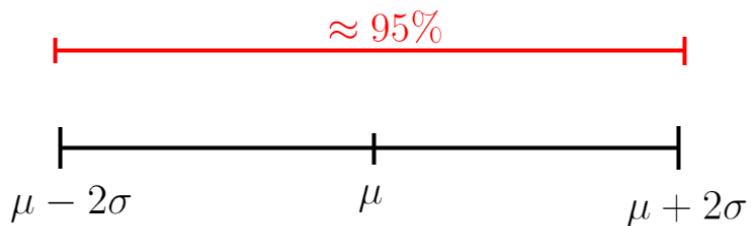
$$\longleftrightarrow \quad n = \frac{Z_{\alpha/2}\sigma^2}{E^2}$$

For desired  $E$ , we need to choose

$$n \geq \left\lceil \frac{Z_{\alpha/2}\sigma^2}{E^2} \right\rceil$$

Another issue: We don't really know  $\sigma_{\hat{\theta}}^2$  and maybe not even  $Z_{\alpha/2}$ .

- ① For  $Z_{\alpha/2}$ , we have seen that the empirical rule  $2\sigma$  usually works for large  $n$ .



- ② For  $\sigma_{\hat{\theta}}$ , either

- a.) Use old data or a “pilot” sample  $S^2$
- b.) Use  $\frac{1}{4}$  the spread of the data set, i.e.,

$$\sigma_{\hat{\theta}} \approx \frac{1}{4} (\max(Y_i) - \min(Y_i))$$

However, for proportions we can do even better.

**Example:** Proportions. For  $\hat{p} = \frac{\sum Y_i}{n}$  (relative frequency),  $\hat{p}$  distributed  $N(p, \frac{pq}{n})$

$$E = Z_{\alpha/2}\sqrt{\frac{pq}{n}} \quad \longleftrightarrow \quad n \geq \left\lceil \frac{Z_{\alpha/2}^2 p(1-p)}{E^2} \right\rceil$$

Again, since  $p(1-p) \leq \frac{1}{4}$  on  $[0, 1]$ , we choose  $n \geq \frac{Z_{\alpha/2}^2}{4E}$

**Example:** Public opinion polls ( $\pm 3\%$ )

For a 95% Confidence interval,

$$Z_{\alpha/2} = 1.960 \quad \text{and} \quad n \geq \left\lceil \frac{1.96^2}{4(0.03)^2} \right\rceil = 1068$$

For a 99% Confidence interval,

$$Z_{\alpha/2} = 2.576 \quad \text{and} \quad n \geq \left\lceil \frac{2.576^2}{4(0.03)^2} \right\rceil = 1844$$

## 8.8 Small-Sample Confidence Intervals for $\mu$ and $\mu_1 - \mu_2$

So for for  $\bar{X}$  when  $n$  is large, we have used

$$\frac{\hat{\theta} - \mu}{\sigma/\sqrt{n}} \approx \frac{\hat{\theta} - \mu}{S/\sqrt{n}}$$

for small  $n$ , no longer precise enough.

But we set ourselves up for this issue back in Chapter 7.

*Recall:* We know  $Z = \frac{\hat{\theta} - \mu}{\sigma/\sqrt{n}}$  is distributed  $N(0, 1)$

and  $V = \frac{(n-1)S^2}{\sigma^2}$  is distributed  $\chi^2$  with  $n-1$  df and  $Z$  and  $V$  are independent. Let

$$\begin{aligned} T &= \frac{\bar{X} - \mu}{S/\sqrt{n}} \\ &= \frac{\bar{X} - \mu}{S/\sqrt{n}} \cdot \frac{1}{s/\sigma} \\ &= \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \cdot \frac{1}{\sqrt{S^2/\sigma^2}} \\ &= Z \cdot \frac{1}{\sqrt{\frac{(n-1)S^2}{\sigma^2}/(n-1)}} \end{aligned}$$

We write  $T = \frac{Z}{\sqrt{V/r}}$  where  $r = (n-1)$ .

$T$  is the product of 2 pivotally distributed independent random variables.

**Definition:** We say  $T$  has a  $t$ -sampling distribution ( $t$ -dist'n) with  $(n-1)$  degrees of freedom and its pdf is given by

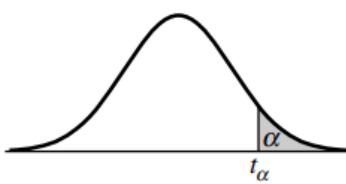
$$f_T(t) = \frac{\Gamma\left(\frac{r+1}{2}\right)}{\sqrt{\pi r} \cdot \Gamma\left(\frac{r}{2}\right)} \cdot \left(1 + \frac{t^2}{r}\right)^{-\frac{r+1}{2}}$$

*Remarks:*

- ① We won't work with  $f_T(t)$  explicitly. Note this distribution is pivotal... Independent of  $\mu$  and  $\sigma^2$ . We pick our  $t_{\alpha/2}$  from Table 5, page 849. Notice  $df = r = n-1 \geq 30$  are  $Z$ -scores:

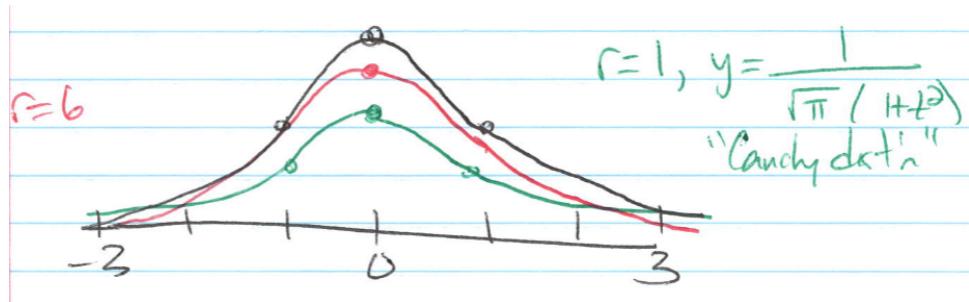
To do so we need  $\alpha$  and d.f. =  $n-1$

**Table 5 Percentage Points of the  $t$  Distributions**



$t_{.100}$	$t_{.050}$	$t_{.025}$	$t_{.010}$	$t_{.005}$	df
3.078	6.314	12.706	31.821	63.657	1
1.886	2.920	4.303	6.965	9.925	2
1.638	2.353	3.182	4.541	5.841	3
1.533	2.132	2.776	3.747	4.604	4
1.476	2.015	2.571	3.365	4.032	5
1.440	1.943	2.447	3.143	3.707	6
1.415	1.895	2.365	2.998	3.499	7
1.397	1.860	2.306	2.896	3.355	8
1.383	1.833	2.262	2.821	3.250	9
1.372	1.812	2.228	2.764	3.169	10
1.363	1.796	2.201	2.718	3.106	11
1.356	1.782	2.179	2.681	3.055	12
1.350	1.771	2.160	2.650	3.012	13
1.345	1.761	2.145	2.624	2.977	14
1.341	1.753	2.131	2.602	2.947	15
1.337	1.746	2.120	2.583	2.921	16
1.333	1.740	2.110	2.567	2.898	17
1.330	1.734	2.101	2.552	2.878	18
1.328	1.729	2.093	2.539	2.861	19
1.325	1.725	2.086	2.528	2.845	20
1.323	1.721	2.080	2.518	2.831	21
1.321	1.717	2.074	2.508	2.819	22
1.319	1.714	2.069	2.500	2.807	23
1.318	1.711	2.064	2.492	2.797	24
1.316	1.708	2.060	2.485	2.787	25
1.315	1.706	2.056	2.479	2.779	26
1.314	1.703	2.052	2.473	2.771	27
1.313	1.701	2.048	2.467	2.763	28
1.311	1.699	2.045	2.462	2.756	29
1.282	1.645	1.960	2.326	2.576	inf.

②  $t$ -distribution looks like “fat tailed” Normal Curves:



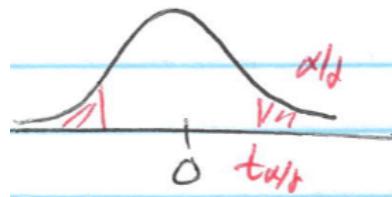
FACT: as  $r \rightarrow \infty$ , the D-distribution  $\rightarrow N(0, 1)$ .

FUN FACT: Cauchy distribution is “famous”.

$$r = 1, \quad \mu_T = 0, \quad \text{but} \quad \sigma_T^2 = \infty$$

**Discussion:** Small  $n$  Confidence Interval.

$$\text{For 2 sided, } T = \frac{\bar{X} - \mu}{S/\sqrt{n}}$$



$$\begin{aligned} \Pr(-t_{\alpha/2} \leq T \leq t_{\alpha/2}) &= 1 - \alpha \\ \implies \bar{X} \pm t_{\alpha/2} \left( \frac{S}{\sqrt{n}} \right) & \end{aligned}$$

**Example:**  $n = 10, \quad \bar{X} = 3.22, \quad S = 1.17$

Find a 95% confidence interval for  $\mu$ .

Table 5: 95%  $\rightarrow t_{0.025}$  with df =  $(n - 1) = 10 - 1 = 9$ .

Use  $t_{0.025} = 2.262$  and then  $\bar{X} \pm t_{0.025} \left( \frac{S}{\sqrt{n}} \right) \implies 3.22 \pm (2.262) \left( \frac{1.17}{\sqrt{10}} \right)$

$$3.22 \pm 0.84 \implies (2.38, 4.06)$$

BONUS MATH: From the t-distribution.

For the math curious, it comes down to a fancy change of variables.

Recall from last semester, if  $X, Y$  are independent, the joint pdf is  $f(x, y) = f_x(x)f_y(y)$ .

Here,  $T = Z \left( \frac{V}{r} \right)^{-\frac{1}{2}}$  and  $Z, V$  are independent (as  $\bar{X}$  and  $S$  are by Fisher's theorem).

To derive  $f_T$ , start with the joint pdf of  $Z$  and  $Y$ , a  $\chi^2$  distribution of  $r$  degrees of freedom. Then,

$$f(y, z) = \frac{1}{\Gamma\left(\frac{r}{2}\right) 2^{r/2}} \cdot y^{\frac{r}{2}-1} \cdot \exp\left[-\frac{r}{2}\right] \cdot \frac{1}{\sqrt{2\pi}} \cdot \exp\left[\frac{Z^2}{2}\right] \quad y \in (0, \infty), \quad Z \in \mathbb{R}$$

Then let  $t = Z \left( \frac{V}{r} \right)^{-\frac{1}{2}}$ ,  $V = y$ .

$$\begin{aligned} \Pr(Y < y, Z < z) &= \int_{-\infty}^y \int_{-\infty}^z f(y, z) dy dz \\ &= \int_{-\infty}^T \int_{-\infty}^V f(y(t, v), z(t, v)) \underbrace{\begin{vmatrix} Y_t & Y_v \\ Z_t & Z_v \end{vmatrix}^{-1}}_{\text{Jacobian}} dt dv \end{aligned}$$

Recall from MTH 420, for  $u = h(x)$ :

$$\int_a^b f(x) dx = \int_{h(a)}^{h(b)} f(h^{-1}(u)) \frac{d}{dx}(h^{-1}(u)) du$$

**Discussion:** Difference in means, small sample.

Start with  $X$ ,  $E[X] = \mu_x$ ,  $V[X] = \sigma_x^2$ , (Normal).  
And a r.v.  $Y$ ,  $E[Y] = \mu_y$ ,  $V[Y] = \sigma_y^2$ , (Normal).

Now we take 2 different iid random samples, which yields

$\bar{X}$  and  $S_{\bar{X}}^2$  with  $n_1$

$\bar{Y}$  and  $S_{\bar{Y}}^2$  with  $n_2$

We have the unbiased estimator  $\bar{X} - \bar{Y}$  with  $V(\bar{X} - \bar{Y}) = \frac{\sigma_x^2}{n_1} + \frac{\sigma_y^2}{n_2}$ . Then,

$$Z = \frac{(\bar{X} - \bar{Y}) - (\mu_x - \mu_y)}{\sqrt{\frac{\sigma_x^2}{n_1} + \frac{\sigma_y^2}{n_2}}} \quad \text{a pivotal quantity.}$$

Again, for small  $n_i$ , using  $S_i \approx \sigma_i$  is not precise enough.

**Assumption:** Let  $\sigma_{\bar{X}} = \sigma^{\bar{Y}}$ . Thus

$$Z = \frac{(\bar{X} - \bar{Y}) - (\mu_x - \mu_y)}{\sigma \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} \quad \text{a pivotal quantity.}$$

Now we need to construct a point estimator for  $\sigma$ .

**Idea:** Try a “pooled” estimator.

$$S_p^2 = \frac{\sum (X_i - \bar{X})^2 + \sum (Y_i - \bar{Y})^2}{r}$$

$$S_p^2 = \frac{(n_1 - 1)S_X^2 + (n_2 - 1)S_Y^2}{r}$$

**Claim:**  $S_p^2$  is unbiased when  $r = [(n_1 - 1) + (n_2 - 1)] = n_1 + n_2 - 2$ . (i.e.  $E[S_p^2] = \sigma^2$ )

**Reason:** It is really the same computation as when we showed that  $E[S^2] = \sigma^2$  when  $S = \frac{1}{n-1} \sum (X_i - \bar{X})^2$

Or, we can think about this as a product distribution of 2 independent  $\chi^2$  distributions. One with  $n_1 - 1$  df and another with  $n_2 - 1$  df. Then (like with mgfs) the product distribution is  $\chi^2$  with  $n_1 - 1 + n_2 - 1$  df.

**Definition:** The pooled estimator is

$$S_p^2 = \frac{\sum (X_i - \bar{X})^2 + \sum (Y_i - \bar{Y})^2}{n_1 + n_2 - 2}.$$

This is unbiased, and now we can define the t-statistic for  $\mu_X - \mu_Y$  as

$$T = \frac{(\bar{X} - \bar{Y}) - (\mu_X - \mu_Y)}{S_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}.$$

Which has  $n_1 + n_2 - 2$  degrees of freedom.

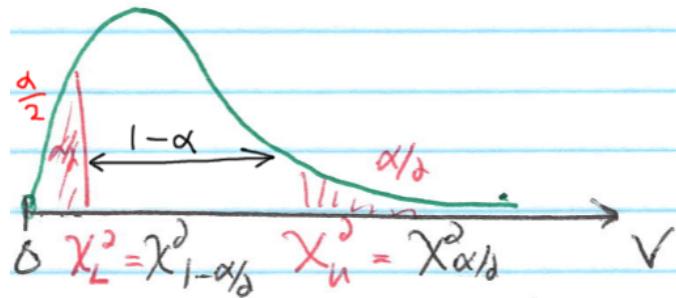
## 8.9 Confidence Intervals for $\sigma^2$

We have  $S = \frac{1}{n-1} \sum (X_i - \bar{X})^2$  is unbiased and  $V = \frac{(n-1)S^2}{\sigma^2}$  is a pivotal quantity with  $\chi^2$  distribution and  $(n-1)$  df.

For a confidence interval, we use

$$\Pr \left( \chi_L^2 \leq \frac{(n-1)S^2}{\sigma^2} \leq \chi_U^2 \right) = 1 - \alpha$$

As with t-distributions, we need both our  $\alpha$  and df  $n-1$ . But  $\chi^2$  itself is not symmetric. The most common C.I. makes symmetric probability in the “tails”



Then

$$\frac{\chi_L^2}{(n-1)S^2} \leq \frac{1}{\sigma^2} \leq \frac{\chi_U^2}{(n-1)S^2}$$

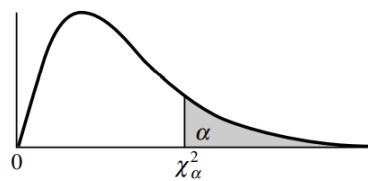
$$\Rightarrow \frac{1}{\chi_L^2} \geq \frac{\sigma^2}{(n-1)S^2} \geq \frac{1}{\chi_U^2}$$

$$\Rightarrow \frac{(n-1)S^2}{\chi_L^2} \geq \sigma^2 \geq \frac{(n-1)S^2}{\chi_U^2}$$

$$\Leftrightarrow \boxed{\frac{(n-1)S^2}{\chi_U^2} \leq \sigma^2 \leq \frac{(n-1)S^2}{\chi_L^2}}$$

Confidence interval for  $\sigma^2$

**Table 6 Percentage Points of the  $\chi^2$  Distributions**



df	$\chi^2_{0.995}$	$\chi^2_{0.990}$	$\chi^2_{0.975}$	$\chi^2_{0.950}$	$\chi^2_{0.900}$
1	0.0000393	0.0001571	0.0009821	0.0039321	0.0157908
2	0.0100251	0.0201007	0.0506356	0.102587	0.210720
3	0.0717212	0.114832	0.215795	0.351846	0.584375
4	0.206990	0.297110	0.484419	0.710721	1.063623
5	0.411740	0.554300	0.831211	1.145476	1.61031
6	0.675727	0.872085	1.237347	1.63539	2.20413
7	0.989265	1.239043	1.68987	2.16735	2.83311
8	1.344419	1.646482	2.17973	2.73264	3.48954
9	1.734926	2.087912	2.70039	3.32511	4.16816
10	2.15585	2.55821	3.24697	3.94030	4.86518
11	2.60321	3.05347	3.81575	4.57481	5.57779
12	3.07382	3.57056	4.40379	5.22603	6.30380
13	3.56503	4.10691	5.00874	5.89186	7.04150
14	4.07468	4.66043	5.62872	6.57063	7.78953
15	4.60094	5.22935	6.26214	7.26094	8.54675
16	5.14224	5.81221	6.90766	7.96164	9.31223
17	5.69724	6.40776	7.56418	8.67176	10.0852
18	6.26481	7.01491	8.23075	9.39046	10.8649
19	6.84398	7.63273	8.90655	10.1170	11.6509
20	7.43386	8.26040	9.59083	10.8508	12.4426
21	8.03366	8.89720	10.28293	11.5913	13.2396
22	8.64272	9.54249	10.9823	12.3380	14.0415
23	9.26042	10.19567	11.6885	13.0905	14.8479
24	9.88623	10.8564	12.4011	13.8484	15.6587
25	10.5197	11.5240	13.1197	14.6114	16.4734
26	11.1603	12.1981	13.8439	15.3791	17.2919
27	11.8076	12.8786	14.5733	16.1513	18.1138
28	12.4613	13.5648	15.3079	16.9279	18.9392
29	13.1211	14.2565	16.0471	17.7083	19.7677
30	13.7867	14.9535	16.7908	18.4926	20.5992
40	20.7065	22.1643	24.4331	26.5093	29.0505
50	27.9907	29.7067	32.3574	34.7642	37.6886
60	35.5346	37.4848	40.4817	43.1879	46.4589
70	43.2752	45.4418	48.7576	51.7393	55.3290
80	51.1720	53.5400	57.1532	60.3915	64.2778
90	59.1963	61.7541	65.6466	69.1260	73.2912
100	67.3276	70.0648	74.2219	77.9295	82.3581

**Table 6 (Continued)**

$\chi^2_{0.100}$	$\chi^2_{0.050}$	$\chi^2_{0.025}$	$\chi^2_{0.010}$	$\chi^2_{0.005}$	df
2.70554	3.84146	5.02389	6.63490	7.87944	1
4.60517	5.99147	7.37776	9.21034	10.5966	2
6.25139	7.81473	9.34840	11.3449	12.8381	3
7.77944	9.48773	11.1433	13.2767	14.8602	4
9.23635	11.0705	12.8325	15.0863	16.7496	5
10.6446	12.5916	14.4494	16.8119	18.5476	6
12.0170	14.0671	16.0128	18.4753	20.2777	7
13.3616	15.5073	17.5346	20.0902	21.9550	8
14.6837	16.9190	19.0228	21.6660	23.5893	9
15.9871	18.3070	20.4831	23.2093	25.1882	10
17.2750	19.6751	21.9200	24.7250	26.7569	11
18.5494	21.0261	23.3367	26.2170	28.2995	12
19.8119	22.3621	24.7356	27.6883	29.8194	13
21.0642	23.6848	26.1190	29.1413	31.3193	14
22.3072	24.9958	27.4884	30.5779	32.8013	15
23.5418	26.2962	28.8454	31.9999	34.2672	16
24.7690	27.5871	30.1910	33.4087	35.7185	17
25.9894	28.8693	31.5264	34.8053	37.1564	18
27.2036	30.1435	32.8523	36.1908	38.5822	19
28.4120	31.4104	34.1696	37.5662	39.9968	20
29.6151	32.6705	35.4789	38.9321	41.4010	21
30.8133	33.9244	36.7807	40.2894	42.7956	22
32.0069	35.1725	38.0757	41.6384	44.1813	23
33.1963	36.4151	39.3641	42.9798	45.5585	24
34.3816	37.6525	40.6465	44.3141	46.9278	25
35.5631	38.8852	41.9232	45.6417	48.2899	26
36.7412	40.1133	43.1944	46.9630	49.6449	27
37.9159	41.3372	44.4607	48.2782	50.9933	28
39.0875	42.5569	45.7222	49.5879	52.3356	29
40.2560	43.7729	46.9792	50.8922	53.6720	30
51.8050	55.7585	59.3417	63.6907	66.7659	40
63.1671	67.5048	71.4202	76.1539	79.4900	50
74.3970	79.0819	83.2976	88.3794	91.9517	60
85.5271	90.5312	95.0231	100.425	104.215	70
96.5782	101.879	106.629	112.329	116.321	80
107.565	113.145	118.136	124.116	128.299	90
118.498	124.342	129.561	135.807	140.169	100

**Example:** Construct a 90% C.I. for  $\mu$  and  $\sigma^2$ .

$$\underbrace{85.4 \quad 86.8 \quad 86.1 \quad 85.3 \quad 84.8 \quad 86.0}_{\text{Data set, } n = 6}$$

$$\bar{X} = 85.7\bar{3} \quad S^2 = 0.502\bar{6} \quad S \approx 0.7089$$

$$90\% \text{ C.I.} \implies \alpha = 0.10, \frac{\alpha}{2} = 0.05 \text{ with df } n - 1 = 5.$$

$$\chi_L^2 = \chi_{0.950}^2 = 1.145476 \quad \chi_U^2 = \chi_{0.050}^2 = 11.0705$$

A common mistake people make is forgetting that the lower bound is  $\chi_U^2$ . For  $\sigma^2$ ,

$$\begin{aligned} \frac{5 \cdot (0.502\bar{6})}{11.0705} &\leq \sigma^2 \leq \frac{5 \cdot (0.502\bar{6})}{1.145475} \\ 0.227 &\leq \sigma^2 \leq 2.194 \\ \sigma^2 &\in (0.227, 2.194) \end{aligned}$$

For  $\mu$ , use a  $t$ -distribution with 5 degrees of freedom.

$$\bar{X} \pm t_{0.050} \cdot S \implies 85.7\bar{3} \pm (2.015)(0.7089)$$

$$\mu \in (84.305, 87.162)$$

# Chapter 9

## Properties of Point Estimators and Methods of Estimation

### 9.2 Relative Efficiency

In general, given two estimators  $\hat{\theta}_1$  and  $\hat{\theta}_2$  of a parameter  $\theta$ , we claim the one with smaller MSE is better.

For unbiased estimators, MSE is variance ( $\sigma^2$ ). So an idea of “better” is if  $\text{Var}(\hat{\theta}_1) < \text{Var}(\hat{\theta}_2)$ , we say  $\hat{\theta}_1$  is the better estimator.

**Definition:** Given unbiased estimators  $\hat{\theta}_1$  and  $\hat{\theta}_2$  of  $\theta$ , the efficiency of  $\hat{\theta}_1$  relative to  $\hat{\theta}_2$  is defined by

$$\text{eff}(\hat{\theta}_1, \hat{\theta}_2) = \frac{\text{Var}(\hat{\theta}_2)}{\text{Var}(\hat{\theta}_1)}.$$

*Note:* If  $\text{eff}(\hat{\theta}_1, \hat{\theta}_2) > 1$ , then  $\frac{\text{Var}(\hat{\theta}_2)}{\text{Var}(\hat{\theta}_1)} > 1 \implies \underbrace{\text{Var}(\hat{\theta}_1) < \text{Var}(\hat{\theta}_2)}_{\hat{\theta}_1 \text{ is “better”}}$

**Example:** Let  $Y_1, \dots, Y_n$  be an iid random sample from  $N(\mu, \sigma^2)$ . Two estimators of  $\sigma^2$  are:

$$\hat{\sigma}_1^2 = S^2 = \frac{1}{n-1} \sum_{i=1}^n (Y_i - \bar{Y})^2$$

$$\hat{\sigma}_2^2 = \frac{1}{2} (Y_1 - Y_2)^2$$

We know  $S^2$  is unbiased.

Show  $\hat{\sigma}_2^2$  is unbiased.

$$\begin{aligned}
E(\hat{\sigma}_2^2) &= E\left(\frac{1}{2}(Y_1^2 - 2Y_1Y_2 + Y_2^2)\right) \\
&= \frac{1}{2}\left[E(Y_1^2) - 2E(Y_1Y_2) + E(Y_2^2)\right] \\
&= \frac{1}{2}\left[\text{Var}(Y_1) + (E(Y_1))^2 - 2E(Y_1)E(Y_2) + \text{Var}(Y_2) + (E(Y_2))^2\right] \\
&= \frac{1}{2}\left[\sigma^2 + \mu^2 - 2\mu \cdot \mu + \sigma^2 + \mu^2\right] \\
&= \sigma^2 \quad \text{unbiased.}
\end{aligned}$$

To compute relative efficiency, we need variances.

For  $\text{Var}(S^2)$ , we know  $\frac{(n-1)S^2}{\sigma^2} \sim \chi^2(n-1)$ .

By properties of  $\chi^2$  distribution,  $\text{Var}\left(\frac{(n-1)S^2}{\sigma^2}\right) = 2(n-1)$ .

Then,  $\frac{(n-1)^2}{\sigma^4} \cdot \text{Var}(S^2) = 2(n-1) \implies \text{Var}(S^2) = \frac{2\sigma^4}{n-1}$ .

For  $\text{Var}(\hat{\sigma}_2^2)$ , this takes more work.

$$\begin{aligned}
\text{Var}\left(\frac{1}{2}(Y_1 - Y_2)^2\right) &= E\left(\underbrace{\left(\frac{1}{2}(Y_1 - Y_2)^2\right)^2}_{*}\right) - \underbrace{\left[E\left(\frac{1}{2}(Y_1 - Y_2)^2\right)\right]^2}_{\sigma^4} \\
* &= \frac{1}{4}E((Y_1 - Y_2)^4) \\
&= \frac{1}{4}E(Y_1^4 - 4Y_1^3Y_2 + 6Y_1^2Y_2^2 - 4Y_1Y_2^3 + Y_2^4) \\
&= \frac{1}{4}\left[E(Y_1^4) - 4E(Y_1^3)E(Y_2) + 6E(Y_1^2)E(Y_2^2) - 4E(Y_1)E(Y_2^3) + E(Y_2^4)\right] \\
&= \frac{1}{4}\left[2E(Y^4) - 8E(Y^3)E(Y) + 6E(Y^2)^2\right]
\end{aligned}$$

We need the higher moments  $m'_3$  and  $m'_4$ . Recall the mgf for  $N(\mu, \sigma^2)$  is  $m(t) = \exp\left[\mu t + \frac{t^2\sigma^2}{2}\right]$

$$\begin{aligned}
m'(t) &= (\mu + t\sigma^2)m(t) \\
m'(0) &= \mu m(0) \\
&= \mu \\
&= E(Y)
\end{aligned}$$

$$\begin{aligned}
m''(t) &= \sigma^2 m(t) + (\mu + t\sigma^2)^2 m(t) \\
m''(0) &= \sigma^2 + \mu^2 m(0) \\
&= \sigma^2 + \mu^2 \\
&= E(Y^2)
\end{aligned}$$

$$\begin{aligned}
m^{(3)}(t) &= \sigma^2 m'(t) + 2(\mu + t\sigma^2)\sigma^2 m(t) + (\mu + t\sigma^2)^2 m'(t) \\
m^{(3)}(0) &= \sigma^2 m'(0) + 2(\mu + 0\sigma^2)\sigma^2 m(0) + (\mu + 0\sigma^2)^2 m'(0) \\
&= \sigma^2 \mu + 2\mu\sigma^2 + \mu^3 \\
&= \mu^3 + 3\mu\sigma^2 \\
&= E(Y^3)
\end{aligned}$$

$$\begin{aligned}
m^{(4)}(t) &= \dots \\
m^{(4)}(0) &= \dots \\
&= \mu^4 + 6\mu^2\sigma^2 + 3\sigma^4 \\
&= E(Y^4)
\end{aligned}$$

Substituting back into \*, we get  $* = \frac{1}{4}(12\sigma^4) = 3\sigma^4$ . Thus,

$$\text{Var}(\widehat{\sigma}_2^2) = \text{Var}\left(\frac{1}{2}(Y_1 - Y_2)^2\right) = \underbrace{3\sigma^4}_{*} - \underbrace{\sigma^4}_{\sigma^4} = 2\sigma^4$$

$$\text{Hence, } \text{eff}(S^2, \widehat{\sigma}_2^2) = \frac{2\sigma^4}{\frac{2\sigma^4}{n-1}} = n-1$$

## 9.3 Consistency

**Discussion:** “Convergence in probability”

This is a different type of analysis than that of Calculus.

**Example:** MTH 420, “pointwise” convergence:

$$f_n(x) = x^n, \quad x \in [0, 1] \quad \lim_{n \rightarrow \infty} f_n(x) = \begin{cases} 0 & x \in [0, 1) \\ 1 & x = 1 \end{cases}$$

Let  $x \in (0, 1)$ . To show  $x^n \rightarrow 0$ :

For any  $\epsilon > 0$ , need an index  $N$  such that when  $n > N$  we get  $|x^n - 0| < \epsilon$

$$\text{Choose } N > \left\lceil \frac{\ln \epsilon}{\ln x} \right\rceil$$

In probability, convergence is about the probabilistic measure of an event.

**Example:** Law of Large Numbers:

$A$ , an event associated with an event  $E$ .  $\Pr(A) = p$ .

Do  $n$  iid repetitions of  $E$  and count  $n_A :=$  number of times  $A$  occurs.

The relative frequency  $f_A = \frac{n_A}{n}$ . **The Law of Large Numbers says  $f_A \rightarrow p$**

What does this mean exactly? Conclusion:  $\underbrace{\Pr(|f_A - p| < \epsilon)}_{*} \geq 1 - \frac{p(1-p)}{n\epsilon^2}$

\* This means the probability that  $p$  is in the confidence interval  $f_A - \epsilon \leq p \leq f_A + \epsilon$ .

$$\text{Convergence? } \lim_{n \rightarrow \infty} \Pr(|f_A - p| < \epsilon) \geq \lim_{n \rightarrow \infty} \left( 1 - \frac{p(1-p)}{n\epsilon^2} \right) = 1$$

$$\implies \lim_{n \rightarrow \infty} \Pr(|f_A - p| < \epsilon) = 1$$

**This is convergence in probability!**

In probability, convergence is about the probabilistic measure of an event (measure theory).

**Definition:** Estimator  $\hat{\theta}$  is consistent if  $\hat{\theta} \rightarrow \theta$  in probability.

That is  $\hat{\theta}_n$  is a consistent estimator of  $\theta$  if for any  $\epsilon > 0$

$$\lim_{n \rightarrow \infty} \Pr \left( |\hat{\theta}_n - \theta| \leq \epsilon \right) = 1$$

$$\text{or } \lim_{n \rightarrow \infty} \Pr \left( |\hat{\theta}_n - \theta| > \epsilon \right) = 0$$

**Discussion:** A tool for showing consistency:

**Theorem** An unbiased estimator  $\hat{\theta}_n$  for  $\theta$  is consistent if

$$\lim_{n \rightarrow \infty} \text{Var}(\hat{\theta}_n) = 0$$

*Proof.* Let  $\epsilon > 0$ , then consider,

$$\Pr(|\hat{\theta}_n - \theta| > \epsilon)$$

Recall Chebyshev's Theorem

$$\Pr(|Y - \mu| > k\sigma) \leq \frac{1}{k^2}$$

$$\text{Here } \epsilon = k\sigma_{\hat{\theta}_n} \implies k = \frac{\epsilon}{\sigma_{\hat{\theta}_n}}, \quad 0 \leq \Pr(|\hat{\theta}_n - \theta| > \epsilon) \leq \frac{1}{(\epsilon/\sigma_{\hat{\theta}_n})^2}$$

Then

$$\Pr(|\hat{\theta}_n - \theta| > \epsilon) \leq \frac{\sigma_{\hat{\theta}_n}^2}{\epsilon^2} = \frac{\text{Var}(\hat{\theta}_n)}{\epsilon^2}$$

As  $n \rightarrow \infty$ ,

$$\begin{aligned} 0 &\leq \lim_{n \rightarrow \infty} \Pr(|\hat{\theta}_n - \theta| > \epsilon) \leq \lim_{n \rightarrow \infty} \frac{\text{Var}(\hat{\theta}_n)}{\epsilon^2} = 0 \\ &\implies \Pr(|\hat{\theta}_n - \theta| > \epsilon) = 0 \end{aligned}$$

Hence,  $\hat{\theta}_n \rightarrow \theta$  in probability. □

**Example:** (Common) Consistent Estimators

①  $\bar{Y}$  know unbiased for  $\mu$  and

$$\text{Var}(\bar{Y}) = \frac{\sigma^2}{n}$$

Provided that  $\sigma^2$  is finite,  $\lim_{n \rightarrow \infty} \text{Var}(\bar{Y}) = 0$

$\bar{Y}$  is consistant, that is,  $\bar{Y} \rightarrow \mu$

②  $\hat{p}$  unbiased for  $p$ .

$$\begin{aligned} \text{Var}(\hat{p}) &= \frac{pq}{n} \quad \text{and} \quad \lim_{n \rightarrow \infty} \text{Var}(\hat{p}) = 0 \\ &\implies \hat{p} \rightarrow p \text{ i.e. consistent} \end{aligned}$$

**Example: (9.17 & 9.18)**

$X_1, \dots, X_n$  and  $Y_1, \dots, Y_n$  are iid random samples with  $\mu_x$  and  $\mu_y$  but  $\text{Var}(X) = \text{Var}(Y) = \sigma^2$ .

a.) Show that  $\bar{X} - \bar{Y}$  is a consistent estimator of  $\mu_x - \mu_y$ .

**Solution:** We showed in Chapter 8 that  $E(\bar{X} - \bar{Y}) = \mu_x - \mu_y$  is unbiased.

$$\text{Then } \text{Var}(\bar{X} - \bar{Y}) = \text{Var}(\bar{X}) + (-1)^2 \text{Var}(\bar{Y}) = \frac{2\sigma^2}{n}.$$

Because  $\lim_{n \rightarrow \infty} \text{Var}(\bar{X} - \bar{Y}) = 0$ , we have consistency.

b.) Show that the pooled estimator,

$$S_p^2 = \frac{\sum_1^n (X_i - \bar{X})^2 + \sum_1^n (Y_i - \bar{Y})^2}{2n - 2}$$

is a consistent estimator of  $\sigma^2$  when  $X, Y \sim N(\mu, \sigma^2)$ .

**Solution:** To show unbiased we need  $\text{Var}(S_p^2)$ .

“Trick,” convert to a distribution we understand well.

$$\begin{aligned} \underbrace{\frac{(2n-2)S_p^2}{\sigma^2}}_{\chi^2(2n-2) \text{ By Fisher's}} &= \frac{\sum_{i=1}^n (X_i - \bar{X})^2 + \sum_{i=1}^n (Y_i - \bar{Y})^2}{\sigma^2} \\ &= \underbrace{\sum_{i=1}^n \left( \frac{X_i - \bar{X}}{\sigma} \right)^2 + \sum_{i=1}^n \left( \frac{Y_i - \bar{Y}}{\sigma} \right)^2}_{U_n \sim \chi^2(2n-2)} \end{aligned}$$

$$\text{E} \left( \frac{(2n-2)S_p^2}{\sigma^2} \right) = 2n-2 \text{ by properties of } \chi^2$$

$$\frac{2n-2}{\sigma^2} \text{E}(S_p^2) = 2n-2 \implies \text{E}(S_p^2) = \sigma^2 \text{ unbiased.}$$

For consistency, use same “trick.”

$$\begin{aligned} \text{Var} \left( \frac{(2n-2)S_p^2}{\sigma^2} \right) &= \text{Var}(U_n) \\ \left( \frac{2n-2}{\sigma^2} \right)^2 \text{Var}(S_p^2) &= 2 \cdot (2n-2) \quad \text{by properties of } \chi^2 \\ \text{Var}(S_p^2) &= \frac{2\sigma^4}{2n-2} \\ \lim_{n \rightarrow \infty} \text{Var}(S_p^2) &= 0 \text{ therefore } S_p^2 \text{ is consistent.} \end{aligned}$$

**Theorem:** “The Limit Laws” (For probability convergence)

Let  $\hat{\theta}_n \rightarrow \theta$  and  $\hat{\Psi}_n \rightarrow \Psi$ . Then,

$$\textcircled{1} \quad \hat{\theta}_n + \hat{\Psi}_n \rightarrow \theta + \Psi$$

$$\textcircled{2} \quad \hat{\theta}_n \cdot \hat{\Psi}_n \rightarrow \theta \cdot \Psi$$

$$\textcircled{3} \quad \frac{\hat{\theta}_n}{\hat{\Psi}_n} \rightarrow \frac{\theta}{\Psi}$$

$$\textcircled{4} \quad \text{If } g \text{ is a continuous function at } \theta, \text{ then } g(\hat{\theta}_n) \rightarrow g(\theta)$$

**Example:**  $S_p^2$  again.

$$\begin{aligned} S_p^2 &= \frac{\sum_1^n (X_i - \bar{X})^2 + \sum_1^n (Y_i - \bar{Y})^2}{2n - 2} \\ &= \frac{(n-1)S_{\bar{X}}^2 + (n-1)S_{\bar{Y}}^2}{2n-2} \\ &= \frac{S_{\bar{X}}^2 + S_{\bar{Y}}^2}{2} \end{aligned}$$

But we already know  $S_{\bar{X}}^2 \rightarrow \sigma^2$  and  $S_{\bar{Y}}^2 \rightarrow \sigma^2$

Hence, by the Limit laws,

$$S_p^2 \rightarrow \frac{\sigma^2 + \sigma^2}{2} = \sigma^2$$

**Discussion:** Consistency and large  $n$  confidence intervals

We have  $\bar{Y}_n \pm Z_{\alpha/2} \frac{S_n}{\sqrt{n}}$  By consistency,  $\bar{Y}_n \rightarrow \mu$  and  $S_n \rightarrow \sigma$

Hence, when  $\sigma$  is finite,  $\bar{Y}_n \pm Z_{\alpha/2} \frac{S_n}{\sqrt{n}} \rightarrow \mu$

$$\text{Hence } \lim_{n \rightarrow \infty} \Pr \left( \left| \bar{Y}_{(n)} - \mu \right| < Z_{\alpha/2} \frac{S_n}{\sqrt{n}} \right) = 1$$

That is,  $\bar{Y}_{(n)} \pm Z_{\alpha/2} \frac{S_n}{\sqrt{n}} \rightarrow \mu$  in probability.

“Consistency justified our working assumption in Chapter 8”

**Example: (9.24)**

Let  $Z_1, \dots, Z_n$  be an iid random sample with  $Z \sim N(0, 1)$

$$\text{Let } U_n = \sum_{i=1}^n Z_i^2 \quad Z = \frac{Z_i - \mu}{\sigma}$$

We “know”  $U_n \sim \chi^2(n)$

Showed in old HW that  $Z^2 \sim X^2(n)$

In notes, we used the product of  $n$  mgfs to show that  $U_n \sim \chi^2(n)$ .

Hence,  $E(U_n) = n$  and  $\text{Var}(U_n) = 2n$ . Let

$$W_n := \frac{1}{n} U_n$$

Note that  $E(W_n) = \frac{1}{n} E(U_n) = 1$

Hence,  $W_n$  is an unbiased estimator when  $\sigma^2 = 1$

$$\begin{aligned} \text{Var}(W_n) &= \text{Var}\left(\frac{1}{n} U_n\right) \\ &= \frac{1}{n^2} \text{Var}(U_n) \\ &= \frac{1}{n^2} \cdot 2n \\ &= \frac{2}{n} \end{aligned}$$

Since  $\lim_{n \rightarrow \infty} \text{Var}(W_n) = 0$ , we get  $W_n \rightarrow 1$  in probability.

## 9.4 Sufficiency

**Discussion:** Among a collection of estimators, how do we choose which one to work with and why?

Conceptually: A statistic is deemed **sufficient** if it “contains all of the available information about a parameter,” given some random sample.

**Example:**

Statistician 1 has the data  $X_1, \dots, X_n$  and computes an estimator  $\hat{\theta}$ .

Statistician 2 has only the stat  $T = \tau(X_1, \dots, X_n)$  estimating  $\theta$ .

Sufficient implies both statisticians can make equally correct estimates about  $\theta$ . Equivalently, a sufficient estimator  $\hat{\theta}$  utilizes all of the information in a sample relating to  $\theta$ .

In either case, both are able to make the “same” conclusions about  $\theta$ .

**Definition:** Let  $Y_1, \dots, Y_n$  be a sample from a probability distribution that is known up to parameter  $\theta$ .

The statistic  $T = \tau(Y_1, \dots, Y_n)$  is **sufficient for  $\theta$**  if the conditional distribution of  $Y_1, \dots, Y_n$  given  $T$  does not depend on  $\theta$ .

*Remark: **Sufficiency Principle:***

Given 2 data sets  $X_1, \dots, X_n$  and  $Y_1, \dots, Y_n$  where the stats

$$\tau(X_1, \dots, X_n) = \tau(Y_1, \dots, Y_n) = T$$

Any inference about  $\theta$  should be the same regardless of the data set.

**Example:** (One time only via the definition)

Let  $Y \sim \text{Bern}(p)$

i.e.  $\Pr(Y = 1) = p$ ,  $\Pr(Y = 0) = 1 - p$ .

Let  $Y_n$  be a sequence of observations. We have  $\hat{p} = \frac{\sum Y_i}{n}$ .

**Sufficiency Principle:**

Any sequence with  $k$  1's and  $n - k$  0's is going to yield the same  $\hat{p}$ . In other words, same inferences about  $p$ .

Let  $S = \sum Y_i$

To show  $S$  is sufficient for  $p$ , we must show that the conditional distribution of the data  $Y_1, \dots, Y_n$  given  $S$  does not depend on  $p$ .

For  $k \in (0, 1, 2, \dots, n)$ ,  $S = Y_1 + \dots + Y_n = k$ . We have,

$$\begin{aligned}
\Pr(Y_1 = y_1, \dots, Y_n = y_n \mid S = k) &= \frac{\Pr(Y_1 = y_1, \dots, Y_n = y_n \text{ and } S = k)}{\Pr(S = k)} \\
&= \begin{cases} 0 & \sum Y_i \neq k \\ \frac{\Pr(Y_1 = y_1, \dots, Y_n = y_n \text{ and } S = k)}{\Pr(S = k)} & \sum Y_i = k \end{cases} \\
&= \underbrace{\prod_{i=1}^n \Pr(Y_i = y_i)}_{\substack{\text{independent} \\ \text{Binomial, } k \text{ "successes" in } n \text{ trials}}} \\
&= \frac{p^{y_1}(1-p)^{1-y_1} \times \dots \times p^{y_n}(1-p)^{1-y_n}}{\binom{n}{k} p^k (1-p)^{n-k}} \\
&= \frac{p^{\sum_1^n Y_i} (1-p)^{n - \sum_1^n Y_i}}{\binom{n}{k} p^k (1-p)^{n-k}} \\
&= \frac{p^k (1-p)^{n-k}}{\binom{n}{k} p^k (1-p)^{n-k}} \\
&= \begin{cases} 0 & S \neq k \\ \frac{1}{\binom{n}{k}} & S = k \end{cases} \quad \leftarrow \text{ independent of } p
\end{aligned}$$

By definition,  $S = \sum Y_i$  is sufficient.

*Remark:* Started with  $\hat{p} = \frac{S}{n}$ .

**Note:** Any one-to-one transformation (injection) of a sufficient stat will again be sufficient.

**Example:**  $f(x) = \frac{x}{n}$  is an injection (has invertible form).

$$S \text{ sufficient} \implies \frac{S}{n} \text{ sufficient}$$

Thus  $\hat{p}$  is also sufficient.

**Example:**  $\sqrt{x}$  is also an injection on  $[0, 1]$ .

Thus  $\sqrt{S}$  is sufficient.

**Topic:** The Factorization Theorem:

In practice we do not use the definition.

**Theorem: The Factorization Theorem**

Let  $X_1, \dots, X_n$  denote a random sample from a distribution with pdf  $f(x | \theta)$

(distribution depends on the unknown parameter  $\theta$ )

The statistic  $T = \tau(x_1, \dots, x_n)$  is a sufficient statistic for  $\theta$  if and only if the conditional pdf can be factored as follows.

$$L(x_1, \dots, x_n | \theta) = g(\tau(x_1, \dots, x_n) | \theta) \cdot h(x_1, \dots, x_n)$$

Where

- $g$  is a function that depends on the data only through the stat  $T$
- $h$  does not depend on  $\theta$

*Remarks:*

- ① Via independence (of data)

$$L(x_1, \dots, x_n | \theta) = f(x_1 | \theta) \cdot f(x_2 | \theta) \times \dots \times f(x_n | \theta)$$

The product of the marginals.

- ② The function  $L$  is called the likelihood of the sample.

**Example:** Let  $X_1, \dots, X_n$  denote a random sample from a Poisson distribution with parameter  $\lambda > 0$ . Find a sufficient statistic for the parameter  $\lambda$ .

Here,

$$f(x) = \frac{\lambda^x e^{-\lambda}}{x!}, \quad x = 0, 1, 2, \dots$$

$$\begin{aligned} L(x_1, x_2, \dots, x_n | \lambda) &= f(x_1 | \lambda) \cdot f(x_2 | \lambda) \times \dots \times f(x_n | \lambda) \\ &= \frac{\lambda^{x_1} e^{-\lambda}}{x_1!} \cdot \frac{\lambda^{x_2} e^{-\lambda}}{x_2!} \times \dots \times \frac{\lambda^{x_n} e^{-\lambda}}{x_n!} \\ &= \underbrace{\lambda^{x_1 + \dots + x_n} e^{-\lambda n}}_{S := \sum X_i} \cdot \underbrace{\frac{1}{x_1! \cdot x_2! \cdot \dots \cdot x_n!}}_{h(x_1, \dots, x_n)} \end{aligned}$$

With  $S = \sum X_i$ ,

$$g(S | \lambda) = e^{-\lambda n} \lambda^S$$

and

$$h(x_1, \dots, x_n) = \frac{1}{x_1! \cdot x_2! \cdot \dots \cdot x_n!}$$

By Factorization Theorem,  $S$  is a sufficient statistic.

**Note:**  $\bar{X} = \frac{1}{n} \sum X_i$  or  $n\bar{X} = S = \sum X_i$ . So,

$$g(S \mid \lambda) = e^{-\lambda} \lambda^{n\bar{X}} = g(\bar{X} \mid \lambda)$$

Therefore  $\bar{X}$  is also a sufficient statistic.

**Example:** The  $\hat{p}$  example again.  $Y \sim \text{Bern}(p)$ .

$$\begin{aligned} L(y_1, \dots, y_n \mid p) &= f(y_1 \mid p) \times \dots \times f(y_n \mid p) \\ &= p^{y_1}(1-p)^{1-y_1} \times \dots \times p^{y_n}(1-p)^{1-y_n} \\ &= p^{\sum y_i}(1-p)^{n-\sum y_i} \quad S := \sum y_i \\ &= \underbrace{p^S(1-p)^{n-S}}_{g(S \mid p)} \cdot \underbrace{1}_{h(\vec{y})} \end{aligned}$$

Where  $\vec{y}$  is the  $n$  tuple  $y_1, y_2, \dots, y_n$ .

By Factorization Theorem,  $S$  is sufficient.

**Yesterday:** Sufficiency  $\underbrace{f(x_1, \dots, x_n \mid S)}_{S \text{ is sufficient for } \theta} = \frac{f(x_1, \dots, x_n \text{ and } S(p) = s)}{\underbrace{f(s)}_{\text{conditional prob is independent of } \theta}}$

Factorization theorem:

$$L(x_1, \dots, x_n \mid \mu) = \underbrace{g(\mu, s)}_{S \text{ is sufficient for } \theta} \cdot h(\vec{x})$$

**Example:** Let  $X_1, \dots, X_n$  be a random sample from  $N(\mu, 1)$ . (unknown  $\mu$ ,  $\sigma^2 = 1$ .)

$$f(x) = \frac{1}{\sqrt{2\pi} \exp \left[ \frac{-(x-\mu)^2}{2} \right]}$$

Find a sufficient statistic for  $\mu$ .

Construct Likelihood function,

$$\begin{aligned}
L(x_1, \dots, x_n | \mu) &= f(x_1 | \mu) \cdots f(x_n | \mu) \\
&= \frac{1}{\sqrt{2\pi}} \exp \left[ \frac{-(x_1 - \mu)^2}{2} \right] \times \cdots \times \frac{1}{\sqrt{2\pi}} \exp \left[ \frac{-(x_n - \mu)^2}{2} \right] \\
&= \left( \frac{1}{\sqrt{2\pi}} \right)^n \exp \left[ -\frac{1}{2} \sum (x_i - \mu)^2 \right] \\
&= \frac{1}{(2\pi)^{n/2}} \exp \left[ -\frac{1}{2} \sum (x_i^2 - 2x_i\mu + \mu^2) \right] \\
&= \frac{1}{(2\pi)^{n/2}} \exp \left[ -\frac{1}{2} \sum (-2x_i\mu + \mu^2) \right] \cdot \exp \left[ -\frac{1}{2} \sum x_i^2 \right] \\
&= \frac{1}{(2\pi)^{n/2}} \exp \left[ -\frac{1}{2} \left( -2\mu \sum x_i + n\mu^2 \right) \right] \cdot \exp \left[ -\frac{1}{2} \sum x_i^2 \right] \\
&= \frac{1}{(2\pi)^{n/2}} \underbrace{\exp \left[ \mu \sum x_i - \frac{1}{2} n\mu^2 \right]}_{\text{Define } S := \sum x_i} \exp \left[ -\frac{1}{2} \sum x_i^2 \right]
\end{aligned}$$

$$g(S | \mu) = \frac{1}{(2\pi)^{n/2}} \exp \left[ \mu S - \frac{n\mu^2}{2} \right]$$

$$h(\vec{x}) = \exp \left[ -\frac{1}{2} \sum x_i^2 \right]$$

By Factorization Theorem,  $S$  is a sufficient statistic for  $\mu$ .

**Example:** (9.49)

Suppose  $X_i \sim \text{Unif}(0, \theta)$ ,  $\theta$  unknown. Then,

$$f(x) = \frac{1}{\theta}, \quad x \in [0, \theta]$$

(In this, we use  $X_{(n)} = \max\{x_1, \dots, x_n\}$  as an estimator of  $\theta$ .)

$$f(x_1 | \theta) \times \cdots \times f(x_n | \theta) = \frac{1}{\theta^n} \text{ independent of any } x_i \text{'s.}$$

To bring the  $x_i$ 's into the story, we use the indicator function

$$I_{[0, \theta]}(x) = \begin{cases} 1 & x \in [0, \theta] \\ 0 & \text{else} \end{cases}$$

So,

$$f(x) = \frac{I_{[0, \theta]}(x)}{\theta}$$

With this,

$$\begin{aligned}
L(x_1, \dots, x_n | \theta) &= \frac{I_{[0, \theta]}(x_1)}{\theta} \times \cdots \times \frac{I_{[0, \theta]}(x_n)}{\theta} \\
&= \frac{1}{\theta^n} I_{[0, \theta]}(x_i \in [0, \theta], i = 1, 2, \dots, n)
\end{aligned}$$

Now,  $x_i \leq \theta$  for all  $i$  if and only if  $\max\{x_1, \dots, x_n\} \leq 0$ . So,

$$L(x_1, \dots, x_n \mid \theta) = \frac{1}{\theta^n} I_{[0, \theta]}(\max\{x_1, \dots, x_n \leq \theta\})$$

Define statistic  $T := \max\{x_1, \dots, x_n\}$  and then

$$g(T \mid \theta) = \frac{1}{\theta^n} I_{[0, \theta]}(T)$$

$$h(x_1, \dots, x_n) = 1$$

By the Factorization Theorem,  $T = \max(x_1, \dots, x_n)$  is a sufficient statistic.

**Note:**  $\bar{X}$  is not sufficient since there is no way to bring  $\bar{X}$  into  $L$ .

(Of course, no reason we should be able to do so since  $\bar{X}$  is also biased for  $\theta$ ... and clearly not consistent.  $\bar{X} \rightarrow \frac{\theta}{2}$  and these  $\bar{X} \rightarrow \theta$  in probability)

## 9.5 Rao-Blackwell Theorem and Minimum-Variance Unbiased Estimation

**Theorem: The Rao-Blackwell Theorem**

Let  $X$  and  $Y$  be a random variable such that  $E(Y) = \mu$  and  $\text{Var}(Y) = \sigma^2 y$ .

Let  $\phi(x) = E(Y | X)$ .

Then  $E(\phi(x)) = \mu$  and  $\sigma_{\phi(x)}^2 \leq \sigma^2 y$

*Proof.* Have  $f(x, y)$ ,  $f_1(x)$ ,  $f_2(y)$ , and

$$h(y | x) = \frac{f(x, y)}{f_1(x)} *$$

$$\begin{aligned} \phi(x) &= E(\phi(x)) \\ &= \int_{\mathbb{R}} y h(y | x) dy. \\ &= \int_{\mathbb{R}} y \frac{f(x, y)}{f_1(x)} dy \\ &= \frac{1}{f_1(x)} \int_{\mathbb{R}} y f(x, y) dy \end{aligned}$$

Then

$$f_1(x)\phi(x) = \int_{\mathbb{R}} y f(x, y) dy$$

$$\begin{aligned} E(\phi(x)) &= \int_{\mathbb{R}} \phi(x) f_1(x) dx \\ &= \int_{\mathbb{R}} \left[ \int_{\mathbb{R}} y f(x, y) dy \right] dx \\ &= \int_{\mathbb{R}} \int_{\mathbb{R}} f(x, y) dy dx \\ &= \int_{\mathbb{R}} y \left[ \int_{\mathbb{R}} f(x, y) dx \right] dy \\ &= \int_{\mathbb{R}} y f_2(y) dy \\ &= E(Y) \\ &= \mu. \end{aligned}$$

Note that  $\sigma_{\phi(x)}^2 = \text{E}((\phi(x) - \mu)^2)$ .

$$\begin{aligned}
\sigma_y^2 &= \text{E}((y - \mu)^2) \\
&= \text{E}((y - \phi(x))(\phi(x) - \mu)^2) \\
&= \text{E}((y - \phi(x))^2 + 2(y - \phi(x))(\phi(x) - \mu) + (\phi(x) - \mu)^2) \\
&= \underbrace{\text{E}((y - \phi(x))^2)}_{\geq 0} + \underbrace{2\text{E}((y - \phi(x))(\phi(x) - \mu))}_{*} + \underbrace{\text{E}((\phi(x) - \mu)^2)}_{\sigma_{\phi(x)}^2} \\
(*) &= \int_{\mathbb{R}} \int_{\mathbb{R}} (y - \phi(x))(\phi(x) - \mu) f(x, y) dy dx \\
&= \int_{\mathbb{R}} (\phi(x) - \mu) \left[ \int_{\mathbb{R}} (y - \phi(x)) f(x, y) dy \right] dx \\
(*) &= \int_{\mathbb{R}} (\phi(x) - \mu) \left[ \int_{\mathbb{R}} (y - \phi(x)) f_1(x) h(y | x) dy \right] dx \\
&= \int_{\mathbb{R}} (\phi(x) - \mu) f_1(x) \left[ \underbrace{\int_{\mathbb{R}} (y - \phi(x)) h(y | x) dy}_{\text{note}**} \right] dx \\
&= 0
\end{aligned}$$

**Note\*\*:** This is zero as

$$\phi = \text{E}(y | x) = \int y h(y | x) dy$$

So  $\sigma_y^2 = \text{E}((y - \phi(x))^2) + \sigma_{\phi(x)}^2 \geq \sigma_{\phi(x)}^2$

Aside: This is almost always a strict inequality. For equality,

$$\Pr(Y - \phi = 0) = 0$$

□

**Discussion:** Using the RBT (Rao-Blackwell Theorem) to construct “best” statistics.

Big Idea: We can take any unbiased estimator of  $\theta$  and combine with a sufficient stat for  $\theta$  to get a “better” estimator.

Let  $x_1, \dots, x_n$  denote a random sample of  $f(x | \theta)$ . We have  $T = \tau(x_1, \dots, x_n)$  sufficient for  $\theta$ . Let  $U = u(x_1, \dots, x_n)$  be an unbiased stat for  $\theta$  which is not itself a function of  $T$ . Consider  $\text{E}(u | T)$ :

Since  $T$  is sufficient, the conditional probability of  $U$  given  $T = \tau$  does not depend upon  $\theta$ .

Define a new stat

$$\phi(\tau) = \text{E}(U | T)$$

which is a function of  $\tau$  alone. Hence,  $\phi(T)$  is a statistic (does not depend on  $\theta$ ). (i.e. only a function of the data.)

By RBT,  $\phi(T)$  is an unbiased statistic for  $\theta$  with the guarantee that

$$\sigma_{\phi(T)}^2 < \sigma_u^2$$

Theorem: Rao-Blackwell Theorem (Textbook version):

Let  $\hat{\theta}$  be an unbiased estimator for  $\theta$  such that  $\text{Var}(\hat{\theta}) < \infty$ .

If  $T$  is a sufficient statistic for  $\theta$ , define  $\hat{\theta}^* = E(\hat{\theta}^* | T)$ .

Then, for all  $\theta$ ,  $E(\hat{\theta}^*) = \theta$  and

$$\text{Var}(\hat{\theta}^*) \leq \text{Var}(\hat{\theta})$$

**Discussion:** Minimum Variance Unbiased Estimator (MVUE)

If  $\hat{\theta}$  is an MVUE, then

$$\text{Var}(\hat{\theta}) \leq \text{Var}(\hat{\Psi}) \quad \text{for any other estimator } \hat{\Psi}.$$

Typically (almost always) the process we described above lends to an MVUE.

①  $\hat{\theta}$  unbiased estimator for  $\theta$

②  $T$  sufficient for  $\theta$

③  $E(\hat{\theta}^* | T)$  is a MVUE

Remark: All of yesterdays sufficiency examples indicate that all traditional estimators are MVUE.

**Example:**  $\text{Bern}(p)$

Here unbiased estimator  $\hat{p} = \frac{S}{n}$ .

Here  $S = \sum Y_i$  is sufficient.

$$E(\hat{p} | S) = \frac{S}{n} \text{ sufficient stat, has no } p.$$

Rao-Blackwell  $\implies \hat{p}$  is an MVUE.

Also note,  $\text{Var}(\hat{p}) < \text{Var}(S)$ , is much better  $\frac{p}{n} < np$ .

**Example:** Last day we showed that  $S = \sum Y_i$  is a sufficient stat for  $\lambda$  in  $\text{Pois}(\lambda)$ .

$$f(y) = \frac{\lambda^y e^\lambda}{y!}, \quad y = 0, 1, 2, \dots$$

We seek an MVUE of  $e^{-\lambda}$ . Need an unbiased estimator of  $Y_1, \dots, Y_n$  iid.

a) Show that

$$W = \begin{cases} 1 & Y_1 = 0 \\ 0 & Y_1 \neq 0 \end{cases}$$

is unbiased for  $e^{-\lambda}$ . (New idea: Not just looking for  $\lambda$ , but a function of  $\lambda$ .)

**Solution:**  $E(W) = 1 \cdot \Pr(Y_1 = 0) + 0 \cdot \Pr(Y_1 \neq 0) = \Pr(Y_1 = 0) = \frac{\lambda^0 e^{-\lambda}}{0!} = e^{-\lambda}$

b) Use RBT to find a MVUE. Compute  $E(W | S)$

$$\begin{aligned} E(W | S) &= 1 \cdot \Pr(Y_1 = 0 | Y_1 + Y_2 + \dots + Y_n = S) + 0 \cdot \Pr(Y_1 \neq 0 | S) \\ &= \frac{\Pr(Y_1 = 0 \text{ and } Y_1 + Y_2 + \dots + Y_n = S)}{\Pr(Y_1 + Y_2 + \dots + Y_n = S)} \\ &= \frac{\Pr(Y_1 = 0 \text{ and } Y_2 + \dots + Y_n = S)}{\Pr(Y_1 + Y_2 + \dots + Y_n = S)} \quad \text{top now independent events} \\ &= \frac{\Pr(Y_1 = 0) \cdot \Pr(Y_2 + \dots + Y_n = S)}{\Pr(Y_1 + Y_2 + \dots + Y_n = S)} \\ &= \frac{e^{-\lambda} \cdot \Pr(Y_2 + \dots + Y_n = s)}{\Pr(Y_1 + Y_2 + \dots + Y_n = s)} \end{aligned}$$

**325 FACT:** How is  $Y_1 + \dots + Y_n$  distributed?  $Y_i \sim \text{Pois}(\lambda)$ .

Via moment generating functions, we proved that  $Y_1 + \dots + Y_n \sim \text{Pois}(n\lambda)$

Hence  $Y_2 + \dots + Y_n \sim \text{Pois}((n-1)\lambda)$ , so

$$\begin{aligned} E(W | S) &= \frac{e^{-\lambda} \left( \frac{((n-1)\lambda)^S e^{-(n-1)\lambda}}{S!} \right)}{\frac{(n\lambda)^S e^{-n\lambda}}{S!}} \\ &= e^{-\lambda} \left( \frac{(n-1)\lambda}{n\lambda} \right)^S \cdot \frac{e^{-(n-1)\lambda}}{e^{-n\lambda}} \\ &= e^{-\lambda} \cdot \left( \frac{n-1}{n} \right)^S e^\lambda \\ &= \left( \frac{n-1}{n} \right)^S \end{aligned}$$

By RBT,  $\phi(S) = \left( \frac{n-1}{n} \right)^S$  is an MVUE of  $e^{-\lambda}$ ,  $S = \sum Y_i$ .

Aside:

$$e^{-\lambda} = \lim_{n \rightarrow \infty} \left( 1 - \frac{1}{n} \right)^n$$

$$\phi(S) = \lim_{n \rightarrow \infty} \left( 1 - \frac{1}{n} \right)^S$$

## 9.6 The Method of Moments

This is our oldest estimation technique. If there are  $k$  parameters that have to be estimated, set the  $1^{\text{st}}$   $k$  population moments (given in terms of the parameters) and set them equal to the sample moments.

Recall:

$$\underbrace{\mu'_k = \mathbb{E}(Y^k)}_{\text{population}} \quad \text{and} \quad \underbrace{m'_k = \frac{1}{n} \sum Y_i^k}_{\text{data set}}$$

**Example:** Let  $Y \sim N(\mu, \sigma^2)$ . Then

$$\begin{aligned}\mu'_1 &= \mathbb{E}(Y) = \mu \\ \mu'_2 &= \mathbb{E}(Y^2) = \sigma^2 + \mu^2\end{aligned}$$

$$\begin{aligned}m'_1 &= \frac{1}{n} \sum Y_i = \bar{Y} \\ m'_2 &= \frac{1}{n} \sum Y_i^2\end{aligned}$$

Set equal and solve for parameters, then

$$\mu = \bar{Y}$$

$$\sigma^2 + \mu^2 = \frac{1}{n} \sum Y_i^2$$

Want  $\sigma^2$  in terms of sample moments.

$$\sigma^2 + \bar{Y}^2 = \frac{1}{n} \sum Y_i^2$$

$$\sigma^2 = \frac{1}{n} \sum Y_i^2 - (\bar{Y})^2$$

Yields 2 estimators:

$$\begin{aligned}\hat{\theta} &= \bar{Y} \quad \text{for } \mu \\ \hat{\Psi} &= \frac{1}{n} \sum Y_i^2 - (\bar{Y})^2 \quad \text{for } \sigma^2\end{aligned}$$

Note:  $\hat{\Psi}$  is our “old” biased estimator for  $\sigma^2$ . While both of these are consistent stats,  $\hat{\Psi}$  is biased as  $E(\hat{\Psi}) = \frac{n-1}{n} \sigma^2$ .

**Example:** Suppose  $Y_1, \dots, Y_n \sim \text{Gamma}(\alpha, \beta)$ . Use MoM to find estimators for  $\alpha$  and  $\beta$ .

Recall for Gamma,  $E(Y) = \alpha\beta$ ,  $\text{Var}(Y) = \alpha\beta^2$ . So,

$$\begin{aligned}
\mu'_1 &= \alpha\beta \\
m'_1 &= \bar{Y} \\
\mu'_2 &= \frac{1}{n} \sum Y_i^2 \\
&= E(Y^2) \\
&= \text{Var}(Y) + (E(Y))^2 \\
&= \alpha\beta^2 + (\alpha\beta)^2 \\
&= \beta^2(\alpha + \alpha^2) \\
m'_2 &= \frac{1}{n} \sum Y_i^2
\end{aligned}$$

**System:**  $\alpha\beta = \bar{Y}$      $\beta^2(\alpha + \alpha^2) = m'_2$      $\alpha = \frac{\bar{Y}}{\beta}$     Then,

$$\begin{aligned}
m'_2 &= \beta^2 \left( \frac{\bar{Y}}{\beta} + \frac{\bar{Y}^2}{\beta^2} \right) \\
&= \bar{Y}^2 + \beta\bar{Y} \\
\iff \beta\bar{Y} &= m'_2 - \bar{Y}^2 \\
\iff \beta &= \frac{m'_2 - \bar{Y}^2}{\bar{Y}} \\
\implies \alpha &= \frac{\bar{Y}^2}{m'_2 - \bar{Y}^2}
\end{aligned}$$

Define  $\alpha$ 's estimator  $\hat{\theta} := \frac{\bar{Y}^2}{m'_2 - \bar{Y}^2}$ .

and  $\beta$ 's estimator  $\hat{\Psi} := \frac{m'_2 - \bar{Y}^2}{\bar{Y}}$ .

The following is the last example using MoM:

**Example:**  $Y_1, \dots, Y_n \sim \text{Unif}(0, \theta)$ , with  $\theta$  unknown.

$$\mu'_1 = E(Y) = \frac{\theta}{2} \quad m'_1 = \frac{1}{n} \sum Y_i = \bar{Y}$$

Mom set equal solution for population parameters.

$$\frac{\theta}{2} = \bar{Y} \implies \theta = 2\bar{Y}$$

Define  $\hat{\theta} := 2\bar{Y}$ .

So this is another estimator for  $\theta$ . Is it any good?

Well, it is unbiased since  $E(2\bar{Y}) = 2 \cdot \frac{\theta}{2} = \theta$  and it is consistent.

Use variance trick to show...

$$\text{Var}(2\bar{Y}) = 4 \text{Var}(\bar{Y}) = 4 \frac{\text{Var}(Y)}{n} = 4 \cdot \frac{\theta^2/12}{n} = \frac{\theta^2}{3n}$$

Note:

$$\lim_{n \rightarrow \infty} \text{Var}(2\bar{Y}) = \lim_{n \rightarrow \infty} \frac{\theta^2}{3n} = 0$$

Consistent by theorem. However, earlier we showed  $\hat{\Psi} = \max\{Y_1, \dots, Y_n\}$  is sufficient. Consider  $\text{Var}(\hat{\Psi})$ .

Like in Homework,  $F(y) = \int_0^y \frac{1}{\theta} dy = \frac{y}{\theta}$

$$\Pr(Y_{(n)} \leq y) = \left(\frac{y}{\theta}\right)^n = \frac{y^n}{\theta^n}.$$

And  $f_{(n)}(y) = \frac{ny^{n-1}}{\theta^n}$ , thus

$$E(\hat{\Psi}) = \int_0^\theta y \cdot f_{(n)}(y) dy = \frac{n}{n+1} \theta.$$

$\hat{\Psi}$  is biased but it is consistent since  $\lim_{n \rightarrow \infty} E[\hat{\Psi}] \rightarrow \theta$ . For  $\text{Var}(\hat{\Psi}) = E(\hat{\Psi}^2) - E(\hat{\Psi})^2$ :

$$E(\hat{\Psi}^2) = \int_0^\theta y^2 \cdot \frac{ny^{n-1}}{\theta^n} dy = \frac{n}{n+2} \theta^2$$

So...

$$\text{Var}(\hat{\Psi}) = \frac{n}{n+2} \theta^2 - \left(\frac{n}{n+1} \theta\right)^2 = \frac{n\theta^2}{(n+1)^2(n+2)}$$

And this variance is much better than  $\text{Var}(\hat{\theta})$ .

$$\begin{aligned} \text{bias}(\hat{\theta}) &= 0 \\ \text{bias}(\hat{\Psi}) &= \frac{n}{n+1} \theta - \theta = -\frac{\theta}{n+1} \\ \text{MSE}(\hat{\theta}) &= \text{Var}(\hat{\theta}) + \text{bias}(\hat{\theta})^2 = \frac{\theta^2}{3n} \\ \text{MSE}(\hat{\Psi}) &= \frac{n\theta^2}{(n-2)(n+1)^2} + \frac{\theta^2}{(n+1)^2} \end{aligned}$$

We see both  $\text{Var}(\hat{\Psi}) < \text{Var}(\hat{\theta})$  and  $\text{MSE}(\hat{\Psi}) < \text{MSE}(\hat{\theta})$

Of course, if we were to define

$$\tilde{\Psi} := \frac{n+1}{n} \hat{\Psi}$$

Now we have an unbiased consistent statistic. Last day, we showed that  $\hat{\Psi}$  is sufficient. So by the Rao-Blackwell Theorem,  $\tilde{\Psi}$  will be an MVUE. That is,  $\text{Var}(\tilde{\Psi}) \leq \text{Var}(\hat{\Psi})$ .

## 9.7 The Method of Maximum Likelihood

Recall, the likelihood function is simply the joint pdf with the interpretation that a parameter(s) is unknown.

$$Y_1 \sim \text{via } f(y | \theta)$$

$Y_i$ 's are iid random sample.

If  $\theta$  is known, the joint pdf is

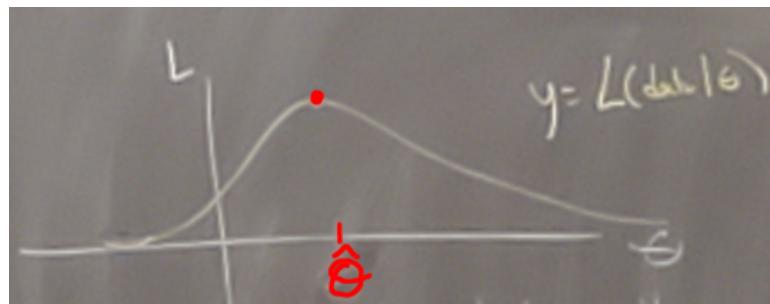
$$\begin{aligned} f(y_1, \dots, y_n) &= \underbrace{f_1(y_1) \cdot f_2(y_2) \cdots f_n(y_n)}_{\substack{\text{domain } f \subseteq \mathbb{R}^n \\ \text{product of marginals}}} && \text{by independence} \\ &= f(y_1) \cdot f(y_2) \cdots f(y_n) && \text{by i.i.d.} \\ &= \prod_{i=1}^n f(y_i) \end{aligned}$$

If  $\theta$  unknown, we emphasize this via the **likelihood function**.

$$L(Y_1, \dots, Y_n | \theta) = \underbrace{f(y_1 | \theta) \cdot f(y_2 | \theta) \times \cdots \times f(y_n | \theta)}_{\substack{\text{domain } L \subseteq \mathbb{R}^{n+1} \\ \text{really the same}}}$$

**Discussion:** The reason this is called a “likelihood” function. Given a set of observations  $Y_i$ , what is the most likely value of  $\theta$ .

Fix  $(Y_1, \dots, Y_n)$ ,  $L$  at this number is a 1 variable function in terms of  $\theta$ .



Given this data set, what is the most likely value that  $\theta$  takes? That would be  $\theta$  corresponding to the largest value of pdf  $\prod_i^n f(y_i | \theta)$ .

How do we find the largest? (Optimize with respect to  $\theta$ ). That is, set  $\frac{\partial L}{\partial \theta} = 0$  solve and verify.

**Definition:** The solution to  $L_\theta = 0$  defines an estimator  $\hat{\theta}$  of  $\theta$ , called the maximum likelihood estimator (MLE).

*Remark:* In computations, verify the critical point *is* a maximum.

**Example:**  $X \sim \text{Bern}(p) = \text{Binom}(1, p)$

Given  $X_1, \dots, X_n$  an iid sample.

$$\begin{aligned} L(X_1, \dots, X_n \mid p) &= p^{X_1}(1-p)^{1-X_1} \cdots p^{X_n}(1-p)^{1-X_n} \\ &= p^S(1-p)^{n-S}, \quad S := \sum X_i \end{aligned}$$

Then,

$$\begin{aligned} \frac{\partial L}{\partial p} &= Sp^{S-1}(1-p)^{n-S} + p^S(n-s)(1-p)^{n-S-1}(-1) \\ &= p^{S-1}(1-p)^{n-S-1}[S(1-p) - p(n-S)] \end{aligned}$$

And  $\frac{\partial L}{\partial p} = 0$  when

$$\begin{aligned} S - Sp - pn + Sp &= 0 \\ S - pn &= 0 \\ p = \frac{S}{n} &= \bar{X} \end{aligned}$$

Is this a max? Yes, use the 1<sup>st</sup> derivative test:  $\frac{\partial L}{\partial p} = p^{S-1}(1-p)^{n-S-1}[S - pn]$ .

If  $p < \frac{S}{n}$ , then  $\frac{\partial L}{\partial p} > 0$ .

If  $p > \frac{S}{n}$ , then  $\frac{\partial L}{\partial p} < 0$

Therefore  $\bar{X}$  is an MLE for  $p$ .

**Discussion:** The log-likelihood function.

$$L(\vec{X} \mid \theta) = \prod_{i=1}^n f(X_i \mid \theta)$$

Always an  $n$ -product. Depending on  $f$ , computing  $\frac{\partial L}{\partial \theta}$  can be difficult.

Since  $f(X_i \mid \theta) > 0$  and we have a function whose existence is to turn products into sums...

**Definition:** The log-likelihood function:

$$\ln L(\vec{X} \mid \theta)$$

Claim: The MLE of  $L(\vec{X} \mid \theta)$  also maximizes  $\ln L(\vec{X} \mid \theta)$

Reason:

$$\frac{\partial}{\partial \theta} (\ln L(\vec{X} \mid \theta)) = \frac{\frac{\partial L}{\partial \theta}(\vec{X} \mid \theta)}{L(\vec{X} \mid \theta)} = 0 \iff \frac{\partial}{\partial \theta} (L(\vec{X} \mid \theta)) = 0$$

**Example:** The last one again...

$$L(X_1, \dots, X_n | p) = p^S(1-p)^{n-S}, \quad S := \sum X_i$$

$$\ln L(\vec{x} | p) = S \ln p + (n - S) \ln(1 - p)$$

$$\frac{\partial}{\partial p} (\ln L) = \frac{S}{p} + (n - S) \cdot \frac{-1}{1 - p}$$

$$\frac{\partial}{\partial p} (\ln L) = 0 \iff S(1 - p) + (n - S)(-p) = 0$$

$$\iff S - Sp - np + Sp = 0$$

$$\iff p = \frac{S}{n}$$

**Topic:** More population parameters.

This idea scales. If  $f$  depends upon  $k$  unknowns  $\theta_1, \dots, \theta_k$ , then define

$$L(X_1, \dots, X_n | \theta_1, \dots, \theta_k) = \prod_{i=1}^n f(x_i | \theta_1, \dots, \theta_k)$$

Optimizing this is a calculus problem.

*Remark:* We could have also extended the idea of the Factorization Theorem the same way with sufficient statistics to get

$$L(x_1, \dots, x_n | \theta_1, \dots, \theta_k) = g(S_1, \dots, S_k, \theta_1, \dots, \theta_k) \cdot h(x_1, \dots, x_n)$$

**Example:** Let  $X_1, \dots, X_n$  be iid from  $N(\mu, \sigma^2)$  with both parameters unknown. Find the MLE for  $\mu$  and  $\sigma^2$ .

For ease of computation,  $N(\mu, \theta)$  such that  $\theta := \sigma^2$ .

$$L(x_1, \dots, x_n | \mu_1 = \theta) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\theta}} \exp \left[ -\frac{(X_i - \mu)^2}{2\theta} \right]$$

Fixing  $\vec{x}$ , we consider  $L_{\vec{x}}(\mu, \theta)$ . . . I have no intentions of  $\nabla L_{\vec{x}}$  under a product sign...

$$\begin{aligned} \ln L &= \sum_{i=1}^n \left( \ln \left( \frac{1}{\sqrt{2\pi\theta}} \right) + \frac{-(X_i - \mu)^2}{2\theta} \right) \\ &= -\frac{1}{2} \ln(2\pi\theta)n - \frac{1}{2\theta} \sum_{i=1}^n (X_i - \mu)^2 \end{aligned}$$

$$\begin{aligned} \frac{\partial}{\partial \mu} (\ln L) &= 0 - \frac{1}{2\theta} \sum_{i=1}^n 2(X_i - \mu)(-1) \\ &= \frac{1}{\theta} \sum_{i=1}^n (X_i - \mu) \\ &= 0 \quad \text{when} \quad \sum X_i - n\mu = 0 \rightarrow \hat{\mu} = \bar{X} \end{aligned}$$

$$\text{i.e. } \mu = \frac{1}{n} \sum X_i = \bar{X}$$

$$\begin{aligned}\frac{\partial}{\partial \theta} (\ln L) &= -\frac{n}{2} \cdot \frac{2\pi}{2\pi\theta} + \frac{1}{2\theta^2} \sum_{i=1}^n (X_i - \mu)^2 \\ &= -\frac{n}{2\theta} + \frac{1}{2\theta^2} \sum_{i=1}^n (X_i - \mu)^2 \\ &= 0 \quad \text{At our potential critical point } \mu = \bar{X} \\ \iff -n\theta + \sum_{i=1}^n (X_i - \bar{X})^2 &= 0 \\ \text{or } \hat{\theta} &= \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2 \\ &\text{our old biased estimator for } \sigma^2\end{aligned}$$

Is this a max? We need to construct the Hessian Matrix...

$$\begin{aligned}H &= \begin{bmatrix} \ln L_{\mu\mu} & \ln L_{\mu\theta} \\ \ln L_{\theta\mu} & \ln L_{\theta\theta} \end{bmatrix} \\ &= \begin{bmatrix} -\frac{n}{\bar{\theta}} & -\frac{\sum X_i - n\mu}{\hat{\theta}^2} \\ -\frac{\sum X_i - n\mu}{\hat{\theta}^2} & \frac{n}{2\hat{\theta}^2} - \frac{1}{\hat{\theta}^3} \sum (X_i - \mu)^2 \end{bmatrix}\end{aligned}$$

At  $(\hat{\mu}, \hat{\theta})$

$$H(\mu, \theta) = \begin{bmatrix} -\frac{n}{\bar{\theta}} & 0 \\ 0 & \frac{n}{2\hat{\theta}^2} - \frac{1}{\hat{\theta}^3} \sum (X_i - \mu)^2 \end{bmatrix}$$

For a maximum, by the 2<sup>nd</sup> derivitive test, we need  $(\ln L)_{\mu\mu} = -\frac{n}{\bar{\theta}} < 0$ , which it is! And  $\det H > 0$ .

$$\text{Note, } (\ln L)_{\theta\theta} = \frac{\hat{\theta}n - 2\sum(X_i - \bar{X})^2}{2\hat{\theta}^3}.$$

$$\text{But } n\hat{\theta} = \sum(X_i - \bar{X})^2, \text{ so } (\ln L)_{\theta\theta} = \frac{\hat{\theta}n - 2\hat{\theta}n}{2\hat{\theta}^3} = -\frac{n}{2\hat{\theta}^2} < 0.$$

So  $\det H = \frac{-n}{\bar{\theta}} \cdot \frac{-n}{2\hat{\theta}^2} = \frac{n^2}{2\hat{\theta}^3} > 0$ . Hence  $(\bar{X}, \hat{\theta})$  is the location of a MLE by the 2nd derivitive test.

**Discussion:** Why MLEs? Because they have nice properties.

① If  $U$  is a sufficient statistic of  $\theta$ , then the MLE is a function of  $U$ .

Reason: If  $U$  is sufficient, then  $L(\vec{x} \mid \theta)$  is factorable:

$$L(x_1, \dots, x_n \mid \theta_1, \dots, \theta_k) = g(U, \theta) \cdot h(x_1, \dots, x_n)$$

This  $\frac{\partial L}{\partial \theta} = \frac{\partial g}{\partial \theta} \cdot h(\vec{x})$

$$\text{i.e. } \underbrace{\frac{\partial L}{\partial \theta}}_{\substack{\text{sol'n to this is} \\ \hat{\theta} \text{ (the MLE)}}} = 0 \iff \underbrace{\frac{\partial g}{\partial \theta}}_{\substack{\text{sol'n to this is} \\ \hat{\theta} \text{ a function of } U}} = 0$$

## ② The invariance properties of the MLE.

Suppose we want to estimate a function of the parameter (say  $t(\theta)$ ). If  $t$  is a one-to-one (injective) function (i.e. invertible), the MLE of  $t(\theta)$  will be simply  $t(\hat{\theta})$ , where  $\hat{\theta}$  is the MLE of  $\theta$ .

Reason: Let  $t(\theta)$  be injective. Assume  $\hat{\theta}$  the MLE of  $\theta$   $t(\theta)$  invertible  $\implies \Psi = t(\hat{\theta})$  and  $t^{-1}(\Psi) = \theta$

Then  $L(\vec{x} | \theta)$  is maximized at  $\theta = \hat{\theta}$ . Then  $L(\vec{x} | t^{-1}(\Psi))$  is maximized at the same  $\hat{\theta}$ . Hence  $t^{-1}(\Psi) = \hat{\theta}$  or  $\hat{\Psi} = t(\hat{\theta})$ .

**Example:** (#1) Show  $\bar{X}$  is MLE for  $\lambda$  when  $X \sim \text{Pois}(\lambda)$ .

In showing  $S := \sum X_i$  is sufficient, we had factored the likelihood function

$$L(\vec{x} | \lambda) = \underbrace{e^{-n\lambda} \lambda^S}_{g(S, \lambda)} \cdot \frac{1}{x_1! \cdot x_2! \cdots x_n!}$$

To find the MLE  $\frac{\partial L}{\partial \lambda} = 0$  when  $\frac{\partial g}{\partial \lambda} = 0$ .

$$\begin{aligned} \frac{\partial g}{\partial \lambda} &= -ne^{n\lambda} \lambda^S + e^{-n\lambda} s \lambda^{S-1} \\ &= e^{-n\lambda} \lambda^{S-1} [-n\lambda + S] \\ &= 0 \\ \implies \lambda &= \frac{S}{n} \quad \text{i.e. } \hat{\lambda} = \bar{x} \end{aligned}$$

As before, can show a max by the 1st derivative test.

**Example:** Lots of bits together

We have shown that  $\chi^2 \sim \text{Pois}(\lambda)$  with  $S = \sum x_i$  is a sufficient statistic for  $\lambda$ . And we just showed that  $\hat{\theta} = \bar{x}$  is MLE for  $\lambda$ .

Hence, by the Rao-Blackwell Theorem,  $E(\bar{X} | S) = \bar{X}$ , then  $\bar{X}$  is an MVUE for  $\lambda$ . At the end of RBT, via

$$w = \begin{cases} 1 & x_1 = 0 \\ 0 & x_1 \neq 0 \end{cases} \quad (w \text{ unbiased for } e^{-\lambda})$$

Combining  $w$  with  $S$  we get the MVUE  $\widehat{\Psi} = \left(\frac{n-1}{n}\right)^S = \left(\frac{n-1}{n}\right)^{n\bar{X}}$  for  $e^{-\lambda}$ . (also  $\text{Var}(\widehat{\Psi}) < \text{Var}(w)$ ) by RBT.

Now, since  $\bar{X}$  is MLE for  $\lambda$ , by the invariance properties of MLE,

$$\theta^* = e^{-\bar{X}} \text{ is MLE for } e^{-\lambda}.$$

Reason:  $g(x) = e^{-x}$  is clearly injective on  $\mathbb{R}$ .

Moreover, by RBT  $\text{Var}(\widehat{\Psi}) < \text{Var}(e^{-\bar{X}})$ . And we have a conjecture that  $e^{-\bar{X}}$  is actually biased.

While  $e^{-\bar{X}}$  is the “natural” estimator for  $e^{-\lambda}$ , the RBT says  $\widehat{\Psi} = \left(\frac{n-1}{n}\right)^{n\bar{X}}$  is “better” as  $\widehat{\Psi}$  is unbiased and  $\text{Var}(\widehat{\Psi}) < \text{Var}(\theta^*)$ .

**Q: Can we show this directly?**

Maybe. If  $\theta^*$  is unbiased...

$$\text{Var}(\widehat{\Psi}) < \text{Var}(\theta^*) \iff E[(\widehat{\Psi})^2] < E[(\theta^*)^2]$$

FACT:  $\lim_{n \rightarrow \infty} \left(\frac{n-1}{n}\right)^n = \lim_{n \rightarrow \infty} \left(1 - \frac{1}{n}\right)^n = e^{-1}$ . (This is via Calc I application of L'Hopitals)

Moreover, via the same ln technique, it can be shown that  $f(x) = \left(\frac{x-1}{x}\right)^x$  is increasing on  $x \geq 2$  (see addendum). Hence,

$$\begin{aligned} \left(\frac{n-1}{n}\right)^n &< e^{-1} \\ \implies \left(\frac{n-1}{n}\right)^{2n} &< e^{-2} \\ \implies \left(\frac{n-1}{n}\right)^{2n\bar{X}} &< e^{-2\bar{X}} \\ \implies \left(\left(\frac{n-1}{n}\right)^{n\bar{X}}\right)^2 &< e^{-\bar{X}^2} \\ \implies E(\widehat{\Psi}^2) &< E((\theta^*)^2) \\ \implies \text{Var}(\widehat{\Psi}) &< \text{Var}(\theta^*) \\ &\quad (\text{if } \theta^* \text{ is unbiased}) \end{aligned}$$

**Addendum:**

Claim:  $f(x) = \left(\frac{x-1}{x}\right)^x$  is increasing on  $x \geq 2$ .

Reason:

$$\begin{aligned}\ln f &= x \ln(x-1) - x \ln x \\ \frac{f'}{f} &= \ln(x-1) + \frac{x}{x-1} - \ln x - 1 \\ &= \ln\left(\frac{x-1}{x}\right) + \frac{1}{x-1} \\ f'(x) = 0 &\implies \ln\left(\frac{x-1}{x}\right) + \frac{1}{x-1} = 0\end{aligned}$$

But this does not happen on  $(1, \infty)$ . Equivalently,  $\ln\left(\frac{x-1}{x}\right)^{1-x} = 1 \implies \frac{x-1}{x} = 1$  has no solutions.

Since  $f'(2) = \ln\left(\frac{1}{2}\right) + 1 = 1 - \ln 2 > 0$ , we get  $f'(x) > 0$  on  $(1, \infty)$ .

# Chapter 10

## Hypothesis Testing

### Motivational example:

I claim that I am a 75% free-throw shooter. You need to verify, and make me shoot 20 free throws. I make 8.

You are suspicious of my claim.

Is that suspicion founded?

We can think of a free-throw as  $\text{Bern}(0.75)$ .

The probability of getting 8 is binomial.

$$\Pr(X = 8) = \binom{20}{8} \left(\frac{3}{4}\right)^8 \left(\frac{1}{4}\right)^{12} = 0.00075$$

This implies that if  $p = \frac{3}{4}$  is true, we observed a “rare event.”

But, another way to think about this is if  $p = \frac{3}{4}$ , what is the probability that I would shoot 8 or worse?

$$\Pr(X \leq 8) = \sum_{i=0}^8 \binom{20}{i} \left(\frac{3}{4}\right)^i \left(\frac{1}{4}\right)^{20-i} < 0.001$$

The entire array,  $0, 1, \dots, 8$  is less than  $\frac{1}{1000}$  probability. We are “safe” to conclude that the true free-throw percentage is less than 75%.

This is what hypothesis testing is all about.

- Assume something is true
- Observe and compute the probability of this outcome, or more extreme (rarer) events.
- Decide if there is evidence against your assumption.

Flip-side: In chapter 8 we would have constructed a confidence interval for the true  $p$ .

$n = 20$  is small... we need  $t(\text{df} = 19)$  distribution.

A 99% C.I. for  $p$

Uses  $t_{0.005}(19) = 2.861$  standard errors.

$$\hat{p} \pm t_{0.005} \sqrt{\frac{\hat{p}\hat{q}}{20}} \approx 0.4 \pm 0.256 \implies p \in (0.144, 0.656)$$

$p = \frac{3}{4}$  is very far outside. Again, very unlikely that  $p = \frac{3}{4}$ .

\*C.I. and the Hypothesis Test are 2-sides of the same coin.

Hypothesis Testing Vocab

① The Hypotheses.

- **The null hypothesis**  $H_0$  : Assumption of no change.

Assumed true unless “proven” otherwise.

**Example:**  $H_0 : p_0 = 0.75$

- **The alternative hypothesis**  $H_a$  : The claim we aim to test.

$H_a : p < 0.75$

**Example (One sided):**  $H_a : p < p_0$  or  $H_a : p > p_0$

**Example (Two sided):**  $H_a : p \neq p_0$

② **The decision rule.** This is the cutoff  $\alpha$  we use to accept or reject  $H_a$ . Usually stated beforehand.

③ **The decision.** After a probability calculation, either

- (i) Reject the null hypothesis  $H_0$  or
- (ii) Not enough evidence to reject  $H_0$  (accept  $H_a$ )

④ **p-value.** The probability of seeing as much, or more evidence for  $H_a$  than we saw in the data.

$$p = \Pr(X \leq 8) < 0.001$$

\*The smaller the  $p$ -value, the more support for  $H_a$ .

Test Results are called **statistically significant** if  $H_0$  is rejected.

### Topic: Error Types

2 possible decisions  $\implies$  2 possible mistakes.

The “truth”	$H_0$	$H_a$
Test Supports $H_0$	good	Type II
Test Supports $H_a$	Type I	good

**Type I:** We reject  $H_0$  when  $H_0$  is the truth.

\*This error comes from the decision rule.

**Example:** I actually am a 75% free-throw shooter, but had a bad day.

**Type II:** We accept  $H_0$  when  $H_a$  is the truth.

We will see that we can measure these error types in some sense, but it depends on our decision rule  $\alpha$  and the unknown, true population parameter  $\theta$ .

$$\Pr(\text{Type I error}) = \alpha.$$

$$\Pr(\text{Type II error}) = \beta.$$

**Example:** Back to the free throws.

*Decision Rule:* Beforehand, you decide if I make 10 or less, you will reject  $H_0 : p = 0.75$ .

The text calls the set of outcomes the **rejection region** (RR). If  $T \in \text{RR}$  then we reject  $H_0$  (support the alternative).

Here,  $RR = \{0, 1, 2, \dots, 10\}$  and we can exactly calculate this:

$$\begin{aligned}\alpha &= \Pr(\text{Type I error}) \\ &= \Pr(\text{rejecting } H_0 \text{ when } H_0 \text{ is true}) \\ &= \Pr\left(0 \leq T \leq 10 \text{ and } p = \frac{3}{4}\right) \\ &= \text{pbinom}(10, 20, 0.75) \\ &= 0.013\end{aligned}$$

Hence, it is very unlikely to get a Type I error.

Computing Type II error probabilities requires a guess for the true  $H_a$ .

$$\begin{aligned}\beta &= \Pr(\text{Type II error}) \\ &= \Pr(\text{accept } H_0 \text{ when } H_a \text{ is true}) \\ &= \Pr\left(11 < T \leq 20 \text{ when } p = \theta < \frac{3}{4}\right) \\ &= 1 - \Pr\left(0 \leq T \leq 10 \text{ when } p = \theta < \frac{3}{4}\right) \\ \beta(\theta) &= 1 - \text{pbinom}(10, 20, \theta) \text{ a function of } \theta\end{aligned}$$

Here, if  $\theta = 0.6$ , then  $\underbrace{\beta = 0.755}_{\text{This is a lot, 75\%}}$

If  $\theta = 0.5$ ,  $\beta = 0.4119$

Note the larger the true difference between  $\theta$  and  $p_0 = \frac{3}{4}$  is, the smaller the Type II error.

FACT:  $\alpha$  and  $\beta$  are inversely related! If we increase  $\alpha$  (Probability of Type I error), we see a decrease in  $\beta$  (Probability of Type II error) (and vice versa).

**Definition:** The power function of a test is defined as

$$\text{power}(\theta) = 1 - \beta(\theta)$$

and this measures Type II error.

**Example:** New ultrasound machine: Claimed to detect tumors better than the old machine.

Hospital designs a test: Take a known patient with a known tumor distribution “tumor set.” Scan each patient in both machine: record the proportion of known tumors detected with each machine.

Let  $p_0$  be the proportion found of old machine, and  $p_1$  be the new.

- (a)  $H_0 : p_0 = p_1$       and       $H_a : p_0 < p_1$
- (b) A Type I error occurs when we decide the new machine is better, when it is actually worse.

*Real world consequences:* Results in investment in “better” equipment when is not better and does not help your patients

- (c) A Type II error occurs when we accept  $H_0$  when  $H_a$  is true. (We think that the old equipment is better when it isn’t.)

*Real world consequences:* We could have had better machines, detected more cancer, and saved more lives (but didn’t).

### 10.3 Z-tests (large samples)

**Example (10.18):** The hourly wages in a particular industry is distributed  $N(13.20, 2.50)$ . A company in this industry employs 40 workers, paying them an average of \$12.20 per hour. Can this company be accused of paying in substandard wages. Use  $\alpha = 0.01$ .

**Solution:** Recall  $n = 40 > 30$  is considered “large”, Hence

$$\begin{aligned}\bar{X} &\sim N\left(\mu, \frac{\sigma^2}{n}\right) \\ &= N\left(13.20, \frac{2.50}{40}\right) \\ &= N(13.20, 0.0625)\end{aligned}$$

Our test:

$$H_0 : \mu_c = 13.20 \quad H_a : \mu_c < 13.20$$

Where  $\mu_c$  is the company average. Here,  $\alpha = 0.01$  is the decision rule, also the significance level, and also the probability of making a Type I error (reject  $H_0$  when  $H_0$  is true).

$$\begin{aligned}\alpha &= \Pr(\text{Type I error}) \\ &= \Pr(\text{rejecting } H_0 \text{ when it is true}) \\ &= 0.01\end{aligned}$$

Compute the  $p$ -value.

$$p = \Pr(\bar{X} \leq 12.20 \mid \mu = 13.20)$$

Convert to a Z-score (like we did in 325).

$$\begin{aligned}&= \Pr\left(Z \leq \frac{12.20 - 13.20}{0.25}\right) \\ Z &= \frac{\bar{X} - \mu}{\sigma - \bar{X}} \quad \sigma_{\bar{X}} \sqrt{\frac{2.50}{40}} = \frac{1}{4}\end{aligned}$$

$$\begin{aligned}p &= \Pr(Z \leq -4) \\ &= \Pr(Z \geq 4) \\ &= 0.0000317 \text{ by table 4} \\ &= \text{pnorm}(12.20, 13.20, 0.25)\end{aligned}$$

Decision and conclusion: Since p-value is much less than  $0.01 = \alpha$ , we reject  $H_0$  (accept the alternative). Therefore we conclude that, yes, the company appears to be systematically underpaying its employees in relation to the rest of the industry.

### Remark:

①  $Z = \frac{\bar{X} - \mu}{\sigma - \bar{X}} = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}}$  is called the **test statistic** or **Z-statistic**. E.g.  $Z = -4$ .

② R command: `pnorm(0, mu, sigma)` and is a left tail calculator (opposite of Table 4).

**Example (10.21):** Shear strength measurements are derived from unconfined compression tests for two types of soils.

Soil 1	Soil 2
$n_1 = 30$	$n_2 = 35$
$\bar{Y}_1 = 1.65$	$\bar{Y}_2 = 1.43$
$S_1 = 0.26$	$S_2 = 0.22$

Tons per square foot (unit).

Do the soils appear to differ with respect to average shear strength at the 1% significance level ( $\alpha = 0.01$ ).

Note that  $n_1, n_2 \geq 30$ . This implies that we can use “large sample” assumptions.

i.e.  $\sigma_1 = S_1$  and  $\sigma_2 = S_2$  without loss of precision (i.e. no need for t-distribution).

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

$$H_a : \mu_1 \neq \mu_2 \leftarrow \text{This is called a 2-sided test}$$

In Chapter 8, we saw  $\bar{Y}_1 - \bar{Y}_2$  is distributed

$$N\left(\mu_1 - \mu_2, \frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}\right)$$

Under null hypothesis,

$$\begin{aligned} & N\left(0, \frac{0.26^2}{30} + \frac{0.22^2}{35}\right) \\ & = N(0, 0.00363) \end{aligned}$$

and  $\sigma_{\bar{Y}_1 - \bar{Y}_2} = 0.0603$

Note that  $\bar{Y}_1 - \bar{Y}_2 = 1.65 - 1.43 = 0.22$

$$\begin{aligned} p &= \Pr(|\bar{Y}_1 - \bar{Y}_2 - 0| > 0.22) \\ &= \Pr(\bar{Y}_1 - \bar{Y}_2 < -0.22) + \Pr(\bar{Y}_1 - \bar{Y}_2 > 0.22) \\ &= 2 \Pr(\bar{Y}_1 - \bar{Y}_2 > 0.22) \\ &= 2 \Pr\left(Z > \frac{0.22}{0.0603}\right) \quad \text{not on table 4} \\ &< 2 \Pr(Z > 3.5) \\ &= 2 \cdot 0.000233 \\ &= 0.000466 \\ &< \alpha \end{aligned}$$

Conclusion: This is statistically significant. I.e. supports  $H_a$ . The sheer strengths are different.

**Remark:**

$$\textcircled{1} \quad p = 2 \text{pnorm}(-0.22, 0, 0.0603) = 0.000265$$

- $$\textcircled{2} \quad \text{In the last 2 examples, using the empirical rule (68-95-99.7), we could have concluded "Reject } H_0 \text{" simply on the } Z\text{-score alone. (Once we're more than } 3\sigma \text{ away from } \mu, p < 100\% - 99.7\% = 0.3\%).}$$

$$\Pr(Z < Z_*) = \alpha \quad \text{and} \quad H_0 : \mu' = \mu \quad \text{and} \quad H_a : \mu' < \mu.$$

$$-Z_* = \frac{C^* - \mu}{\sigma/n} \iff C^* = \mu - Z_* \frac{\sigma}{\sqrt{n}}$$

The rejection region RR is  $\bar{X} \leq C^*$

$\mu_0 \pm Z^* \frac{\sigma}{\sqrt{n}}$  is a  $1 - \alpha$  level C.I. If  $\bar{X}$  lands in this interval, it supports  $H_0$ , else if it lands outside, reject  $H_0$ .

$$\left( \bar{X} - \mu \right) / \frac{\sigma}{\sqrt{n}}$$

## 10.4 More about errors and sample size

**Motivational example:**  $X$  equals the breaking strength of a steel bar. If the bar is manufactured by Process I, it is known  $X \sim N(50, 36)$ .

We now have Process II and it is hoped that the steel is 10% stronger. i.e.  $X \sim N(55, 36)$ .

Our test?

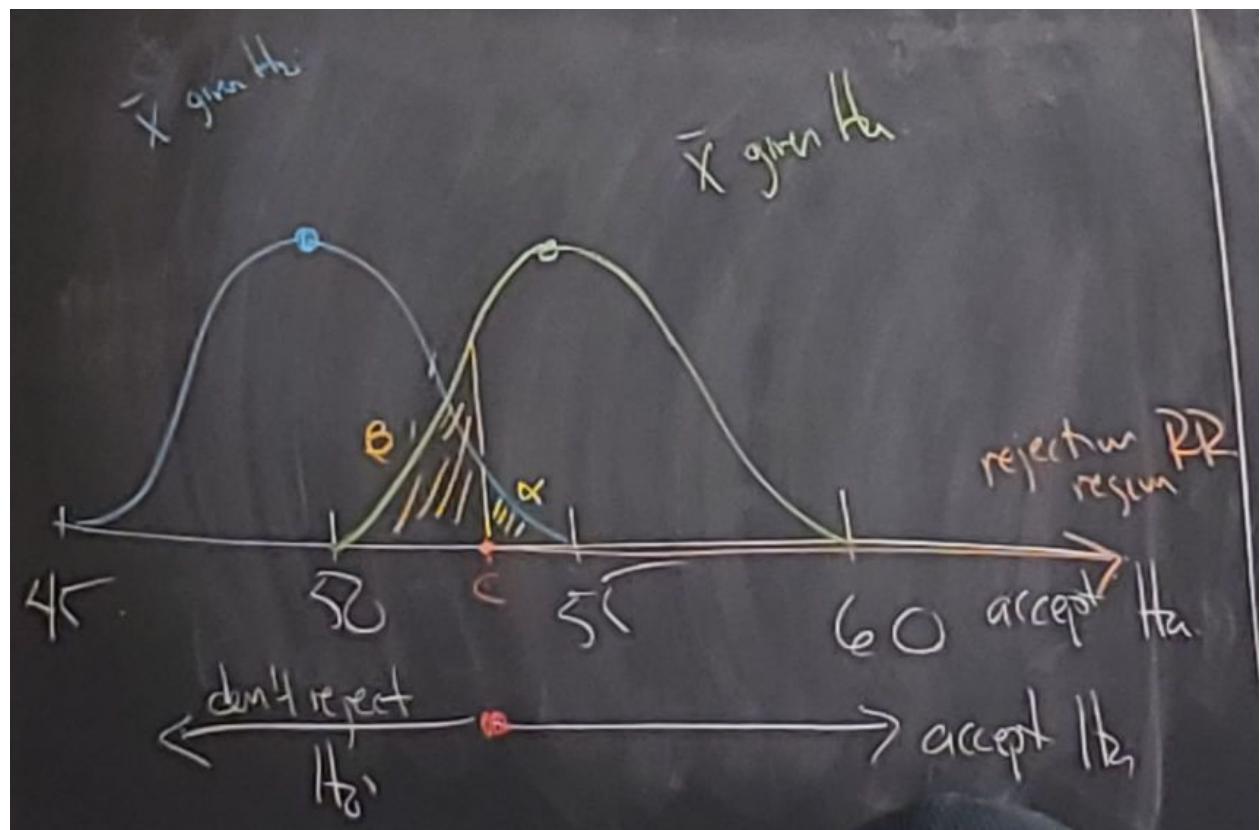
$$H_0 : \mu_I = 50 \quad H_a : \mu_{II} = 55$$

Okay... we can't really test "this"

But we can construct a hypothetical test where if  $H_a$  is true, we can minimize (or control) both the Type I and Type II errors.

For the sake of concrete-ness, set  $n = 16$ . Then,

$$\sigma_{\bar{X}}^2 = \frac{36}{16} \quad \text{and} \quad \sigma_{\bar{X}} = 1.5$$



$$\alpha = \Pr(\text{Type I}) \quad H_0 : \mu = 50 \quad H_a : \mu < 50$$

On the other hand, given  $\alpha$ , we can also see

$$\beta = \Pr(\text{Type II}) = \Pr(\text{accept } H_0 \text{ when } H_a \text{ is true})$$

Given  $\alpha$ , we can find  $c$ .

$$\begin{aligned}\alpha &= \Pr(\bar{X} > C \mid H_0) \\ &= \Pr\left(\frac{\bar{X} - 50}{1.5} > \frac{C - 50}{1.5}\right)\end{aligned}$$

and

$$\begin{aligned}\beta &= \Pr(\bar{X} < C \mid H_a) \\ &= \Pr\left(\frac{\bar{X} - 55}{1.5} < \frac{C - 55}{1.5}\right)\end{aligned}$$

For fixed  $n = 16$ , usually choose  $\alpha$  small.

$$\alpha = \Pr(\bar{X} - 50 > 2\sigma_{\bar{X}}) = 0.025$$

Then  $C = 50 + 2(1.5) = 53$  and

$$\begin{aligned}\beta &= \Pr(\bar{X} < 53 \mid H_a) \\ &= \Pr\left(\frac{\bar{X} - 55}{1.5} < 1.33\right) \\ &= 0.0918\end{aligned}$$

Note almost four times as likely to make a Type II error than a Type I error. Of course, decreasing  $\alpha$  with increase  $\beta$ .

$$\alpha = 0.01 \implies Z\text{-score} = 2.33 \quad (Z_{0.98})$$

$$c = 50 + 2.33(1.5) = 53.495$$

$$\begin{aligned}\beta &= \Pr(\bar{X} < 53.495 \mid H_a) \\ &= \Pr(Z < -1.003) \\ &= 0.1587\end{aligned}$$

Again, we note that the only way to decrease *both*  $\alpha$  and  $\beta$  is to crank up the  $n$ .

**Discussion:** Choosing sample size  $n$ . We consider the 1-sided test

$$H_0 : \mu = \mu_0 \quad H_a : \mu > \mu_0$$

Fix  $\alpha$  at the start.

$$\begin{aligned}\alpha &= \Pr(\bar{X} > C \text{ when } \mu = \mu_0) \\ &= \Pr\left(Z > \frac{C - \mu_0}{\sigma/\sqrt{n}}\right) \\ &= \Pr(Z > Z_\alpha) \quad Z_\alpha = \frac{C - \mu_0}{\sigma/\sqrt{n}}\end{aligned}$$

But, as with the power function, we need to choose specific  $\mu_a$ 's to work with (e.g.  $\mu_a = 55$ ).

$$\begin{aligned}\beta &= \Pr(\bar{X} < C \text{ when } \mu = \mu_a) \\ &= \Pr\left(Z < \frac{C - \mu_a}{\sigma/\sqrt{n}}\right) \\ &= \Pr(Z < -Z_\beta) \quad \text{when} \quad -Z_\beta = \frac{C - \mu_a}{\sigma/\sqrt{n}}\end{aligned}$$

2 equations in the 2 unknowns,  $C, n$

$$\begin{aligned}C &= \mu_0 + Z_\alpha \frac{\sigma}{\sqrt{n}} \quad \text{and} \quad C = \mu_a - Z_\beta \frac{\sigma}{\sqrt{n}} \\ \mu_0 + Z_\alpha \frac{\sigma}{\sqrt{n}} &= \mu_a - Z_\beta \frac{\sigma}{\sqrt{n}}\end{aligned}$$

Solve for  $n$ :

$$\begin{aligned}(Z_\alpha + Z_\beta) \frac{\sigma}{\sqrt{n}} &= \mu_a - \mu_0 \\ n &= \frac{(Z_\alpha + Z_\beta)^2 \sigma^2}{(\mu_a - \mu_0)^2}\end{aligned}$$

Remark:

- ① Of course all of this the fudge factor that we don't really know  $\mu_a$ .
- ② If we did  $H_0 : \mu_0 = \mu_a$  and  $H_a : \mu_0 > \mu_a$  we get the same formula for  $n$ . "sample size estimator for a one-sided  $\alpha$ -level test"

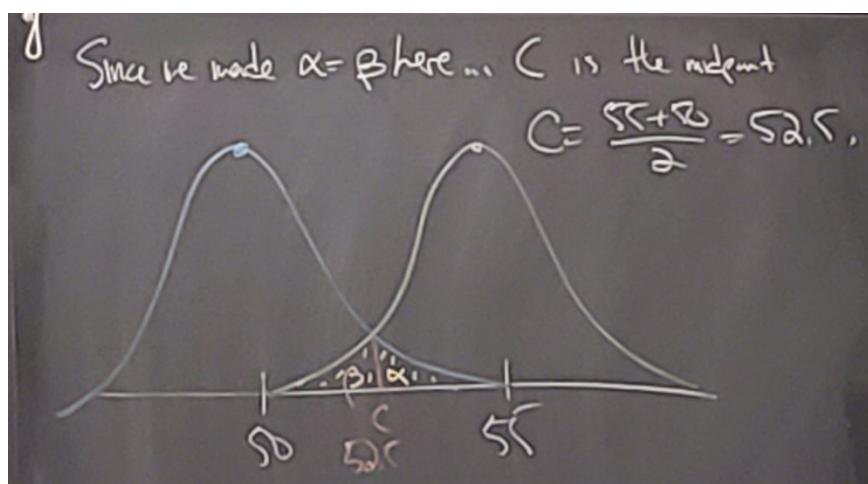
**Example:** Back to the steel example

If we decided at the start that we want  $\alpha = \beta = 0.05$ , what  $n$  should we choose?

For  $\alpha = 0.05 \implies Z_\alpha = Z_{0.05} = 1.645$ . Similarly for  $\beta$ , we need  $Z_\beta = 1.645$ . Into our formula,

$$n = \left( \frac{1.645 + 1.645}{55 - 50} \right)^2 \cdot 36 = \lceil 15.5867 \rceil = 16$$

Since we made  $\alpha = \beta$  here,  $C$  is the midpoint  $C = \frac{55 + 50}{2} = 52.5$ .



## 10.8 T tests

Recall that for a small sample  $n < 30$ . We need to use the  $t$ -distribution.

Chapter 8: t-stat confidence interval  $\bar{X} \pm t_{\alpha/2}(\text{df}) \sqrt{\frac{s^2}{n}}$

**Example:** 100 mL sample of water from swimming areas are tested for fecal coliform bacteria. It is considered to be safe if the level of bacteria is less than 400 in 3.3 oz.

20 Samples are taken: Found  $\bar{X} = 1231$  and  $S = 1038$ .

Construct a test to determine if it's safe to swim. Our test:

$$\begin{aligned} H_0 : \mu_0 &= 400 \\ H_a : \mu_a &> 400 \end{aligned}$$

Test statistic  $T = \frac{\bar{X} - \mu_0}{S/\sqrt{n}} = \frac{1231 - 1038}{1038/\sqrt{20}} = 0.350$  is the test statistic.

The degrees of freedom is  $n - 1 = 19$ . Using table 5,  $t_{0.005}(19) = 2.861$  The  $p$ -value is:

$$\begin{aligned} p &= \Pr(T \geq 3.580 \mid \mu_0 = 400) \\ &< \Pr(T \geq 2.861) \\ &= 0.005 \end{aligned}$$

Reject  $H_0$ : There is poop in the water; stay out!

*Remark:* In general:

$$H_a := \begin{cases} \mu > \mu_0 \\ \mu < \mu_0 \\ \mu \neq \mu_0 \end{cases} \implies \text{RR} := \underbrace{\begin{cases} t > t_\alpha \\ t < -t_\alpha \\ |t| > t_{\alpha/2} \end{cases}}_{\text{the } t\text{-stat}}$$

**Example:** Lifestyle comparison. Monitoring the active time (in minutes per day) between 2 populations: obese and lean.

Group	Count	Standing/Walking	$S$
Obese	$n = 10$	373.269	67.498
Lean	$n = 13$	525.751	107.121

Question: Are these two groups significant different? This is a “are population means different” question.

$$\begin{aligned} H_0 : \mu_L &= \mu_0 \implies \mu_L - \mu_0 = 0 \\ H_a : \mu_L &\neq \mu_0 \implies \mu_L - \mu_0 \neq 0 \end{aligned}$$

Questions about “different” results in a 2-sided test

Confidence interval:  $\bar{X}_1 - \bar{X}_2 \pm t_{\alpha/2}(\text{df}) \left( S_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}} \right)$

$$S_{\text{pool}} = \sqrt{\frac{(n_1 - 1)S_1^2 + (n_2 - 1)S_2^2}{n_1 + n_2 - 2}} \quad \text{with df} = n_1 + n_2 - 2$$

The difference in means  $t$  test statistic is

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

Our  $\bar{X}_L - \bar{X}_0 = 152.482$ .  $H_0 : \mu_L - \mu_0 = 0$  is the null hypothesis.  $\text{df} = 10 + 13 - 2 = 21$ .

$$\begin{aligned} S_p &= \sqrt{\frac{(13 - 1)107.121^2 + (10 - 1)67.498^2}{13 + 10 - 2}} \\ &= 92.247. \end{aligned}$$

$$\begin{aligned} T &= \frac{152.481}{92.247 \sqrt{\frac{1}{10} + \frac{1}{13}}} \\ &= 3.93 \end{aligned}$$

$$\begin{aligned} p &= \Pr(|T| > 3.93) \\ &= 2 \cdot \Pr(T > 3.93) \\ &< 2 \cdot 0.0005 && \text{(by table 5)} \\ &= 0.001 \end{aligned}$$

So  $p$  is super small ( $p < \alpha$ ) which implies that we reject the null hypothesis. That is,  $\mu_0 < \mu_L$  is true.

Our goal this week is to justify hypothesis testing.

Last day: Pearson Neyman Lamma

$X_1, \dots, X_n \sim$  via PDF  $f(x | \theta)$  when  $\theta_0$  and  $\theta_a$  are two possible values of  $\theta$ . If there exists a constant  $k$  and a subset  $C$  of the sample space such that

1.  $P\{(x_1, \dots, x_n) \in C | \theta_0\} = \alpha$
2.  $\frac{L(\theta_0)}{L(\theta_a)} \leq k$  for  $x_1, \dots, x_n \in C$ .
3.  $\frac{L(\theta_0)}{L(\theta_a)} \geq k$  for  $x_1, \dots, x_n \in \bar{C}$

Then  $C$  is a best critical region of size  $\alpha$  for testing  $H_0 : \theta = \theta_0$  versus  $H_a : \theta = \theta_a$ .

Before the proof... example.

**Example:** Let  $Y_1, \dots, Y_n$  from  $f(y | \theta) = \frac{2}{\theta}ye^{-y^2/\theta}$  and  $y > 0$ . The Rayleight distribution.

Want to test  $H_0 : \theta = \theta_0$  versus  $H_a : \theta = \theta_a$ .

$$\text{Note } L(Y_1, \dots, Y_n | \theta) = \frac{2}{\theta}y_1e^{-y_1^2/\theta} \cdots \frac{2}{\theta}y_ne^{-y_n^2/\theta}.$$

$$\underbrace{\left(\frac{2}{\theta}\right)^n \exp\left[-\left(\sum_1^n y_i^2\right)/\theta\right]}_{g(S, \theta)} \underbrace{\prod_1^n y_i}_h$$

We will come back to this and connect to sufficiency next day.

The N-P ratio of likelihood function becomes

$$\begin{aligned} \frac{L(\theta_0)}{L(\theta_a)} &= \frac{\left(\frac{2}{\theta_0}\right)^n \exp\left[-(\sum_1^n y_i^2)/\theta_0\right] \prod_1^n y_i}{\left(\frac{2}{\theta_a}\right)^n \exp\left[-(\sum_1^n y_i^2)/\theta_a\right] \prod_1^n y_i} \\ &= \left(\frac{\theta_a}{\theta_0}\right)^n \exp\left[-\sum y_i^2 \cdot \left(\frac{1}{\theta_0} - \frac{1}{\theta_a}\right)\right] \end{aligned}$$

We want  $C$  to relate to our rejection region RR.

So (2), implies “reject  $H_0$  if”

$$\left(\frac{\theta_a}{\theta_0}\right)^n \exp\left[-\sum y_i^2 \cdot \left(\frac{1}{\theta_0} - \frac{1}{\theta_a}\right)\right] \leq k$$

This looks scary, but  $\theta_a$  and  $\theta_0$  are fixed and  $\theta_0 < \theta_a$  so  $\left(\frac{1}{\theta_0} - \frac{1}{\theta_a}\right) > 0$ . In the end, making  $\frac{L(\theta_0)}{L(\theta_a)}$  “small enough” for  $k$  simlifies does to the condition that  $\sum_i^n y_i^2$  is large enough. i.e. we need  $\sum_1^n y_i^2 > k'$  for some  $k'$  (dependent upon  $k$ ).

New problem: To determine an appropriate  $k'$ , we need to know how  $S = \sum_1^n y_i^2$  is distributed.

Using CDF method (§ 6.3) we can show that the distribution of  $y^2$  is exponention with mean  $\theta$ . Then, via products of moment generating functions.

$$S = \sum_1^n y_i^2 \sim \text{Gamma}\left(n, \frac{1}{\theta}\right)$$

Stop! What just happened.

The N-P lemma tells us that the most powerful test of  $H_0$  vs  $H_a$  is to use the statistic  $S = \sum y_i^2$ . Then, given any  $\alpha$  level significance (5%, 1%, whatever), we use the test on  $S$  to be if  $S$  is larger than  $100(1 - \alpha)$  percentile of the Gamma  $(n, \frac{1}{\theta})$  distribution... we reject  $H_0$ .

Picture with some  $\gamma_*$  such that  $\int_{\gamma_*}^{\infty} = \alpha$ .

$S > \gamma_*$ , we reject the null hypothesis. Moreover, by the NP lemme, we don't have to compare the powe rof other possible tests because the lemma says any data

$$(Y_1, \dots, Y_n) \in C \iff S \geq \gamma_*$$

$$C := \left\{ (Y_1, \dots, Y_n) \mid \sum y_i^2 \geq \gamma_* \right\}$$

which is itself a consequence of our significance level  $\alpha$  results in the most powerful test.

$$P(C \mid \theta_a) \geq P(D \mid \theta_a)$$

where  $D$  is another critical region ( $P(D \mid \theta_a) = \alpha$ ).

**Definition:** The test using the best critical region is called the most powerful test.

Recap: We wanted to construct a hypothesis test at  $\alpha$  significance level. All in one fell swoop, the N-P lemma says

1. Here is the stat to use ( $S = \sum y_i^2$  in last example)
2. Here is your rejection region ( $S \geq \gamma_a$  in last example)
3. in fact, this is the most powerful test possible.

(This is kinda awesome)

Theorem: NP lemma;

*Proof.* (Continuous)

Let  $B \subseteq \mathbb{R}^n$  and define  $\underbrace{\int_B L(\theta)}_{ntimes} = \int \cdots \int L(x_1, \dots, x_n \mid \theta) dx_1 dx_2, \dots, dx_n$ .

Assume there exists a critical region  $C$  satisfying bullets 1 2 and 3.

$\alpha$  fixed,  $\alpha = \int_C L(\theta_0)$  (same as  $P(C \mid \theta_0)$ )

Need to prove “best”, i.e.  $P(C \mid \theta_a) \geq P(D \mid \theta_a)$

Assume  $D$  is another critical region,

$$\alpha = \int_D L(\theta)_0 = P(D \mid \theta_0)$$

$$\text{So } 0 = \int_C L(\theta_0) - \int_D L(\theta_0)$$

$$\begin{aligned} 0 &= \int_{C \cap D^C} L(\theta_0) + \int_{C \cap D} L(\theta_0) - \left[ \int_{C \cap D} L(\theta_0) + \int_{C^C \cap D} L(\theta_0) \right] \\ &= \int_{C \cap D^C} L(\theta_0) - \int_{C^C \cap D} L(\theta_0) \end{aligned}$$

By (2), there exists  $k$  such that  $\frac{L(\theta_0)}{L(\theta_a)} \leq k$  i.e.  $L(\theta_0) \leq kL(\theta_a)$  at every point in  $C$ .

$$\leq k \int_{C \cap D^C} L(\theta_a) - \int_{C^C \cap D} L(\theta_0)$$

By (3) there exists  $k$  such that  $\frac{L(\theta_0)}{L(\theta_a)} \geq k$  i.e.  $L(\theta_0) \geq kL(\theta_a)$

$$\leq k \left[ \int_{C \cap D^C} L(\theta_a) - \int_{C^C \cap D} L(\theta_a) \right]$$

$$= k \left[ \int_{C \cap D^C} L(\theta_a) + \int_{C \cap D} L(\theta_a) - \left[ \int_{C \cap D} L(\theta_a) + \int_{C^C \cap D} L(\theta_a) \right] \right]$$

$$= k \left[ \int_C L(\theta_a) - \int_D L(\theta_a) \right]$$

$$\implies \int_C L(\theta_a) - \int_D L(\theta_a) \geq 0$$

$$\implies \int_C L(\theta_a) \geq \int_D L(\theta_a)$$

$$\implies P(C \mid \theta_a) \geq P(D \mid \theta_a)$$

Here,  $C$  is a best critical region of size  $\alpha$  by definition. □

**Example:** Back to our old sample size example.

$$X_1, \dots, X_n \sim N(\mu, 36)$$

We played with  $H_0 : \mu = 50$  and  $H_a : \mu = 55$ .

Consider the ratio of the likelihood functions.

$$\begin{aligned}
\frac{L(50)}{L(55)} &= \frac{(72\pi)^{-n/2} \exp \left[ -\left(\frac{1}{72}\right) \sum_1^n (X_i - 50)^2 \right]}{(72\pi)^{-n/2} \exp \left[ -\left(\frac{1}{72}\right) \sum_1^n (X_i - 55)^2 \right]} \\
&= \exp \left[ -\frac{1}{72} \sum_1^n [(X_i - 50)^2 - (X_i - 55)^2] \right] \\
&= \exp \left[ -\frac{1}{72} \sum_1^n (10X_i - 525) \right] \\
&= \exp \left[ -\frac{5}{36} \sum_1^n X_i + \frac{175n}{24} \right] \leq k \text{ to satisfy (2)} \\
&\implies -\frac{5}{36} \sum_1^n X_i + \frac{175n}{24} \leq \ln k \\
&\implies \sum_1^n X_i \geq \frac{105n}{2} - \frac{36}{5} \ln k \\
&\implies \bar{X} = \frac{1}{n} \sum_1^n X_i \geq \frac{105}{2} - \frac{36 \ln k}{5n}
\end{aligned}$$

We have our stat  $\bar{X} \geq k'$  thus will define our rejectoin rejection.

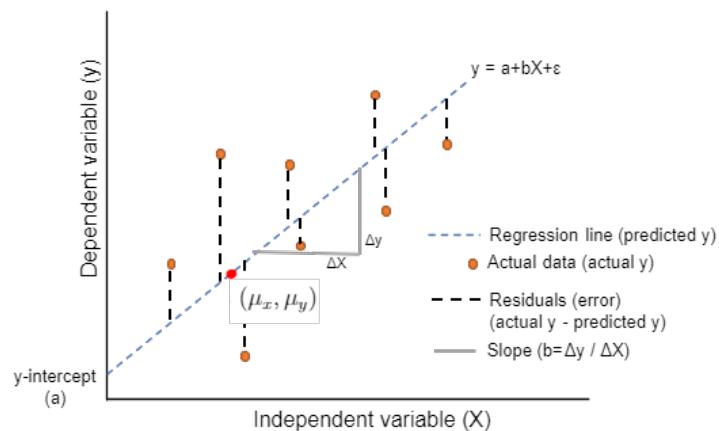
# Chapter 11

## Linear Models and Estimation by Least Squares

### 11.1 Introduction

We are going to do best-fit lines in 2 separate ways.

Consider some sample (data)  $S = \{(x_i, y_i)\}$ . In the background, there can be an underlying joint PDF  $f(x, y)$



### 11.2 Linear Statistical Models

Assumptions

- ①  $Y$  is dependent upon  $X$ . (i.e. not an independent random variable.)
- ② In general,  $y = mx + b$ , the general relationship between  $X$  and  $Y$ .

## 11.3 The Method of Least Squares

### A.) The probabilistic construction:

Given all possible lines  $y = mx + b$  that *could* describe the data, the “most likely” one will be a line that intersects the point  $(\mu_x, \mu_y)$ .

That is  $y - \mu_y = m(x - \mu_x)$ . Because we are interpreting  $y$  and a function of  $x$ , the “best line” should be one that minimizes the vertical distance (error (residuals)) in the  $y$ -direction.

For a point  $(x_k, y_k)$  in  $S$ , the vertical distance is

$$\text{dist} = |y_k - (m(x_k - \mu_x) + \mu_y)|.$$

Minimizing absolute values is a pain, but calculus I suggests we square this,

$$\text{dist}^2 = ((y_k - \mu_y) - m(x_k - \mu_x))^2.$$

As  $x, y$  are distributed via some PDF, the best way to compute our minimization problem is to minimize expectation with respect to slope  $m$ :

$$K(m) := E(\text{dist}^2(m)) = \underbrace{E\left(\left((y_k - \mu_y) - m(x_k - \mu_x)\right)^2\right)}_{**}.$$

**\*\*** The  $m$  that minimizes is called the solution to our least-squares problem.

$$\begin{aligned} K(m) &= E((y_k - \mu_y)^2 - 2m(x_k - \mu_x)(y_k - \mu_y) + m^2(x_k - \mu_x)^2) \\ &= E((y_k - \mu_y)^2) - 2m E((x_k - \mu_x)(y_k - \mu_y)) + m^2 E((x_k - \mu_x)^2) \\ &= S_y^2 - 2m \text{Cov}(X, Y) + m^2 S_x^2 \end{aligned}$$

Minimizing (with respect to  $m$ )  $\frac{\partial}{\partial m}$ ,

$$K'(m) = -2 \text{Cov}(X, Y) - 2mS_x^2.$$

With  $K'(m) = 0$  then  $m = \frac{\text{Cov}(X, Y)}{S_x^2}$ .

Note that by the second derivative test,  $K''(m) = 2S_x^2 > 0$ . In AP Stats (MTH 111), the slope is moved to a “prettier” form:

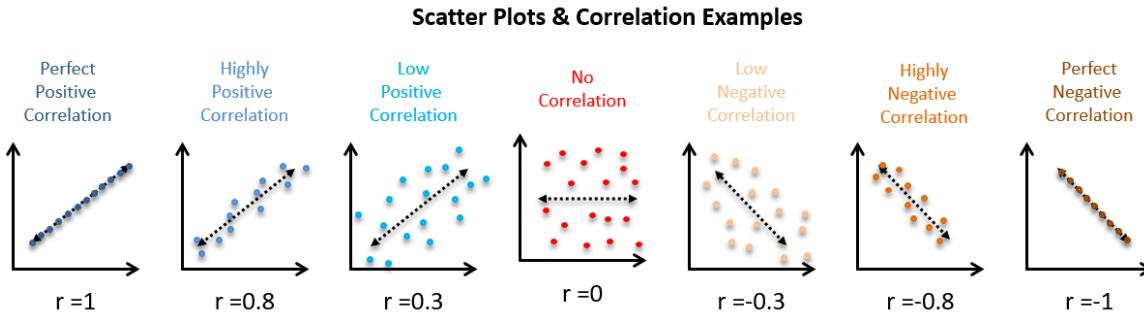
$$\begin{aligned} m &= \frac{\text{Cov}(X, Y)}{S_x^2} \\ &= \frac{\text{Cov}(X, Y)}{S_x S_y} \cdot \frac{S_y}{S_x} \\ &= \rho \cdot \frac{S_y}{S_x} \quad \rho := \frac{\text{Cov}(X, Y)}{S_x S_y} \end{aligned}$$

**FACT #1** : The sign of the slope depends entirely upon  $\rho$ , where the sign is dependent upon  $\text{Cov}(X, Y)$

**FACT #2** :  $-1 \leq \rho \leq 1$ .

(For my MTH 325 class, we proved this fact in full generality).

**FACT #3** : The closer  $|\rho|$  is to 1, the better the line  $\hat{y} = mx + b$  “fits” the data.



The AP Stats definition for the best fit line is  $\hat{y} = \hat{b}_0 + \hat{b}_1 x$  where  $\hat{b}_1 = \rho \frac{S_y}{S_x}$  and  $\hat{b}_0 = \bar{y} - \hat{b}_1 \bar{x}$ .

### B.) The linear algebra construction:

We still have  $S = \{(x_k, y_k)\}$ . We want  $y = mx + b$ . Using  $S$  always yields an overdetermined (i.e. inconsistent) system

$$\begin{aligned} y_1 &= mx_1 + b \\ y_2 &= mx_2 + b \\ &\vdots \\ y_n &= mx_n + b. \end{aligned}$$

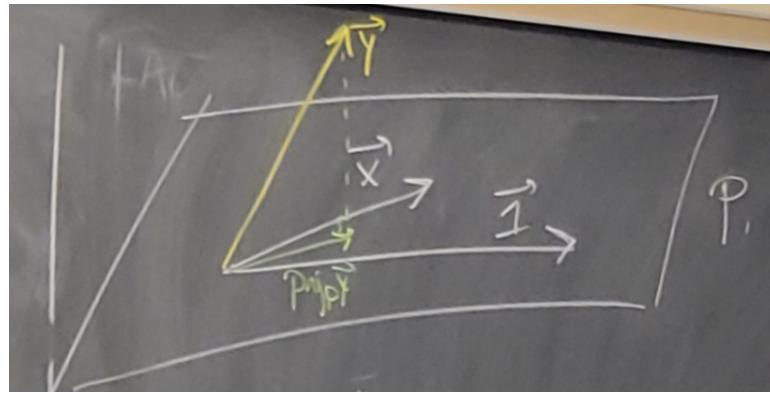
This is  $n$  equations with 2 variables, but we can rewrite as vector equations.

$$\vec{y} = m\vec{x} + b\vec{1} \quad \text{where } \vec{y} = [y_1, \dots, y_n]^T, \vec{x} = [x_1, \dots, x_n]^T, \text{ and } \vec{1} = \underbrace{[1, 1, \dots, 1]^T}_{n \text{ components}}.$$

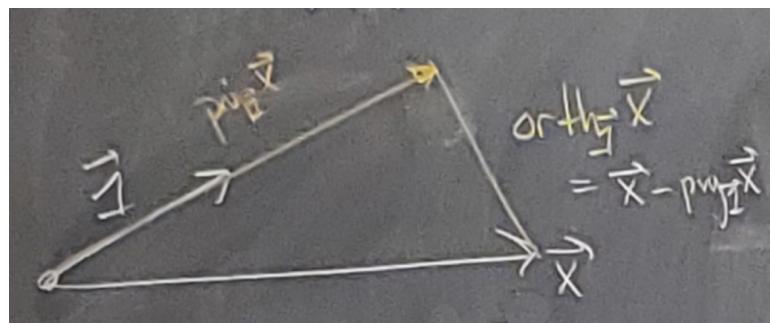
Trying to write  $\vec{y}$  as a linear combination of  $\vec{x}$  and  $\vec{1}$ . Consider the plane  $P = \text{span}\{\vec{x}, \vec{1}\}$ . Inconsistent implies that  $\vec{y} \notin P$ .

But the vector in  $P$  that is closest to  $\vec{y}$  is  $\text{proj}_P \vec{y}$ . (Closest under Euclidean distance), recall  $\langle u, v \rangle = u \cdot v$  in  $\mathbb{R}^n$ .

$$\text{dist}^2 = \sum (y_i - v_i)^2, \quad \vec{v} \in P = \underbrace{\vec{y} \cdot \vec{v} = \langle \vec{y}, \vec{v} \rangle}_{\text{“least squares”}}$$



Easiest way to project onto  $P$  is to have an orthogonal basis for  $P$ . We need an orthogonal basis for  $P$ :



$$P = \text{span} \{1, \text{orth}_1 x\} = \text{span} \{1, v\} \text{ where}$$

$$\begin{aligned} v &= x - \text{proj}_1 x \\ &= x - \frac{\langle x, 1 \rangle}{\langle 1, 1 \rangle} 1 \\ &= x - \frac{\langle x, 1 \rangle}{n} 1 \end{aligned}$$

Then

$$\begin{aligned} \text{proj}_P y &= \frac{\langle y, 1 \rangle}{\langle 1, 1 \rangle} 1 + \frac{\langle y, v \rangle}{\langle v, v \rangle} v \\ &= \frac{\langle y, 1 \rangle}{n} 1 + \frac{\langle y, v \rangle}{\langle v, v \rangle} x - \frac{\langle y, v \rangle}{\langle v, v \rangle} \cdot \frac{\langle x, 1 \rangle}{n} 1 \\ &= \underbrace{\frac{\langle y, v \rangle}{\langle v, v \rangle} x}_{\text{(slope)}} + \underbrace{\left( \frac{\langle y, 1 \rangle}{n} - \frac{\langle y, v \rangle}{\langle v, v \rangle} \cdot \frac{\langle x, 1 \rangle}{n} \right) 1}_{\text{(y-intercept)}} \end{aligned}$$

$$\begin{aligned}
\langle v, v \rangle &= \left\langle x - \frac{\langle x, 1 \rangle}{n} 1, x - \frac{\langle x, 1 \rangle}{n} 1 \right\rangle \\
&= \langle x, x \rangle - 2 \left\langle x, \frac{\langle x, 1 \rangle}{n} 1 \right\rangle + \left\langle \frac{\langle x, 1 \rangle}{n} 1, \frac{\langle x, 1 \rangle}{n} 1 \right\rangle \\
&= \langle x, x \rangle - 2 \frac{\langle x, 1 \rangle}{n} \langle x, 1 \rangle + \frac{\langle x, 1 \rangle^2}{n^2} \langle 1, 1 \rangle \\
&= \langle x, x \rangle - 2 \frac{\langle x, 1 \rangle^2}{n} + \frac{\langle x, 1 \rangle^2}{n} \\
&= \frac{n \langle x, x \rangle - \langle x, 1 \rangle^2}{n}
\end{aligned}$$

$$\begin{aligned}
\langle v, y \rangle &= \langle x, y \rangle - \frac{\langle x, 1 \rangle \langle y, 1 \rangle}{n} \\
&= n \langle x, y \rangle - \langle x, 1 \rangle \langle y, 1 \rangle
\end{aligned}$$

Then

$$\begin{aligned}
m &= \frac{\langle v, y \rangle}{\langle v, v \rangle} \\
&= \frac{n \langle x, y \rangle - \langle x, 1 \rangle \langle y, 1 \rangle}{n \langle x, x \rangle - \langle x, 1 \rangle^2}
\end{aligned}$$

Then

$$\begin{aligned}
b &= \frac{\langle y, 1 \rangle}{n} - m \frac{\langle x, 1 \rangle}{n} \\
&= \bar{y} - m\bar{x}
\end{aligned}$$

Which is what we got from the probabilistic approach.

Claim: The two-sets of formulas are identical

$$\begin{aligned}
m &= \frac{\text{Cov}(X, Y)}{S_x^2} = \rho \frac{S_y}{S_x} \\
&= \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}
\end{aligned}$$

Useful formula:

$$\begin{aligned}
\sum (x_i - \bar{x})(y_i - \bar{y}) &= \sum (x_i y_i - x_i \bar{y} - y_i \bar{x} + \bar{x} \bar{y}) \\
&= \sum x_i y_i - \bar{y} \sum x_i - \bar{x} y_i + \bar{x} \bar{y} \sum 1 \\
&= \sum x_i y_i - \bar{y} n \cdot \frac{1}{n} \sum x_i - \bar{x} n \cdot \frac{1}{n} \sum y_i + n \bar{x} \bar{y} \\
&= \sum x_i y_i - n \bar{x} \bar{y} - n \bar{x} \bar{y} + n \bar{x} \bar{y} \\
&= \sum x_i y_i - n \bar{x} \bar{y}
\end{aligned}$$

$$\begin{aligned}
m &= \frac{\text{Cov}(X, Y)}{S_x^2} = \rho \frac{S_y}{S_x} \\
&= \frac{\sum(x_i - \bar{x})(y_i - \bar{y})}{\sum(x_i - \bar{x})^2} \\
&= \frac{\sum x_i y_i - n \bar{x} \bar{y}}{\sum x_i^2 - n \bar{x}^2} \\
&= \frac{\langle x, y \rangle - n \left( \frac{\langle x, 1 \rangle}{n} \right) \left( \frac{\langle y, 1 \rangle}{n} \right)}{\langle x, x \rangle - n \left( \frac{\langle x, 1 \rangle}{n} \right)^2} \\
&= \frac{n \langle x, y \rangle - \langle x, 1 \rangle \langle y, 1 \rangle}{n \langle x, x \rangle - \langle x, 1 \rangle^2} \\
&= \frac{n \sum x_i y_i - (\sum x_i)(\sum y_i)}{n \sum x_i^2 - (\sum x_i)^2} \\
&= m
\end{aligned}$$

## 11.4 Properties of the Least-Squares Estimators: Simple Linear Regression

The Coefficents of the Best-Fit Line are Estimators. We have shown the best-fit line to be  $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$ .

Where  $\hat{\beta}_1 = \frac{S_{xy}}{S_{xx}}$  with  $S_{xy} = \sum(x_i - \bar{x})(y_i - \bar{y})$  and  $S_{xx} = \sum(x_i - \bar{x})^2$ . (i.e.  $\text{Cov}(X, Y)$  and  $\text{Var}[X]$  via use of  $\bar{x}, \bar{y}$  for  $\mu_x, \mu_y$  ).

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

We recognize  $\hat{\beta}_0$  and  $\hat{\beta}_1$  as stats dependent upon  $\bar{x}, S_x^2$ , and  $\bar{y}$ .

But what exactly are they estimating? In theory,  $X \times Y$  distributed via pdf  $f(x, y)$  and there is some “best” linear relationship over the same probabiltiy space.

$$y = \beta_0 + \beta_1 x.$$

In other words, when we are “at”  $x_i$ ,

$$\text{E}[Y] = \beta_0 + \beta_1 x_i.$$

A modern approach to the distribution of  $Y$  at  $x$  is to add an error parameter  $\varepsilon$ . That is,

$$y = \underbrace{\beta_0 + \beta_1 x}_{\substack{\text{deterministic} \\ \text{component} \\ \text{of } Y}} + \underbrace{\varepsilon}_{\text{“random” component}}$$

Still want  $E[Y] = \beta_0 + \beta_1 x$ . By linearity,

$$\begin{aligned} E[Y] &= E[\beta_0 + \beta_1 x + \varepsilon] \\ &= E[\beta_0 + \beta_1 x] + E[\varepsilon] \\ \implies E[\varepsilon] &= 0 \quad \text{error averages to zero} \end{aligned}$$

We make the additional assumption that the variance of  $\varepsilon$  is independent of  $x$ . That is  $\text{Var}[Y] = \text{Var}[\varepsilon] = \sigma^2$ . The “value” of  $y$  depends on  $x$ , but the spread of the  $y$ -values does not.

**Proposition:**  $\hat{\beta}_0, \hat{\beta}_1$  are unbiased estimators of  $\beta_0, \beta_1$  where  $Y = \beta_0 + \beta_1 x + \varepsilon$ .

**Reason:**

$$\begin{aligned} E(\hat{\beta}_1) &= E\left[\frac{S_{xy}}{S_{xx}}\right] \\ &= E\left[\frac{\sum(x_i - \bar{x})(y_i - \bar{y})}{S_{xx}}\right] \\ &= E\left[\frac{\sum(x_i - \bar{x})y_i - \sum(x_i - \bar{x})\bar{y}}{S_{xx}}\right] \\ &= E\left[\frac{\sum(x_i - \bar{x})y_i - \bar{y}\overbrace{\sum(x_i - \bar{x})}^{0 \text{ by def'n of } \bar{x}}}{S_{xx}}\right] \\ &= E\left[\frac{\sum(x_i - \bar{x})y_i}{S_{xx}}\right] \\ &= \frac{\sum(x_i - \bar{x})E(y_i)}{S_{xx}} \\ &= \frac{\sum(x_i - \bar{x})(\beta_0 + \beta_1 x_i)}{S_{xx}} \\ &= \frac{\overbrace{\beta_0 \sum(x_i - \bar{x})}^0 + \beta_1 \sum(x_i - \bar{x})x_i}{S_{xx}} \\ &= \frac{\beta_1 \sum(x_i^2 - x_i \bar{x})}{S_{xx}} \\ &= \frac{\beta_1 (\sum x_i^2 - \bar{x} \sum x_i)}{S_{xx}} \\ &= \frac{\beta_1 (\sum x_i^2 - n\bar{x}^2)}{S_{xx}} \\ &= \beta_1 \frac{S_{xx}}{S_{xx}} \\ &= \beta_1 \end{aligned}$$

Therefore our estimator is unbiased.

The point of all that:  $E(\hat{\beta}_1) = \hat{\beta}_1$ .

For  $E(\hat{\beta}_0) = E[\bar{y} - \hat{\beta}_1 \bar{x}] = E[\bar{y}] - \bar{x} E(\hat{\beta}_1)$ .

But  $\bar{y} = \frac{1}{n} \sum y_i = \frac{1}{n} \sum (\beta_0 + \beta_1 x + \varepsilon) = \beta_0 + \beta_1 \bar{x} + \varepsilon$ .

and  $E[\bar{y}] = \beta_0 + \beta_1 \bar{x}$ . and  $E[\hat{\beta}_0] = \beta_0 + \beta_1 \bar{x} - \beta_1 \bar{x} = \beta_0$ .

**Corollary:**  $\text{Var}(\hat{\beta}_1) = \frac{\sigma^2}{S_{XX}}$ ,  $\text{Var}(\hat{\beta}_0) = \frac{\sigma^2 \sum x_i^2}{n S_{XX}}$ .

**Reason**

$$\begin{aligned}\text{Var}(\hat{\beta}_1) &= \text{Var}\left(\frac{\sum(x_i - \bar{x})(y_i - \bar{y})}{\sum(x_i - \bar{x})^2}\right) \\ &= \text{Var}\left(\frac{\sum(x_i - \bar{x})^2 y_i}{S_{XX}}\right) \\ &= \frac{1}{S_{XX}^2} \text{Var}\left(\sum(x_i - \bar{x}) y_i\right) \\ &= \frac{\sum(x_i - \bar{x})^2 \text{Var}(Y_i)}{S_{XX}^2}\end{aligned}$$

assumption:  $\text{Var}(Y) = \text{Var}(\varepsilon) = \sigma^2$

$$\begin{aligned}&= \frac{\sum(x_i - \bar{x})^2 \sigma^2}{S_{XX}^2} \\ &= \frac{\sigma^2 \sum(x_i - \bar{x})^2}{S_{XX}^2} \\ &= \frac{\sigma^2 S_{XX}}{S_{XX}^2} \\ &= \frac{\sigma^2}{S_{XX}}.\end{aligned}$$

$$\text{Var}(\hat{\beta}_0) = \text{Var}(\bar{y} - \hat{\beta}_1 \bar{x})$$

Note in our  $\varepsilon$  set up,  $\bar{y}$  and  $\hat{\beta}_1$  both functions of  $Y_i$ 's... may be independent.

$$\begin{aligned}&= \text{Var}(\bar{y}) - \text{Var}(\hat{\beta}_1 \bar{x}) - 2 \text{Cov}(\bar{y}, \hat{\beta}_1 \bar{x}) \\ &= \text{Var}(\bar{y}) + \bar{x}^2 \text{Var}(\hat{\beta}_1) - 2 \bar{x}^2 \text{Cov}(\bar{y}, \hat{\beta}_1)\end{aligned}$$

And i.)  $\text{Var}(\bar{y}) = \text{Var}(\bar{\varepsilon}) = \frac{1}{n} \text{Var}(\varepsilon) = \frac{1}{n} \sigma^2$ . And ii.)

$$\begin{aligned}\text{Cov}(\bar{y}, \hat{\beta}_1) &= \text{Cov}\left(\frac{1}{n} \sum y_i, \sum \frac{(x_i - \bar{x})y_i}{S_{XX}}\right) \\ &= \sum \frac{(x_i - \bar{x}) \text{Var}(Y_i)}{n S_{xx}} + \sum \sum_{i \neq j} \frac{(x_j - \bar{x})(x_i - \bar{x})}{n S_{XX}} \underbrace{\text{Cov}(Y_i, Y_j)}_{\text{independent} = 0} \\ &= \frac{\sigma^2}{n S_{XX}} \sum (x_i - \bar{x}) \\ &= 0.\end{aligned}$$

Then

$$\begin{aligned}\text{Var}(\hat{\beta}_0) &= \text{Var}(\bar{y}) + \bar{x}^2 \text{Var}(\hat{\beta}_1) - 2\bar{x}^2 \text{Cov}(\bar{y}, \hat{\beta}_1) \\ &= \frac{\sigma^2}{n} + \bar{x}^2 \cdot \frac{\sigma^2}{S_{XX}} \\ &= \sigma^2 \left( \frac{S_{XX} + n\bar{x}^2}{n S_{XX}} \right)\end{aligned}$$

$$\begin{aligned}S_{XX} &= \sum (x_i - \bar{x})^2 \\ &= \sum (\bar{x}^2 - 2\bar{x}) + \sum x_i^2 \\ &= n\bar{x}^2 - 2\bar{x} \sum x_i \\ &= n\bar{x}^2 - 2\bar{x}n \left( \frac{\sum x_i}{n} \right) + \sum x_i^2 \\ &= \sum x_i^2 - n\bar{x}^2 \quad \text{i.e. } n\bar{x}^2 = \sum x_i^2 - S_{XX}\end{aligned}$$

$$\begin{aligned}\text{Var}(\hat{\beta}_0) &= \sigma^2 \left( \frac{S_{XX} + \sum x_i^2 - S_{XX}}{n S_{XX}} \right) \\ &= \frac{\sigma^2 \sum x_i^2}{n S_{XX}}\end{aligned}$$

Corollary:  $\text{Cov}(\hat{\beta}_0, \hat{\beta}_1) = -\frac{\bar{x}\sigma^2}{S_{XX}}$ . Proof in textbook. Note:  $\hat{\beta}_0, \hat{\beta}_1$  guaranteed to be dependent when  $\bar{x} = 0$ .

**Topic:** estimating  $\sigma^2$ .

We have working with the assumption that  $\text{Var}(y) = \text{Var}(\varepsilon) = \sigma^2$ , but this is usually unknown.

In the past, we used  $\text{Var}(Y) = \frac{1}{n-1} \sum (y_i - \bar{y})^2$ .

In our new setup, our point estimator for  $y_i$  is no longer  $\bar{y}$ . Our estimator of  $Y_i$  is the best-fit line:  $E(Y_i) = \hat{\beta}_0 + \hat{\beta}_1 x$ .

We define  $S^2$  via sum square error

$$\text{SSE} = \sum_{i=1}^n (y_i - \hat{y}_i)^2.$$

The term inside the summation is the residual; it is the error between the prediction and observed.

We need to define an unbiased estimator for  $\sigma^2$ . Let

$$\begin{aligned}\hat{\theta} &:= K \cdot \text{SSE} \\ &= K \sum (y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i))^2\end{aligned}$$

By computing  $E(\hat{\theta})$ , we choose  $K$  so that  $E(\hat{\theta}) = \sigma^2$ .

As we did in chapter 9, we can show that  $E(\text{SSE}) = (n - 2)\sigma^2$ . So,

$$S^2 = \frac{1}{n-2} \text{SSE} = \frac{1}{n-2} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

**Proposition:**  $E(S^2) = \sigma^2$ , (unbiased estimator).

(aside: It would be nice to have a computation corollary like we did, i.e.  $\text{Var}(x) = E(x^2) - E(x)^2$ )

**Corollary:** The computation corollary

$$\begin{aligned}\text{SSE} &= S_{YY} - \hat{\beta}_1 S_{XY} \\ \text{where } S_{YY} &= \sum_{i=1}^n (y_i - \bar{y})^2.\end{aligned}$$

**Example:** (11.16)

Potency of antibiotic ( $y$ ): (38, 43, 29), (32, 26, 33), (19, 27, 27), (14, 19, 21)

After storage at  $x^\circ\text{F}$ : 30, 50, 70, 90

$$\bar{y} = 27, \quad \bar{x} = 60$$

$$S_{xy} = \sum (x_i - \bar{x})(y_i - \bar{y}) = -1900$$

$$S_{xx} = \sum (x_i - \bar{x})^2 = 6000$$

$$S_{yy} = \sum (y_i - \bar{y})^2 = 792$$

Then  $\hat{\beta}_1 = \frac{S_{xy}}{S_{xx}} = -\frac{19}{60}$  and

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} = 27 - \left( \frac{-19}{60} \right) 60 = 46.$$

So the best fit line is  $\hat{y} = 46 - \frac{19}{60}x$ . For

$$\begin{aligned} S^2 &= \frac{1}{n-2} \text{SSE} \\ &= \frac{1}{n-2} (S_{yy} - \hat{\beta}_1 S_{xy}) \\ &= \frac{1}{12-2} \left( 792 - \left( \frac{-19}{60} \right) (-1900) \right) = \frac{571}{30} \\ &\approx 19.0\bar{3} \end{aligned}$$

**Discussion:** Do we know how  $S^2$  is distributed? This depends on our assumption on how  $\varepsilon$  is distributed.

$$y = \beta_0 + \beta_1 x + \varepsilon.$$

We have  $E(\varepsilon) = 0$  and  $\text{Var}(\varepsilon) = \sigma^2$ , and thus are independent of  $X$ .

Note that the distribution of the point estimators  $\hat{\beta}_0$  and  $\hat{\beta}_1$  are  $\sigma^2$  dependent (which we will estimate with  $S^2$ ).

However, the common assumption is that “everything” is normal! We assume  $\varepsilon \sim N(0, \sigma^2)$ .

We can now prove a Fisher’s Theorem like result for the new  $S^2$ .

Namely... Theorem:  $\frac{(n-2)S^2}{\sigma^2} = \frac{\text{SSE}}{\sigma^2} \sim \chi^2(n-2)$ , moreover  $S^2$  is independent of  $\hat{\beta}_0$  and  $\hat{\beta}_1$ . Proof: Omitted for mental health reasons.

## 11.5 Inferences concerning the point estimators

Under the assumption  $\varepsilon \sim N(0, \sigma^2)$ , we have that  $\hat{\beta}_0$  and  $\hat{\beta}_1$  are normal  $N(\beta_i, \text{Var}(\beta_i))$ . Thus we can do confidence intervals for the true  $\hat{\beta}_0$  and  $\hat{\beta}_1$ .

$$y = \beta_0 + \beta_1 x + \varepsilon$$

Use the same  $\mathcal{Z}/\mathcal{T}$  rules as before, ( $n \leq 30$  use  $\mathcal{T}$  and  $n > 30$  use  $\mathcal{Z}$ ).

As always, the standardization is

$$\frac{\hat{\theta} - \theta}{\sqrt{\text{Var}(\hat{\theta})}} = \frac{\hat{\theta} - \theta}{\sqrt{\text{Var}(\hat{\theta})}}$$

**Example:** Given  $Z_{\alpha/2}$  ( $n > 30$ ),  $\hat{\beta}_1 \pm Z_{\alpha/2} \sqrt{\text{Var}(\hat{\beta}_1)}$  is an  $\alpha$ -level 2-sided C.I. for the true  $\beta_i$ .

**Example:** Find a 90% C.I. for the slope  $\hat{\beta}_1$  of the potency example. We have  $\hat{\beta}_1 = \frac{-19}{60} \approx$

$-0.3166\bar{7}$ .

$$\begin{aligned}\text{Var}(\hat{\beta}_1) &= \frac{\sigma^2}{S_{xx}} \\ &\approx \frac{S^2}{S_{xx}} \\ &= \frac{19.03}{6000} \\ &= 0.00317 \\ \implies S_{\hat{\beta}_1} &= \sqrt{\text{Var}(\hat{\beta}_1)} = 0.0563\end{aligned}$$

Using  $S^2$  implies we use  $S^2$  degrees of freedom,  $n = 12$  implies 10 degrees of freedom. Thus  $S^2 \sim \chi^2(10)$ . We need  $t_{0.05}(10) = 1.812$  by table and

$$\begin{aligned}\hat{\beta}_1 &\pm t_{0.05}(10)S_{\hat{\beta}_1} \\ &- 0.31667 \pm (1.812)(0.0563) \\ &- 0.31667 \pm 0.1020 \\ &(-0.4187, -0.214)\end{aligned}$$

### Example: (Potency again)

We will assume that it is commonly known that penicillin derivatives are considered “effective” if when stored in a deep freezer ( $0^\circ\text{F}$ ) the potency is 50 or better. We wish to test if our new drug is “effective”.

We have  $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x = 46 - \frac{19}{60}x$ . When  $\hat{y}(0) = \hat{\beta}_0 = 46$ . We construct the test: If the true  $\beta_0 = 50$ , given our sample observation  $\hat{\beta}_0 = 46 \dots$  is this rare or not.

$$H_0 : \hat{\beta}_0 = 50 \quad H_a : \hat{\beta}_0 < 50$$

Again,  $S^2 \sim \chi^2(10), \dots$  we need to use t-test. We are assuming  $\hat{\beta}_0 \sim N(\beta_0, \text{Var}(\beta_0)) \approx N(\beta_0, \text{Var}(\hat{\beta}_0))$ . p-value:

$$\begin{aligned}p &= \Pr(\hat{\beta}_0 \leq 46 \mid \hat{\beta}_0 = 50) \\ &= \Pr\left(\frac{\hat{\beta}_0 - \beta_0}{\sqrt{\text{Var}(\hat{\beta}_0)}} \leq \frac{46 - 50}{\sqrt{\text{Var}(\hat{\beta}_0)}}\right)\end{aligned}$$

Pausing to compute some of these

$$\text{Var}(\hat{\beta}_0) = \frac{\sigma^2 \sum x_i^2}{n S_{xx}} \sim \frac{S^2 \sum x_i^2}{n S_{xx}}$$

Here  $S^2 = 19.03$ ,  $n = 12$ , and  $S_{xx} = 6000$ .

$\sum x_i^2 = 3 \cdot (30^2 + 50^2 + 70^2 + 90^2) = 49200$  and  $\text{Var}(\hat{\beta}_0) = \frac{19.03 \cdot 49200}{12 \cdot 6000} = 13.004$  and  $\sqrt{\text{Var}(\hat{\beta}_0)} = 3.606$ . Then,

$$p = \Pr\left(T \leq \frac{46 - 50}{3.606}\right) = P(T \leq -1.09) = 0.150641$$

Cannot reject the null hypothesis! So maybe it *is* actually 50. Let's sell it!

## 11.6 Predictions via least squares regression line

Using  $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$ . Our emphasis so far has been on the coefficients  $\hat{\beta}_0$  and  $\hat{\beta}_1$  and their distributions. The goal is to understand  $y$  as a function of  $x$ .  $y$  is the exact, deterministic expectation whereas  $\hat{y}$  is our best attempt to approximate something that isn't exact.

**Example:** (Potency Example Again)

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x = 46 - \frac{19}{60}x$$

If we store the antibiotic at temperature  $x = 20^{\circ}\text{F}$ , what would we expect the potency to be?

$$\hat{y} = 46 - \frac{19}{60}(20) = 39.\bar{6}$$

But, what does this really mean? This is the expected value of  $y$  when  $x = 20$ . If we store a bunch of samples at  $20^{\circ}\text{F}$ , we would expect the sample mean to be approximately  $39.\bar{6}$ . In other words,  $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$  is a point estimator of  $E(Y)$ .

$\hat{y}$  is a statistic of its own. And, where there exists point estimators, there exists interval estimators (confidence intervals). We have that

$$\hat{\beta}_i \sim N(\beta_i, \text{Var}[\beta_i]).$$

So,  $\hat{y}$  is also normal for every fixed value of  $x$ .

To avoid formula confusion, we use  $x^*$  for the fixed  $x$ . Note:

$$E(\hat{y}(x^*)) = E\left(\hat{\beta}_0 + \hat{\beta}_1 x^*\right) = E\left(\hat{\beta}_0\right) + x^* E\left(\hat{\beta}_1\right) = \beta_0 = \beta_1 x^*.$$

We now need the variance,  $\text{Var}[\hat{y}]$ .

$$\begin{aligned} \text{Var}[\hat{y}] &= \text{Var}\left[\hat{\beta}_0 + \hat{\beta}_1 x^*\right] \\ &= \text{Var}\left[\hat{\beta}_0\right] + (x^*)^2 \text{Var}\left[\hat{\beta}_1\right] + 2x^* \text{Cov}(\hat{\beta}_0, \hat{\beta}_1) \\ &= \frac{\sigma^2 \sum x_i^2}{n S_{xx}} + (x^*)^2 \frac{\sigma^2}{S_{xx}} + 2x^* \left( \frac{-\bar{x}\sigma^2}{S_{xx}} \right) \\ &= \frac{\sigma^2}{S_{xx}} \left( \frac{\sum x_i^2 + n(x^*)^2 - 2x^*\bar{x}n}{n} \right) \end{aligned}$$

$$\begin{aligned} \text{Note: } S_{xx} &= \sum x_i^2 - n\bar{x}^2 \\ &= \frac{\sigma^2}{S_{xx}} \left( \frac{S_{xx} + n\bar{x}^2 + n(x^*)^2 - 2x^*\bar{x}n}{n} \right) \\ &= \frac{\sigma^2}{S_{xx}} \left( \frac{S_{xx} + n[\bar{x}^2 + (x^*)^2 - 2x^*\bar{x}]}{n} \right) \\ &= \sigma^2 \left( \frac{1}{n} + \frac{(x^* - \bar{x})^2}{S_{xx}} \right) \end{aligned}$$

When estimating  $\sigma^2$ , make sure to use the new one:

$$S^2 = \frac{1}{n-2} \text{SSE}.$$

Recap:  $\hat{y}$  estimator:

$$E(\hat{y}) = \beta_0 + \beta_1 x^*$$

and

$$\text{Var}[\hat{y}] = \sigma^2 \left( \frac{1}{n} + \frac{(x^* - \bar{x})^2}{S_{xx}} \right).$$

Standardize (as always):

$$\frac{\hat{\theta} - \theta}{\sigma_{\hat{\theta}}} = \frac{\hat{y}(x^*) - (\beta_0 + \beta_1 x^*)}{S \sqrt{\frac{1}{n} + \frac{(x^* - \bar{x})^2}{S_{xx}}}}.$$

This is a test statistic with  $(n - 2)$  degrees of freedom (df).

It follows that  $\alpha(1 - \alpha)$  level confidence interval is given by:

$$(\hat{\beta}_0 + \hat{\beta}_1 x^*) \pm t_{\alpha/2}(n - 2) S \sqrt{\frac{1}{n} + \frac{(x^* - \bar{x})^2}{S_{xx}}}.$$

**Example:** Find a 90% confidence interval for the average potency when stored at 20°F.

$$x^* = 20$$

$$\hat{y}(20) = 39.6$$

$$n = 12 \implies \text{so d.f.} = 10$$

$$\bar{x} = 60$$

$$S^2 = \frac{571}{30}$$

$$S_{xx} = 6000$$

**Answer:**

$$39.6 \pm \underbrace{(1.812)}_{\text{table}} \sqrt{\underbrace{\frac{571}{30}}_{S^2} \left( \frac{1}{12} + \frac{(20 - 60)^2}{6000} \right)}$$

$$39.6 \pm 4.677$$

**Remark:** We can now talk about hypothesis testing.

For example,

$$H_0 : y(x^*) = y_0 \quad \text{or} \quad H_a : y(x^*) \neq y_0$$

Our t-stat

$$\mathcal{T} = \frac{y - 0 - (\hat{\beta}_0 + \hat{\beta}_1 x^*)}{S \sqrt{\frac{1}{n} + \frac{(x^* - \bar{x})^2}{S_{xx}}}}$$

## 11.7 Predictions on $y$

The difference between this section and the previous is perspective. In §11.6 we focused on the spread of the average of  $y$ . In other words,  $E(\hat{y})$ .

What instead, we wanted to focus on the distribution of the values that  $y$  can take when  $x = x^*$ . Recall:

$$y = \underbrace{\beta_0 + \beta_1 x}_{\text{deterministic}} + \underbrace{\varepsilon}_{\text{"spread" of } y}$$

and  $\text{Var}(y) = \text{Var}(\varepsilon) = \sigma^2$ . Our estimator for  $y$  at  $x^*$  is defined

$$\hat{y}^* = \hat{\beta}_0 + \hat{\beta}_1 x^* + \varepsilon.$$

Since  $\varepsilon \sim N(0, \sigma^2)$ .

$$\begin{aligned} E(\hat{y}^*) &= E(\hat{\beta}_0 + \hat{\beta}_1 x^*) + E(\varepsilon) \\ &= E(\hat{y}(x^*)) \\ &= \beta_0 + \beta_1 x^* \end{aligned}$$

Computing variance is surprisingly easy. As before,  $\hat{y}^*$  is a sum of Normal random variables. Hence,  $\hat{y}^*$  is Normally distributed.

Moreover, our assumption on  $\varepsilon$  is that it is independent of  $x$ .

We have that  $S^2$  is independent of  $\hat{\beta}_0$  and  $\hat{\beta}_1$ . So

$$\begin{aligned} \text{Var}(\hat{y}^*) &= \text{Var}(\hat{\beta}_0 + \hat{\beta}_1 x^*) + \text{Var}(\varepsilon) && (\text{by independence}) \\ &= \sigma^2 \left( \frac{1}{n} + \frac{(x^* - \bar{x})^2}{S_{xx}} \right) + \sigma^2 \\ &= \sigma^2 \left( 1 + \frac{1}{n} + \frac{(x^* - \bar{x})^2}{S_{xx}} \right) \end{aligned}$$

This is the only difference from last sections  $\text{Var}(\hat{y})$ .

When  $\sigma^2$  unknown.

$$\text{Var}(\hat{y}^*) = S^2 \left( 1 + \frac{1}{n} + \frac{(x^* - \bar{x})^2}{S_{xx}} \right)$$

For the  $1 - \alpha$  level confidence interval (two sided)

$$\begin{aligned} \Pr(|T| > t_{\alpha/2}(\text{df})) &= 1 - \alpha \\ \mathcal{Z} &= \frac{y^* - E(\hat{y}^*)}{\sigma_{\hat{y}^*}} \end{aligned}$$

or

$$\mathcal{T} = \frac{y^* - (\hat{\beta}_0 + \hat{\beta}_1 x^*)}{S \sqrt{1 + \frac{1}{n} + \frac{(x^* - \bar{x})^2}{S_{xx}}}}$$

t-distributed with  $n - 2$  degrees of freedom and  $1 - \alpha$  level confidence interval is

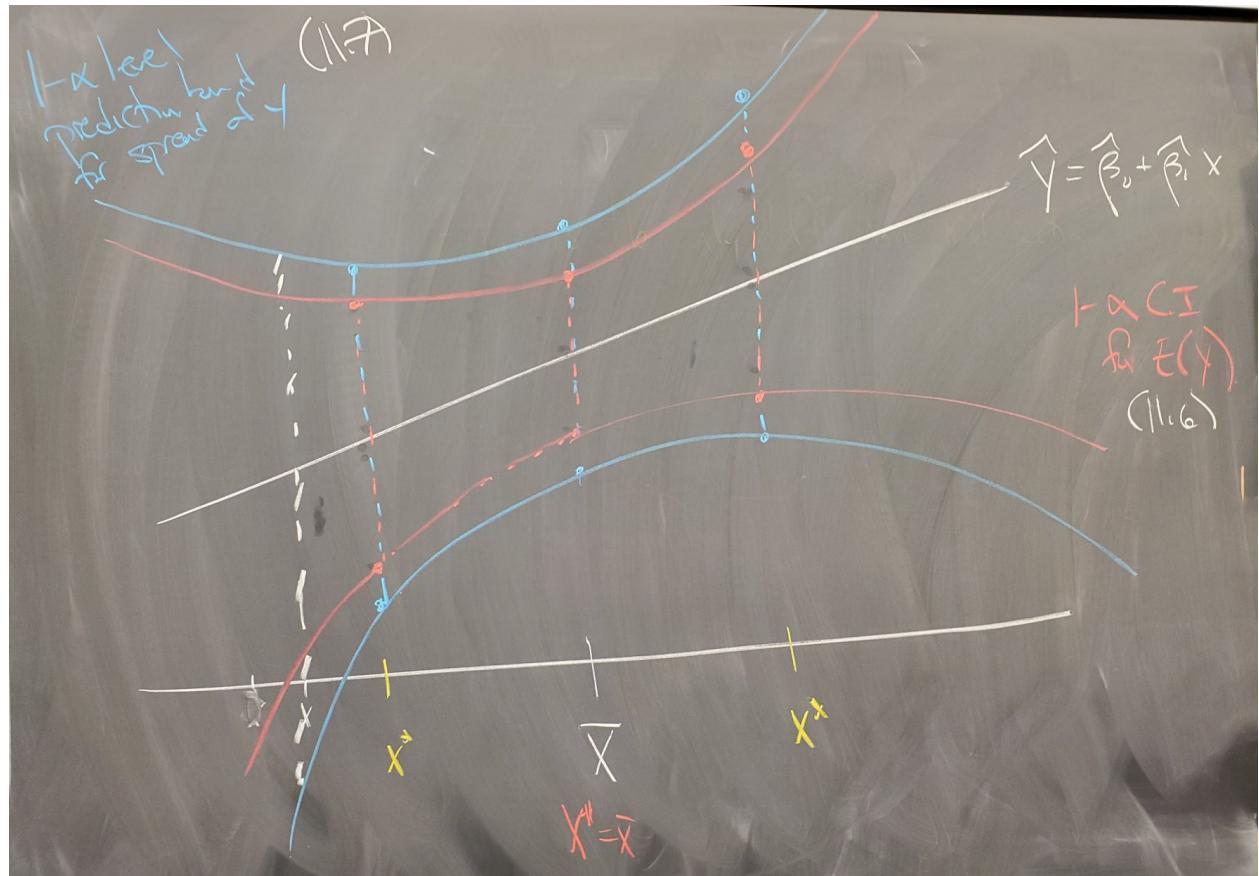
$$(\hat{\beta}_0 + \hat{\beta}_1 x^*) \pm t_{\alpha/2}(\text{df}) S \sqrt{1 + \frac{1}{n} + \frac{(x^* - \bar{x})^2}{S_{xx}}}$$

**Discussion:** Confidence and prediction bands.

In §11.6, we had C.I. for the mean of  $\hat{y}(x^*)$ .

In §11.7, we have C.I. for the spread of  $y$  at  $x^*$ ,

In either case, as  $x^*$  moves away from  $\bar{x}$ , the  $\frac{(x^* - \bar{x})^2}{S_{xx}}$  increases as does the standard error.  
i.e. the intervals get longer.



**Example:** Potency of antibiotic when stored at 20°F. Last day we showed that  $E(\hat{y}(20))$  had 90% confidence interval

$$\underbrace{39.6 \pm 4.677}_{\text{red cross-section}}$$

What spread of potency would we expect to see with 90% confidence if stored at 20°F. The same calculation as the last but with new standard error

$$39.6 \pm 1.812 \sqrt{1 + \frac{1}{12} + \frac{1600}{6000}}$$

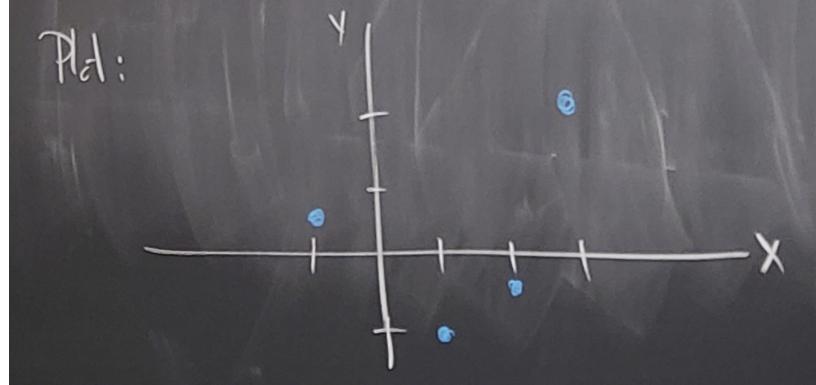
$$\underbrace{39.6 \pm 5.069}_{\text{blue cross-section}}$$

## 11.10 Multiple Linear Regression

**Example:** Consider the data:

$x$	-1	1	2	3
$y$	0.5	-1	-0.5	2

The plot does not look linear:



Seems unlikely that a line will fit well. What about a parabola? Want  $y = \beta_0 + \beta_1 x + \beta_2 x^2$ . Of course, in the probabilistic setup:

$$y = \underbrace{\beta_0 + \beta_1 x + \beta_2 x^2}_{\text{deterministic}} + \underbrace{\varepsilon}_{\text{error}}$$

Using the data:

$$\begin{aligned} \frac{1}{2} &= \beta_0 + \beta_1(-1) + \beta_2(-1)^2 \\ -1 &= \beta_0 + \beta_1(1) + \beta_2(1)^2 \\ -\frac{1}{2} &= \beta_0 + \beta_1(2) + \beta_2(2)^2 \\ 2 &= \beta_0 + \beta_1(3) + \beta_2(3)^2 \end{aligned}$$

But this is the same idea as before, just in higher dimensions (3 coefficients).

$$\begin{aligned} \vec{y} &= \beta_0 \vec{1} + \beta_1 \vec{x} + \beta_2 \vec{x^2} \\ &= \beta_0 \begin{pmatrix} 1 \\ 1 \\ 1 \\ 1 \end{pmatrix} + \beta_1 \begin{pmatrix} -1 \\ 1 \\ 2 \\ 3 \end{pmatrix} + \beta_2 \begin{pmatrix} 1 \\ 1 \\ 4 \\ 9 \end{pmatrix} \end{aligned}$$

Or, as a matrix equation,

$$\begin{bmatrix} 1 & -1 & 1 \\ 1 & 1 & 1 \\ 1 & 2 & 4 \\ 1 & 3 & 9 \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \end{bmatrix} = \begin{bmatrix} 1/2 \\ -1 \\ -1/2 \\ 2 \end{bmatrix}.$$

Remark: This is why this process is called **linear** regression.

The system for the coefficients will be a linear system. It has **nothing** to do with the underlying model... Which here is quadratic  $y = \beta_0 + \beta_1 x + \beta_2 x^2$ .

Just as before, we have an overdetermined system, which implies there is no solution. Hence, the best we can do is find the projection.

We need a way to solve this projection problem in general.

### Topic: The Normal Equations:

In complete generality,

$$y = \beta_0 + \beta_1 f_1(x) + \beta_2 f_2(x) + \cdots + \beta_n f_n(x) + \varepsilon.$$

Let there be  $m$  data points  $(x_i, y_i)$ .

This yields a matrix equation:

$$\begin{bmatrix} 1 & f_1(x_1) & f_2(x_1) & \cdots & f_n(x_1) \\ 1 & f_1(x_2) & f_2(x_2) & \cdots & f_n(x_2) \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & f_1(x_n) & f_2(x_n) & \cdots & f_n(x_n) \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_n \end{bmatrix} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

Write  $\mathbf{X}\vec{\beta} = \vec{y}$ . Then  $\mathbf{X}$  is  $m \times n$ .

Note, this setup really only makes sense when the system is overdetermined (more equations than unknowns). That is, more rows than columns in  $\mathbf{X}$ . i.e.  $m > n$ . We seek the coefficient vectors  $\hat{\beta} = [\hat{\beta}_1, \dots, \hat{\beta}_n]$  such that  $\|\vec{y} - \mathbf{X}\hat{\beta}\|$  is minimized (LSR).

In the general case, there are 2 new issues to contend with:

- ① no reason the column vectors of  $\mathbf{X}$  are actually a basis for  $\text{Col } \mathbf{X}$  (*might just be a spanning set*).
- ② Certainly no reason columns are an orthogonal set.

The implication of ① is that a least squares solution  $\hat{\beta}$  need not be unique. As for ②, ends we don't need this. By definition  $\vec{y} - \mathbf{X}\hat{\beta}$  is orthogonal to  $\text{Col } \mathbf{X}$ .

In other words,  $\vec{y} - \mathbf{X}\hat{\beta}$  must lie in the null space of  $\mathbf{X}^T$ . (Recall  $(\text{Col})^\perp = \text{null } A^T$ ).

Thus,  $\mathbf{X}^T(\vec{y} - \mathbf{X}\hat{\beta}) = \vec{0}$ .

$$\begin{aligned} &\implies \underbrace{\mathbf{X}^T \vec{y}}_{m \text{ vector}} - \underbrace{\mathbf{X}^T \mathbf{X} \hat{\beta}}_{m \text{ vector}} = \underbrace{\vec{0}}_m \\ &\implies \mathbf{X}^T \mathbf{X} \hat{\beta} = \mathbf{X}^T \vec{y} \end{aligned}$$

Definition: The normal equations:

$$\mathbf{X}^T \mathbf{X} \hat{\beta} = \mathbf{X}^T \vec{y}$$

Facts:

- ①  $\mathbf{X}^T \mathbf{X}$  is an  $n \times n$  matrix.

- ② The normal equation is always a consistent system. That is, there is a solution to the least squares problem.
- ③ If the columns of  $\mathbf{X}$  are linearly independent, then  $\mathbf{X}^T \mathbf{X}$  is an invertible square matrix and the unique solution to the least squares problem is

$$\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T y$$

**Example:** Our best fit parabola.

$$\mathbf{X} = \begin{bmatrix} 1 & -1 & 1 \\ 1 & 1 & 1 \\ 1 & 2 & 4 \\ 1 & 3 & 9 \end{bmatrix}, \quad \mathbf{X}^T = \begin{bmatrix} 1 & 1 & 1 & 1 \\ -1 & 1 & 2 & 3 \\ 1 & 1 & 4 & 9 \end{bmatrix}$$

$$\hat{\beta} = [\hat{\beta}_0, \hat{\beta}_1, \hat{\beta}_2], \quad y = \left[ \frac{1}{2}, -1, -\frac{1}{2}, 2 \right]$$

$$\mathbf{X}^T \mathbf{X} = \begin{bmatrix} 4 & 5 & 15 \\ 5 & 15 & 35 \\ 15 & 35 & 99 \end{bmatrix}$$

Note that  $\det(\mathbf{X}^T \mathbf{X}) = 440 \neq 0$  so it is invertible.

$$\begin{aligned} \hat{\beta} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T y \\ &= \begin{bmatrix} -41/44 \\ -379/440 \\ 53/88 \end{bmatrix} \end{aligned}$$

So, the best-fit parabola is

$$\begin{aligned} y &= -\frac{41}{44} - \frac{379}{440}x + \frac{53}{88}x^2 \\ &\approx -0.932 - 0.861x + 0.602x^2 \end{aligned}$$

**Discussion:** Revisit the best-fit line

Given  $n$   $(x_i, y_i)$ 's,  $y = \beta_0 + \beta_1 x + \varepsilon$ .

$$\Rightarrow [\vec{1} \quad \vec{x}] \hat{\beta} = \vec{y} \quad \text{and} \quad \hat{\beta} = \begin{bmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \end{bmatrix}$$

Here  $\mathbf{X} = [\vec{1} \quad \vec{x}]$  is a  $n \times 2$  matrix and  $\mathbf{X}^T = \begin{bmatrix} 1^T \\ x^T \end{bmatrix}$  is  $2 \times n$ .

Then

$$\begin{aligned} \mathbf{X}^T \mathbf{X} &= \begin{bmatrix} 1^T \\ x^T \end{bmatrix} [\vec{1} \quad \vec{x}] \\ &= \begin{bmatrix} 1 \cdot 1 & 1 \cdot x \\ x \cdot 1 & x \cdot x \end{bmatrix} \\ &= \begin{bmatrix} n & \sum x_i \\ \sum x_i & \sum x_i^2 \end{bmatrix} \end{aligned}$$

Normal equation  $X^T X \hat{\beta} = X^T y$ . Can we solve uniquely?

$$\begin{aligned}\det(X^T X) &= n \sum x_i^2 - \left(\sum x_i\right)^2 \\ &= n \sum x_i^2 - n^2 \bar{x}^2 \\ &= n \left(\sum x_i^2 - n \bar{x}\right) \\ &= n S_{xx} \\ &> 0\end{aligned}$$

Therefore  $X^T X$  is invertible! ( $\det A \neq 0 \iff A^{-1}$  exists)

So  $\hat{\beta} = (X^T X)^{-1} X^T y$

$$\text{Recall: } \begin{bmatrix} a & b \\ c & d \end{bmatrix}^{-1} = \frac{1}{\det A} [( \operatorname{tr} A ) I - A] = \frac{1}{ad - bc} \begin{bmatrix} d & -b \\ -c & a \end{bmatrix}.$$

So

$$(X^T X)^{-1} = \frac{1}{n S_{xx}} \begin{bmatrix} \sum x_i^2 & -\sum x_i \\ -\sum x_i & n \end{bmatrix}$$

and

$$X^T y = \begin{bmatrix} 1^T \\ x^T \end{bmatrix} y = \begin{bmatrix} 1 \cdot y \\ x \cdot y \end{bmatrix} = \begin{bmatrix} \sum y_i \\ \sum x_i y_i \end{bmatrix}$$

Then

$$\begin{aligned}\hat{\beta} &= \frac{1}{n S_{xx}} \begin{bmatrix} \sum x_i^2 & -\sum x_i \\ -\sum x_i & n \end{bmatrix} \begin{bmatrix} \sum y_i \\ \sum x_i y_i \end{bmatrix} \\ &= \frac{1}{n S_{xx}} \begin{bmatrix} (\sum x_i^2)(\sum y_i) - (\sum x_i)(\sum x_i y_i) \\ -(\sum x_i)(\sum y_i) + n \sum x_i y_i \end{bmatrix}\end{aligned}$$

Same  $\hat{\beta}_0, \hat{\beta}_1$  as earlier. But wait there's more..., Look at  $(X^T X)^{-1}$ :

$$\begin{aligned}(X^T X)^{-1} &= \frac{1}{n S_{xx}} \begin{bmatrix} \sum x_i^2 & -n \bar{x} \\ -n \bar{x} & n \end{bmatrix} \\ &= \begin{bmatrix} \frac{\sum x_i^2}{n S_{xx}} & -\frac{\bar{x}}{S_{xx}} \\ -\frac{\bar{x}}{S_{xx}} & \frac{1}{S_{xx}} \end{bmatrix}\end{aligned}$$

And

$$\sigma^2 (X^T X)^{-1} = \begin{bmatrix} \operatorname{Var}(\hat{\beta}_0) & \operatorname{Cov}(\hat{\beta}_0, \hat{\beta}_1) \\ \operatorname{Cov}(\hat{\beta}_0, \hat{\beta}_1) & \operatorname{Var}(\hat{\beta}_1) \end{bmatrix}$$

**FACT:** This always happens.  $\sigma^2 (X^T X)^{-1}$  is the table of covariances:

$$\sigma^2 - a_{ij} = \operatorname{Cov}(\hat{\beta}_i, \hat{\beta}_j) \quad \text{where} \quad [a_{ij}] = (X^T X)^{-1}$$

Lastly,

$$S^2 = \frac{\text{SSE}}{n-2} = \frac{\sum(y_i - \hat{y}_i)^2}{n-2}$$

And “by some matrix algebra” (Wackerly),

$$\text{SSE} = y^T y - \hat{\beta}^T \mathbf{X}^T y.$$

**Example:** (Antibiotic again)

$x$	30	30	30	50	50	50	70	70	70	90	90	90
$y$	38	43	29	32	26	33	19	24	23	14	19	21

$n = 12$  and  $\mathbf{X}^T \mathbf{X} = \begin{bmatrix} n & \sum x_i \\ \sum x_i & \sum x_i^2 \end{bmatrix} = \begin{bmatrix} 12 & 720 \\ 720 & 49,200 \end{bmatrix}$ . Then  $\det(\mathbf{X}^T \mathbf{X}) = 72,000$ . Then

$$(\mathbf{X}^T \mathbf{X})^{-1} = \frac{1}{72,000} \begin{bmatrix} 49,200 & -700 \\ -700 & 12 \end{bmatrix}.$$

**Note:**  $\text{Var}(\hat{\beta}_0) = \frac{41}{60}\sigma^2$ , and  $\text{Var}(\hat{\beta}_1) = \frac{\sigma^2}{6000}$ .

Then  $\mathbf{X}^T y = \begin{bmatrix} \sum y_i \\ \sum x_i y_i \end{bmatrix} = \begin{bmatrix} 324 \\ 17,540 \end{bmatrix}$  and  $\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T y = \begin{bmatrix} 46 \\ -19/60 \end{bmatrix}$  (same as before).

Hence

$$\hat{y} = 46 - \frac{19}{60}x$$

$$\text{SSE} = y^T y - \hat{\beta}^T \mathbf{X}^T y = \frac{2961}{15}$$

$$\sigma^2 \approx \frac{\text{SSE}}{n-2} = \frac{2971/15}{10} \approx 19.8067.$$

## 11.11 A big ol theorem

Theorem: Let  $Y_i = \beta_0 + \beta_1 f_1(x_i) + \cdots + \beta_k f_k(x_i) + \varepsilon_i$  where  $\varepsilon_i \sim N(0, \sigma^2)$  (with  $E[\varepsilon] = 0$ ).

Then the least-squares estimates are given by  $\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T y$  provided  $(\mathbf{X}^T \mathbf{X})^{-1}$  exists. Then

- ①  $E(\hat{\beta}_i) = \hat{\beta}_i$  (unbiased)
- ②  $\text{Var}(\hat{\beta}_1) = c_{ii}\sigma^2$  where  $(\mathbf{X}^T \mathbf{X})^{-1} = [c_{ij}]$
- ③  $\text{Cov}(\hat{\beta}_i, \hat{\beta}_j) = c_{ij}\sigma^2$
- ④  $S^2 = \frac{\text{SSE}}{n-(k+1)}$  with  $k+1$   $\hat{\beta}_i$ 's and  $\text{SSE} = y^T y - \hat{\beta}^T \mathbf{X}^T y$  and  $E(S^2) = \sigma^2$  (unbiased)
- ⑤  $\hat{\beta}_i$  is normally distributed.
- ⑥  $\frac{(n-(k+1))S^2}{\sigma^2} \sim \chi^2(n-(k+1))$
- ⑦  $S^2$  and  $\hat{\beta}_1$  are independent for all  $i$ .

## 11.12 Hypothesis Testing C.I.

Motivation: Testing a specific  $\hat{\beta}_i$ .

We have  $\hat{\beta} = (X^T X)^{-1} T y = [\hat{\beta}_0, \dots, \hat{\beta}_k]$ . To pick off a specific  $\hat{\beta}_i$ , we use the dot product.

$$\text{Let } e_i \text{ (standard basis vector), then } \hat{\beta}_i = \underbrace{e_i \cdot \hat{\beta}}_{\text{dot product}} = \underbrace{e_i^T \hat{\beta}}_{\text{matrix multiplication}}$$

Note:  $E(e_i \cdot \hat{\beta}) = E(1 \cdot \hat{\beta}_i) = E(\hat{\beta}_i) = \beta_i$ .

Same with  $\text{Var}(e_i \cdot \hat{\beta}) = \text{Var}(\hat{\beta}_i)$ .

So a test of the form

$$H_0 : \hat{\beta}_i = (\beta_i)_0$$

$$H_a : \hat{\beta}_i \neq (\beta_i)_0$$

yields same old statistics.

$$Z = \frac{\hat{\beta}_i - (\beta_i)_0}{\sqrt{\text{Var}(\beta_i)}} = \frac{\hat{\beta}_i - (\beta_i)_0}{\sqrt{c_{ii}\sigma^2}} \quad \text{with} \quad [c_{ii}] = (X^T X)^{-1}$$

$$T = \frac{\hat{\beta}_i - (\beta_i)_0}{S\sqrt{c_{ii}}} \quad \text{with} \quad n - (k + 1) \text{ df.}$$

**Example:** We fit

x	-1	1	2	3
y	0.5	-1	-0.5	2

with  $y = \beta_0 + \beta_1 x + \beta_2 x^2$  because it “looked” non-linear.



Is there evidence that the data is in fact non-linear? That is, are we reasonably certain  $\beta_\alpha \neq 0$ .

$$H_0 : \beta_\alpha = 0$$

$$H_a : \beta_\alpha \neq 0$$

We had

$$\begin{aligned}\hat{\beta} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T y \\ &= \begin{bmatrix} -41/44 \\ -379/440 \\ 53/88 \end{bmatrix} \\ &= \begin{bmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \\ \hat{\beta}_2 \end{bmatrix}\end{aligned}$$

$e_3 = (0, 0, 1)$  and  $e_3^T \hat{\beta} = \frac{53}{80}$ . We had

$$\mathbf{X}^T \mathbf{X} = \begin{bmatrix} 4 & 5 & 15 \\ 5 & 15 & 35 \\ 15 & 35 & 99 \end{bmatrix}$$

and  $\det(\mathbf{X}^T \mathbf{X}) = 440$ . Then

$$(\mathbf{X}^T \mathbf{X})^{-1} = \frac{1}{440} \begin{bmatrix} 260 & 30 & -50 \\ 30 & 171 & -65 \\ -50 & -65 & 35 \end{bmatrix}$$

With  $\text{Var}(\hat{\beta}_2) = c_{22} \sigma^2 = \frac{35}{440} \sigma^2$ .

For  $\sigma^2$ , we need  $S^2 = \frac{\text{SSE}}{n-2}$ .

$$\text{SSE} = y^T y - \hat{\beta}^T X^T y = \frac{4791}{440}.$$

$$S^2 = \frac{\frac{4791}{440}}{n-2} = \frac{4791}{850} \approx 5.44432.$$

$$\text{Var}(\hat{\beta}_2) \approx c_{22} S^2 = \sqrt{\text{Var}(\hat{\beta}_2)} = 0.690201.$$

So

$$\begin{aligned}T &= \frac{e_3 \cdot \hat{\beta} - (\hat{\beta}_2)_0}{\sqrt{\text{Var}(\hat{\beta}_2)}} \\ &= \frac{\frac{53}{88} - 0}{0.690201} \\ &= 0.872604\end{aligned}$$

$$\Pr(|T| \geq 0.872604 \mid \text{df} = 2) = 0.4748909.$$

This sucks! It supports the null Hypothesis.

**Example:** (Potency again)

$x$	30	30	30	50	50	50	70	70	70	90	90	90
$y$	38	43	29	32	26	33	19	24	23	14	19	21

I tried quadratic  $y = \beta_0 + \beta_1 x + \beta_2 x^2$ .

$$X = [1 \ x \ x^2]$$

$$(X^T X)^{-1} = \frac{1}{1.92 \times 10^6} \begin{bmatrix} 1.092 \times 10^6 & -391200 & 3100 \\ -391200 & 14720 & -120 \\ 3100 & -120 & 1 \end{bmatrix}$$

Note:  $\text{Var}(\hat{\beta}_2) = \frac{1}{1.92 \times 10^6} \sigma^2 \approx \frac{S^2}{1920000}$ .

$$\text{SSE} = y^T y + \hat{\beta}^T X^T y = 189.$$

$$S^2 = \frac{\text{SSE}}{n-2} = \frac{\text{SSE}}{10} = 18.9$$

$$\text{Var}(\hat{\beta}_2) = 9.84775 \times 10^{-6}.$$

$$\sqrt{\text{Var}(\hat{\beta}_2)} = 0.00313748$$

Testing

$$\begin{aligned} H_0 : \beta_2 &= 0 \\ H_a : \beta_2 &\neq 0 \end{aligned}$$

I computed  $\hat{\beta}_2 = \frac{1}{1200}$ .

$$t = \frac{\hat{\beta}_2 - (\beta_2)_0}{\sqrt{\text{Var}(\hat{\beta}_2)}} = 0.2656$$

With  $\text{df} = 10$ ,  $\Pr(|T| > 0.2656) = 0.79$ , which is HUGE again. Hence it supports the null hypothesis.

Of course, this is a math class and we love to generalize. We don't have to use a standard basis vector  $e_i$ . Let  $\vec{a}$  be a constant vector  $a = (a_0, \dots, a_k)$ . Then

$$a^T \hat{\beta} = a_0 \hat{\beta}_0 + \dots + a_k \hat{\beta}_k.$$

Note:

$$\begin{aligned} E(a^T \hat{\beta}) &= \sum i = 1^k a_i E(\hat{\beta}_i) && \text{(by linearity)} \\ &= \sum i = 1^k a_i \beta_i && \text{(by unbiased estimators)} \\ &= a^T \beta \end{aligned}$$

for  $\beta = (\beta_0, \dots, \beta_k)$ . Also,

$$\begin{aligned} \text{Var}(a^T \hat{\beta}) &= \text{Var}(a_0 \hat{\beta}_0 + \dots + a_k \hat{\beta}_k) \\ &= \sum i = 1^k \sum j = 1^k \text{Cov}(a_i \hat{\beta}_i, a_j \hat{\beta}_j) \\ &= \sum i = 1^k \sum j = 1^k a_i a_j \underbrace{\text{Cov}(\hat{\beta}_i, \hat{\beta}_j)}_{\text{These are } c_{ij} \sigma^2 \text{'s from } (X^T X)^{-1}} \end{aligned}$$

This entire double sum can be written as the matrix product

$$\begin{aligned}\text{Var}(a^T \hat{\beta}) &= a^T (X^T X)^{-1} a \cdot \sigma^2 \\ \implies Z &= \frac{a^T \hat{\beta} - \mathbb{E}(a^T \hat{\beta})}{\text{Var}(a^T \hat{\beta})} = \frac{a^T \hat{\beta} - a^T \beta}{\sigma \sqrt{a^T (X^T X)^{-1} a}}\end{aligned}$$

Given  $H_0 : a^T \beta = (a^T \beta)_0$  versus some alternative, use

$$T = \frac{a^T \hat{\beta} - (a^T \beta)_0}{S \sqrt{a^T (X^T X)^{-1} a}}$$

with  $n - (k + 1)$  degrees of freedom.

# Chapter 13

## One-way Analysis of Variance

We know how to test fit the difference in 2 means

$$\begin{aligned} H_0 : \mu_1 = \mu_2 &\quad \text{i.e.} \quad \mu_1 - \mu_2 = 0. \\ H_a : \mu_1 \neq \mu_2 &\quad \text{i.e.} \quad \mu_1 - \mu_2 \neq 0 \end{aligned}$$

Now we endeavor to construct a test to detect difference in means over multiple groups.

**Example:** 24 expert typists test three new keyboard designs. Randomly assign 8 to each type of keyboard and assigned the same document to type up. The time of the task is recorded with the idea that a group with significantly less time on task would imply a better keyboard design

design	times
KB1	364 366 394 386 379 398 371 370
KB2	355 359 374 342 378 355 376 358
KB3	360 345 374 390 386 373 393 366

An obvious stat to look at is means

$$\bar{A} = 378.5 \quad \bar{B} = 362.125 \quad \bar{C} = 373.375$$

$B$  “looks” better, but can we design a test to determine if it is “significantly” so? This is our goal.

The standard null hypothesis is that the true population means are the same for each group,  $H_0 : \mu_1 = \mu_2 = \dots = \mu_k = \mu_0$ . The alternative  $H_a$  is that at least one of the  $\mu_i$ 's are different.

Notation: Let  $i$  indicate the “group”. E.g.  $k_i$  with  $i = 1, 2, 3$ . ( $k$  for keyboard). Let  $j$  denote the data point in that group. i.e.  $Y_{13} = 394$ ,  $Y_{25} = 378$ .

Let  $n_i$  be the number of data points in the  $i^{\text{th}}$  group. Here  $n_1 = n_2 = n_3 = 8$ . Let  $n$  be the total of all data points.  $n = \sum_{i=1}^k n_i$ . Here  $n = 24$ .

Formally, we assume  $Y_{ij} \sim N(\mu_i, \sigma^2)$  where each group has mean  $\mu_i$ , but all groups have the same variance  $\sigma^2$ .

We need a stat.

To find “a” stat, we construct a likelihood ratio test. For  $H_0$  we have our usual MLEs. If there is only  $\mu_0$

$$\hat{\mu}_0 = \frac{1}{n} \sum_{i=1}^k \sum_{j=1}^{n_i} Y_{ij}$$

and

$$S_0^2 = \sum_{i=1}^k \sum_{j=1}^{n_i} (Y_{ij} - \hat{\mu}_0)^2.$$

Under the alternative hypothesis, for the individual means, we again use the standard MLE

$$\hat{\mu}_i = \frac{1}{n_i} \sum_{j=1}^{n_i} Y_{ij} = \bar{Y}_i.$$

Aside: In our motivational example,  $\bar{A} = \bar{Y}_1 = 378.5$ ,  $\bar{B} = \bar{Y}_2 = 362.125$ , etc.

But we are still assuming a unique  $\sigma^2$  for all groups. So the MLE for

$$S_a^2 = \frac{1}{n} \sum_{i=1}^k \sum_{j=1}^{n_i} (Y_{ij} - \hat{\mu}_i)^2.$$

Remark: Proving that this is an MLE is a good review and will be coming to a homework soon.

Construct a likelihood function:

$$\begin{aligned} L(\text{All } Y_{ij} \mid \hat{\mu}_0, S_0^2) &= \prod_{i=1}^k \prod_{j=1}^{n_i} \left( \frac{1}{2\pi S_0^2} \right)^{1/2} \exp \left( \frac{-(Y_{ij} - \hat{\mu}_0)^2}{2S_0^2} \right) \\ &= \left( \frac{1}{2\pi S_0^2} \right)^{n/2} \exp \left( -\frac{1}{2S_0^2} \underbrace{\sum_{i=1}^k \sum_{j=1}^{n_i} (Y_{ij} - \hat{\mu}_0)^2}_{nS_0^2} \right) \\ &= \left( \frac{1}{2\pi S_0^2} \right)^{n/2} \exp \left( -\frac{nS_0^2}{2S_0^2} \right) \\ &= \left( \frac{1}{2\pi S_0^2} \right)^{n/2} e^{-n/2} \end{aligned}$$

Under  $H_a$ , same “algebra” occurs.

$$L(\text{All } Y_{ij} \mid \hat{\mu}_0, S_a^2) = \left( \frac{1}{2\pi S_a^2} \right)^{n/2} e^{-n/2}$$

Then the ratio of likelihood functions yields

$$\frac{L(\text{All } Y_{ij} \mid H_0)}{L(\text{All } Y_{ij} \mid H_a)} = \frac{\left( \frac{1}{2\pi S_0^2} \right)^{n/2} e^{-n/2}}{\left( \frac{1}{2\pi S_a^2} \right)^{n/2} e^{-n/2}} = \left( \frac{S_a^2}{S_0^2} \right)^{n/2} < k$$

Implies the statistic  $T = \frac{S_a^2}{S_0^2}$  and small values of  $T$  support the alternative hypothesis (define our RR). This makes sense because if the true means are different then  $S_a^2$  will be the correct estimator for  $\sigma^2$  and  $S_0^2$  will be larger (on average).

We have a problem... This is an entirely new stat for us.

One thing we should show is that (a)

$$S_a^2 \text{ is an estimator of } \sigma^2$$

moreover, we need it unbiased.

Note that via multiplying by the correct degrees of freedom  $\nu$ , and dividing by  $\sigma^2$ , the numerator and denominator are  $\chi^2$  distributed. This lead to the advent of the  $F$ -distribution.

Topic: The  $F$ -distribution. Suppose  $X_1 \sim \chi^2(m)$  and  $X_2 \sim \chi^2(n)$  and are independent. Define

$$Y_i := \frac{X_1/m}{X_2/n}$$

to be an  $F$  random variable with  $m$  numerator degrees of freedom and  $n$  denominator degrees of freedom. We usually write  $Y_i \sim F(m, n)$ .

The derivation is akin to the derivation of the  $T$ -stat. We used the Jacobian method of transformations. (§6.6 in the text. This is usually skipped in MTH 325.)

Letting  $Y_2 = X_2$ , we can find the joint density function of  $Y_1$  and  $Y_2$ , then we integrate the joint density to get the marginal density of  $Y_i$ ...

$$f(y_1) = \frac{\Gamma\left(\frac{m+n}{2}\right)}{\Gamma\left(\frac{m}{2}\right)\Gamma\left(\frac{n}{2}\right)} \left(\frac{m}{n}\right)^{n/2} \frac{y_1^{m/2-1}}{\left(1 + \left(\frac{my_1}{n}\right)^{(m+n)/2}\right)} \text{ and } y_i > 0.$$

**Last day:** Looking at at designing a test  $H_0 : \mu_1 = \mu_2 = \dots = \mu_k = \mu_0$  versus  $H_a$  not all equal.

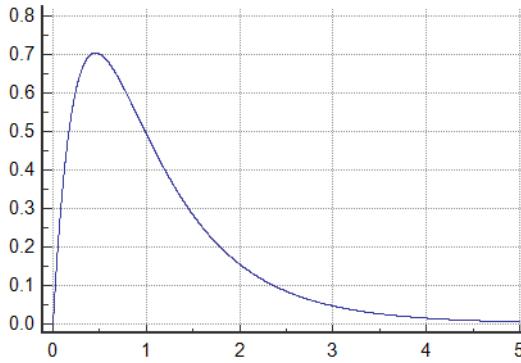
Theorem: Likelihood ration, find a new (for us) statistic that looked like a  $\frac{\chi^2 \text{ r.v.}}{\chi^2 \text{ r.v.}}$ . This lead us to the  $F$ -distribution.

$$Y \sim F(m, n) \quad \text{when} \quad Y = \frac{X_1/m}{X_2/n}$$

and

$$X_1 \sim \chi^2(m) \quad \text{and} \quad X_2 \sim \chi^2(n).$$

Graph  $f(y), y > 0$ :



For  $E(Y)$  and  $\text{Var}(Y)$ .

$$\begin{aligned}
 E(Y) &= E\left(\frac{X_1/m}{X_2/n}\right) \\
 &= E\left(\frac{n}{m} \cdot \frac{X_1}{X_2}\right) \\
 &= \frac{n}{m} E\left(X_1 \cdot \frac{1}{X_2}\right) \\
 &= \frac{n}{m} E(X_1) E\left(\frac{1}{X_2}\right) && (\text{by independence}) \\
 &= \frac{n}{m} \cdot m E\left(\frac{1}{X_2}\right) \\
 &= n E\left(\frac{1}{X_2}\right) && (\text{numerator df has no impact}) \\
 &= n \cdot \underbrace{\frac{1}{n-2}}_{\text{by thm.}} \\
 E(Y) &= \frac{n}{n-2}.
 \end{aligned}$$

Also, with proof omitted,

$$\text{Var}(Y) = \frac{2n^2(m+n-2)}{n(n-2)^2(n-4)}.$$

End  $F$ -distribution background. Then,

$$T = \frac{S_0^2}{S_a^2} = \frac{\frac{1}{n} \sum_{i=1}^k \sum_{j=1}^{n_i} (Y_{ij} - \hat{\mu}_0)^2}{\frac{1}{n} \sum_{i=1}^k \sum_{j=1}^{n_i} (Y_{ij} - \hat{\mu}_i)^2}$$

and the  $\frac{1}{n}$  can cancel. If  $S_0^2$  and  $S_a^2$  were independent, then dividing by their degrees of freedom would yield an  $F$ -distribution. Sadly, they aren't. But we can algebra to

independent ratios. Start with the top,

$$\begin{aligned}
S_0^2 &= \sum_{i=1}^k \sum_{j=1}^{n_i} (Y_{ij} - \hat{\mu}_0)^2 \\
&= \sum_{i=1}^k \sum_{j=1}^{n_i} ((Y_{ij} - \hat{\mu}_i) + (\hat{\mu}_i - \hat{\mu}_0))^2 \\
&= \sum_{i=1}^k \sum_{j=1}^{n_i} [(Y_{ij} - \hat{\mu}_i)^2 + 2(Y_{ij} - \hat{\mu}_i)(\hat{\mu}_i - \hat{\mu}_0) + (\hat{\mu}_i - \hat{\mu}_0)^2] \\
&= \sum_{i=1}^k \sum_{j=1}^{n_i} (Y_{ij} - \hat{\mu}_i)^2 + 2 \sum_{i=1}^k [(\hat{\mu}_i - \hat{\mu}_0) \underbrace{(Y_{ij} - \hat{\mu}_i)}_{=0 \text{ as it's the sum of all deviations for mean within a group from group mean}}] + \sum_{i=1}^k n_i (\hat{\mu}_i - \hat{\mu}_0)^2 \\
&= \sum_{i=1}^k \sum_{j=1}^{n_i} (Y_{ij} - \hat{\mu}_i)^2 + \sum_{i=1}^k n_i (\hat{\mu}_i - \hat{\mu}_0)^2 \\
&= S_a^2 + \sum_{i=1}^k n_i (\hat{\mu}_i - \hat{\mu}_0)^2
\end{aligned}$$

Then

$$\frac{S_0^2}{S_a^2} = 1 + \frac{\sum_{i=1}^k n_i (\hat{\mu}_i - \hat{\mu}_0)^2}{\sum_{i=1}^k \sum_{j=1}^{n_i} (Y_{ij} - \hat{\mu}_i)^2}$$

We are pretty close to independence. Recall we have Fisher's Theorem:  $\chi^2 \sim N(\mu, \sigma^2)$  then  $\bar{X}$  and  $S^2$  are independent random variables and  $\bar{X} \sim N(\mu, \sigma^2/n)$ .

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

Thus the  $\hat{\mu}_i$ 's up top are independent of all the variances  $S_i^2$  in sum of the denominator. Also,  $\hat{\mu}_0$  is a function of  $\hat{\mu}_i$

$$\hat{\mu}_0 = \frac{n_1 \hat{\mu}_1 + n_2 \hat{\mu}_2 + \cdots + n_k \hat{\mu}_k}{n}$$

And the entire numerator is independent of the variance terms in the numerator when  $H_0$  is true. Now, when  $H_0$  is true, all the resultant variable (top and bottom) are  $\chi^2$  distributed. The last thing to do is figure out degrees of freedom. Given  $\mu_0$ ,

$$\frac{\hat{\mu}_1 - \mu_0}{\sigma / \sqrt{n_1}} \sim N(0, 1) \quad \text{and} \quad \frac{n_1 (\hat{\mu}_1 - \mu_0)^2}{\sigma^2} \sim \chi^2(1)$$

So, the sum of  $\chi^2$  random variables is also a  $\chi^2$  rv and the degrees of freedom add. Thus

$$\frac{1}{\sigma^2} \sum_{i=1}^k n_i (\hat{\mu}_i - \mu_0)^2 \sim \chi^2(k).$$

On the other hand, to use these facts, use the say “trick” again

$$\widehat{\mu}_1 - \mu_0 = (\widehat{\mu}_1 - \widehat{\mu}_0) + (\widehat{\mu}_0 - \mu_0)$$

And

$$\sum_{i=1}^k n_i(\widehat{\mu}_i - \mu_0)^2 = \underbrace{\sum_{i=1}^k n_i(\widehat{\mu}_i - \widehat{\mu}_0)^2}_{\text{the numerator}} + n \underbrace{(\widehat{\mu}_0 - \mu_0)^2}_{\substack{\text{divide by } \sigma^2, \\ \text{this is } \chi^2(1)}}$$

In the end

$$\underbrace{\frac{1}{\sigma^2} \sum_{i=1}^k n_i(\widehat{\mu}_i - \mu_0)^2}_{\chi^2(k)} = \underbrace{\frac{1}{\sigma^2} \sum_{i=1}^k n_i(\widehat{\mu}_i - \widehat{\mu}_0)^2}_{\substack{\text{the numerator is} \\ \chi^2(k-1)}} + n \underbrace{\frac{(\widehat{\mu}_0 - \mu_0)^2}{\sigma^2}}_{\chi^2(1)}$$

Denominator is much easier,

$$\sum_{i=1}^k \underbrace{\sum_{j=1}^{n_i} \frac{(Y_{ij} - \mu_i)^2}{\sigma^2}}_{S_i^2} \sim \chi^2(n_i - 1)$$

and by independences, summing all the  $\chi^2(n_i - 1)$  variables yield a

$$\chi^2 \left( \underbrace{\sum_{i=1}^k (n_i - 1)}_{\text{add all df}} \right) = \chi^2(n - k)$$

This is a ton of work, but in the end we have

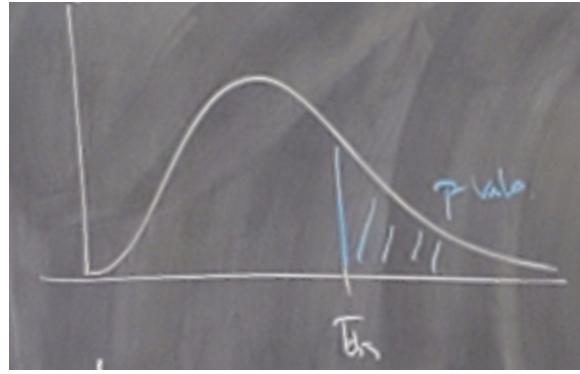
$$F = \frac{\sum_{i=1}^k n_i(\widehat{\mu}_i - \mu_0)^2 / k - 1}{\sum_{i=1}^k \sum_{j=1}^{n_i} (Y_{ij} - \widehat{\mu}_i)^2 / n - k} \sim F(k - 1, n - k)$$

when  $H_0$  is true.

Topic: The ANOVA test.

- ① Compute the null and alternative estimate of all the sample means.
- ② Use these to compute the sum of squares in our  $F$ -statistic
- ③ Construct a table (R or Excel?) to compute the  $p$ -value.  
 $F_{\text{obs}}$  is  $F$  observed. We computed this.

$$\Pr(F > F_{\text{obs}}), F - F(k - 1, n - k)$$



④ Conclusion. There is standard language for all this. In the denominator,

$$SSW = \sum_{i=1}^k \sum_{j=1}^{n_i} (Y_{ij} - \hat{\mu}_i)^2$$

the sum of squares WITHIN a group. And the numerator

$$SSA = \sum_{i=1}^k n_i (\hat{\mu}_i - \mu_0)^2$$

the sum of the squares AMONG groups. A measure of variance among the sample means. Recall originally,

$$\frac{S_0^2}{S_a^2} = \frac{S_a^2 + SSA}{SSW}$$

Thus we showed the factorization yield total variation.

$$S_0^2 = SSW + SSA \quad \text{or} \quad SST = SSW + SSA$$

### ANOVA Test

$$SSA = \sum_{i=1}^k n_i (\hat{\mu}_i - \hat{\mu}_0)^2 \quad (\text{SS among groups})$$

$$SSW = \sum_{i=1}^k \sum_{j=1}^{n_i} (Y_{ij} - \hat{\mu}_i)^2 \quad \text{SS within groups}$$

$$S_0^2 = SST = SSA + SSW$$

$$R^2 = \frac{SSA}{SST} \quad \text{coef. of determination}$$

Discussion: The ANOVA table

Source	SS	df	Mean Squares	$F_{\text{obs}}$	p-value
among	SSA	$k - 1$	$SSA / (k - 1)$	$MSA / MSW$	$\Pr(F(k - 1, n - k) > F_{\text{obs}})$
within	SSW	$n - k$	$SSW / (n - k)$		
total	SST	$n - 1$			

**Example:**

Treatment	Data		
1	2	4	3
2	6	4	
3	3	5	4

Mean squares  $k = 3, n_1 = 3, n_2 = 2, n_3 = 3, n = 8, \hat{\mu}_1 = 3, \hat{\mu}_2 = 5, \hat{\mu}_3 = 4$ .

$$\hat{\mu}_0 = \frac{\sum \sum Y_{ij}}{n} = \frac{n_1 \hat{\mu}_1 + n_2 \hat{\mu}_2 + n_3 \hat{\mu}_3}{n} = \frac{9 + 10 + 12}{8} = 3.875.$$

$$\begin{aligned} \text{SSA} &= \sum_{i=1}^k n_i (\hat{\mu}_i - \hat{\mu}_0)^2 \\ &= 3(3 - 3.875)^2 + 2(5 - 3.875)^2 + 3(4 - 3.875)^2 \\ &= 4.875. \end{aligned}$$

$$\begin{aligned} \text{SSW} &= \sum_{i=1}^k \sum_{j=1}^{n_i} (Y_{ij} - \hat{\mu}_i)^2 \\ &= \underbrace{(2-3)^2 + (4-3)^2 + (3-3)^2}_{i=1} + \underbrace{(6-5)^2 + (4-5)^2}_{i=2} + \underbrace{(3-4)^2 + (5-4)^2 + (4-4)^2}_{i=3} \\ &= 6 \end{aligned}$$

$$\text{SST} = \text{SSW} + \text{SSA} = 10.875.$$

Note that  $R^2 = \frac{\text{SSA}}{\text{SST}} \approx 0.326$ . Make table:

Source	SS	df	Mean Squares	$F_{\text{obs}}$	p-value
among	4.875	2	2.4375	2.03125	0.2261023
within	6	5	1.2		
total	10.875	7			

$$\begin{aligned} p &= \Pr(F(2, 5) > F_{\text{obs}} = 2.03125) \\ &= \text{pf}(F_{\text{obs}}, \text{k}-1, \text{n}-\text{k}, \text{lower.tail}=FALSE) \quad (\text{R code}) \\ &= 0.2261023. \end{aligned}$$

Hence a large  $p$ -value does not support the alternative. No evidence that the means are different.

**Example:** Recall the typists/keyboards designs:

design	times
KB1	364 366 394 386 379 398 371 370
KB2	355 359 374 342 378 355 376 358
KB3	360 345 374 390 386 373 393 366

Then  $\hat{\mu}_1 = 378.5$ ,  $\hat{\mu}_2 = 362.125$ ,  $\hat{\mu}_3 = 373.375$ .  $k = 3$ ,  $n_1 = n_2 = n_3 = 8$ ,  $n = 24$ .

R-code, with little c for column.  $y = c(364, 355, 360, 366, 359, 345, \dots, 370, 358, 366)$ .  
 $\text{typists} = \text{rep}(1:3, 8)$  anova (lm (y factor(typists))). R spits out:

	df	SS	Mean Squares	$F_{\text{obs}}$	p-value
factor(typists)	2	532.0	266.00	1.5181	0.2422
residuals	21	3679.6	175.22		

Conclusion, again, fairly large  $p$  value. Cannot reject  $H_0$ . Practical conclusion: No evidence that our keyboard design is different in use than another. Remark: Missing from R's table is  $\text{SST} = \text{SSW} + \text{SSA} = 4211.6$  and  $R^2 = \text{SSA} / \text{SST} = 0.126$ .

**Example:** Background sounds and impact on memory. 30 students are randomly divided into 3 sets of 10. Silence, classical, and jazz. They are told to “study”. The “data ish”

	Quiet	Classical	Jazz
$\hat{\mu}_i$ sample means	89.5	89.7	79.4
$S_i$ sample standard dev.	8.91	10.25	7.06

Need  $k = 3$  treatments.  $n_1 = n_2 = n_3 = 10$  and  $n = 30$ . Then

$$\hat{\mu}_0 = \frac{n_1\hat{\mu}_1 + n_2\hat{\mu}_2 + \dots + n_k\hat{\mu}_k}{n} = 86.2$$

$$\begin{aligned} \text{SSA} &= \sum_{i=1}^k n_i(\hat{\mu}_i - \hat{\mu}_0)^2 \\ &= 10(89.5 - 86.2)^2 + 10(89.7 - 86.2)^2 + 10(79.4 - 86.2)^2 \\ &= 680.3 \end{aligned}$$

Note we have  $S_i^2 = \frac{\sum_j (Y_{ij} - \hat{\mu}_i)^2}{n_i - 1}$ .

$$\begin{aligned} \text{SSW} &= \sum_{i=1}^k \sum_{j=1}^{n_i} (Y_{ij} - \hat{\mu}_i)^2 \\ &= \sum_{j=1}^{10} (Y_{1j} - \hat{\mu}_1)^2 + \sum_{j=1}^{10} (Y_{2j} - \hat{\mu}_2)^2 + \sum_{j=1}^{10} (Y_{3j} - \hat{\mu}_3)^2 \\ &= 9 \left( \frac{1}{9} \sum_{j=1}^{10} (Y_{1j} - \hat{\mu}_1)^2 + \frac{1}{9} \sum_{j=1}^{10} (Y_{2j} - \hat{\mu}_2)^2 + \frac{1}{9} \sum_{j=1}^{10} (Y_{3j} - \hat{\mu}_3)^2 \right) \end{aligned}$$

So  $\text{SSW} = 9(S_1^2 + S_2^2 + S_3^2)$ . Fact,  $\text{SSW} = \sum_{i=1}^k (n_i - 1)S_i^2$ . Here,  $\text{SSW} = 2108.65$ .

Source	SS	df	Mean Squares	$F_{\text{obs}}$	p-value
among	680.3	2	340.15	4.3598	0.022863
within	2108.65	27	78.09815		
total	2788.95	29			

**Example:** Background sounds (again)

	Quiet	Classical	Jazz
$\hat{\mu}_i$ sample means	89.5	89.7	79.4
$S_i$ sample standard dev.	8.91	10.25	7.06

$n_1 = n_2 = n_3 = 10$  and  $n = 30$  and  $k = 3$ . Then

$$\hat{\mu}_0 = \frac{\sum n_i \hat{\mu}_i}{n} = 86.2$$

$$SSA = \sum_{i=1}^3 n_i (\hat{\mu}_i - \hat{\mu}_0)^2 = 680.3$$

$$SSW = \sum_{i=1}^3 \sum_{j=1}^{10} (Y_{i,j} - \hat{\mu}_i)^2 = \sum_{i=1}^3 0.95_i^2 = 2108.65$$

FACT:  $SSW = \sum_{i=1}^k (n_i - 1) S_i^2$ . Recall:

Source	SS	df	Mean Squares	$F_{\text{obs}}$	p-value
among	680.3	2	340.15	4.3598	0.022863
within	2108.65	27	78.09815		
total	2788.95				

Where  $p$ -value is  $\Pr(F(2, 27) > 4.3598) \approx 2.28\%$ .

At  $\alpha = 0.05$  level, we cannot reject  $H_0$ . Conclusion: “At least one of the group means is different than the others.”

Topic: Confidence interval (§ 13.7)

We may be interested in a CI for the group mean. We clearly have an estimate for  $\mu_i$  and  $\hat{\mu}_i$ .

$$\text{C.I.} \equiv \hat{\theta} \pm t^* \sqrt{\text{Var}(\hat{\theta})}.$$

Big question: What do we use for variance?

We have estimators  $S_0^2, S_a^2$  which are MLE.

Issue: As with the original  $S^2 = \frac{1}{n} \sum (X_i - \bar{X})^2$ , our MLEs are biased.

$H$  is true if  $H_0$  is true. Both alone are “good” estimators. If  $H_a$  is true, it stands to reason that  $S_a^2$  is “better”. Using a bit of algebra, we can make an unbiased estimator of  $S_a^2$ .

Last day: we showed

$$SSW = \sum_{i=1}^k \sum_{j=1}^{n_i} (Y_{i,j} - \hat{\mu}_i)^2 = \sum_{i=1}^k (n_i - 1) S_i^2.$$

This is a weighted sum of unbiased estimators  $S_i^2$  for  $\sigma^2$ . As individually, each  $S_i^2 \sim \sigma^2$  (unbiased). Then the collective average should be an even better approximation.

$$\sigma^2 = \frac{\sum_{i=1}^k (n_i - 1)S_i^2}{(n_1 - 1) + \dots + (n_k - 1)} = \frac{\sum_{i=1}^k (n_i - 1)S_i^2}{n - k} = \text{MSW}$$

**FACT:** We can use the MSW as an unbiased estimator for  $\sigma^2$ .

Aside: Isn't this just the pooled estimator written largely?

$$S_p^2 = \frac{(n_1 - 1)S_1^2 + (n_2 - 1)S_2^2}{n_1 + n_2 + k}$$

Here

$$S^2 = \frac{(n_1 - 1)S_1^2 + \dots + (n_k - 1)S_k^2}{n_1 + n_2 + \dots + n_k - k}$$

We are constructing hypothesis of different of CI. The “t” inherits the df from  $S^2$ . For our CI we need  $\text{df} = n - k$ .

**Example:** Background sounds (again)

- a.) Find a 95% CI for the mean score for someone listening to Jazz.

Here,  $n_3 = 10$  and  $t_{\alpha/2}(\text{df}) = t_{0.025}(27) \approx 2.052$ .

Then,  $\hat{\mu}_3 \pm t_{0.025}(27)\sqrt{S^2/n_3}$  with  $S^2/n$  being the total variance affiliated with sampling dist.

We have  $79.4 \pm 2.052\sqrt{\frac{78.09815}{10}}$  which is

$$79.4 \pm 5.734 \quad \text{or} \quad (73.665, 85.134)$$

Note  $\hat{\mu}_1 = 89.5$  and  $\hat{\mu}_2 = 89.7$  this is still a couple standard errors away.

- b.) Is the proof of the other 2  $\mu_i$ 's are different? Note: We could ANOVA again without Jazz, or do chapter 9 stuff:

Doing Ch. 9 stuff, consider  $\hat{\mu}_1 - \hat{\mu}_2 = -0.2$ . Then to make a CI we use the exact same formula from before (Ch9) but we use the “ANOVA” estimator for  $S^2 = \text{MSE}$  with  $\text{df} = n - k$ .

$$\begin{aligned} \hat{\mu}_1 - \hat{\mu}_2 &\pm t_{0.025}(27)\sqrt{\frac{S^2}{n_1} + \frac{S^2}{n_2}} = \hat{\mu}_1 - \hat{\mu}_2 \pm t_{0.025}(27)S\sqrt{\frac{1}{n_1} + \frac{1}{n_2}} \\ &= -0.2 \pm 20.052\sqrt{\frac{75.09815}{5}} \\ &= -0.2 \pm 7.95 \\ &\equiv (-8.15, 7.75) \end{aligned}$$

And notice that  $0 \in \text{CI}$ .

Summary: 2-sided  $1 - \alpha$  level CI.

$$\widehat{\mu}_i \pm t_{\alpha/2}(\text{df}) \frac{S}{\sqrt{n_i}} = \widehat{\mu}_i - \widehat{\mu}_j \pm t_{\alpha/2}(\text{df}) S \sqrt{\frac{1}{n_i} + \frac{1}{n_j}}$$

Where  $S^2 = \text{MSW} = \frac{\text{SSW}}{n - k}$  with  $\text{df} = n - k$ .

**Example:** Arbitrary numbers (no story)

A	80	85	71	64
B	70	72	75	
C	83	70		

Then  $n_1 = 4$ ,  $n_2 = 3$ ,  $n_3 = 2$ ,  $n = 9$ ,  $\bar{A} = 75$ ,  $\bar{B} = 217/3$ ,  $\bar{C} = 153/2$  and

$$\widehat{\mu}_0 = \frac{4\bar{A} + 3\bar{B} + 2\bar{C}}{9} = \frac{690}{9}.$$

	SS	df	MS	$F_{\text{obs}}$
among	SSA = 23.0556	2	11.5278	0.192575
within	SSW = 359.167	6	39.8611	

The above was computed on mathematica. Using R for the  $p$ -value,  $p \approx 82.97\%$ . We cannot reject  $H_0$ .

The R code:

```
scores = c(80, 85, 71, 64, 70, 72, 75, 83, 70)
treat = c(rep("A", 4), rep("B", 3), rep("C", 2))
table = data.frame(saved, treat)
results = aov(scores ~ treat, data = table)
summary(results)
```

The output is

	df	Sum Square	Mean Square	$F$	$\text{Pr}(> F)$
treat	2	23.1	11.53	0.193	0.83
residual	6	359.2	59.86		