

## 8 Estimation

### 8.1 Introduction

Recall that a statistic is a random variable that is a function of the random sample. We call a statistic an *estimator* if it is supposed to estimate or approximate some feature of the population. A *point estimator* is an estimator whose value is a single number and which estimates a single parameter for the population. For example, the sample mean is a point estimator for the population mean. The sample variance is a point estimator for the population variance.

A point estimator only gives an approximation to the parameter it is supposed to estimate. An important question is how good it is. In the following section we consider two characterizations of goodness - bias and mean square error.

### 8.2 Bias and mean square error for point estimators

Suppose  $\hat{\theta}$  is an estimator for the population parameter  $\theta$ . So  $\hat{\theta}$  is a RV and so it has a mean. This mean may or may not be equal to  $\theta$ .

**Definition 1.**  $\hat{\theta}$  is an unbiased estimator of  $\theta$  if  $E[\hat{\theta}] = \theta$ . If they are not equal we say the estimator is biased. The bias is defined to be

$$B(\hat{\theta}) = E[\hat{\theta}] - \theta \quad (1)$$

The bias is only part of the story. Even if the estimator is unbiased, if the estimator has a large variance it will not usually be close to  $\theta$ . So we also need to consider the variance of the estimator. But if the bias is large and the variance is small, the estimator will not usually be close to  $\theta$ . So rather than just looking at the bias or just looking at the variance, we consider a better measure of goodness - how close  $\hat{\theta}$  is to  $\theta$  on average.

**Definition 2.** The mean square error (MSE) of a point estimator for  $\theta$  is

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

**Proposition 1.**

$$MSE(\hat{\theta}) = Var(\hat{\theta}) + [B(\hat{\theta})]^2 \quad (3)$$

**Proof:** To make the following computation a little easier to follow, let  $c = E[\hat{\theta}]$ . Note that  $c$  is a constant, i.e., non-random.

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

**Definition 3.** *The standard error of a point estimator is the standard deviation of the estimator.*

### 8.3 Some unbiased point estimators

An important problem in statistics is how do you find an estimator for a population parameter. We will study this in chapter 9. Here we look at some intuitively obvious estimators.

**The mean of one population:** Suppose we have a single population with mean  $\mu$ . We draw a random sample  $Y_1, Y_2, \dots, Y_n$ . The sample mean,

$$\bar{Y} = \frac{1}{n} \sum_{i=1}^n Y_i \quad (4)$$

is a natural estimator for  $\mu$ . Since  $E[\bar{Y}] = \mu$ , it is an unbiased estimator. The variance of this estimator is  $\sigma^2/n$  where  $\sigma^2$  is the population variance. (The book says the variance of the sampling distribution of the estimator is  $\sigma^2/n$ .) Note that the MSE is also  $\sigma^2/n$  since the bias is zero.

**The difference of the means of two populations:** Suppose we have two populations. Population 1 has mean  $\mu_1$  and variance  $\sigma_1^2$ . Population 2 has mean  $\mu_2$  and variance  $\sigma_2^2$ . We want to estimate  $\mu_1 - \mu_2$ . We draw a random sample of size  $n_1$  from population 1 and let  $\bar{Y}_1$  denote the sample mean for this sample. We draw an independent random sample of size  $n_2$  from population 2 and let  $\bar{Y}_2$  denote the sample mean for this sample. (This is slightly dangerous notation since we have also used a subscript on  $\bar{Y}$  to

indicate the sample size. Here the subscript just indicates which population.) Then  $\bar{Y}_1 - \bar{Y}_2$  is a point estimator for  $\mu_1 - \mu_2$ . You studied this estimator in homework 2. The mean is  $\mu_1 - \mu_2$ , so it is an unbiased estimator. The variance is

$$Var(\bar{Y}_1 - \bar{Y}_2) = Var(\bar{Y}_1) + (-1)^2 Var(\bar{Y}_2) = \frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2} \quad (5)$$

Since the bias is zero, the MSE is equal to this variance.

---

### End of lecture on Thurs, 1/25

---

**Proportion for one population:** Suppose that the experiment only has two outcomes which we will call success and failure. We let  $p$  be the probability of success. In terms of the population, we are considering a population in which each member of the population has one of two properties which we call success or failure. An example: the population is all registered voters in the US and we look at whether they voted (success) or not (failure) in the most recent election. In this case  $p$  is the fraction of registered voters who voted in the most recent election. Another example: laboratory mice are given a vaccine for a virus and then exposed to the virus. Success is that the mouse gets the virus, failure that the mouse does not. So  $p$  is the probability that a vaccinated mouse gets the virus.

The parameter  $p$  is the proportion of the population that has success. So we refer to  $p$  as the population proportion. We let  $Y$  be 1 when we have success and 0 for failure. Note that the mean of  $Y$  is just  $p$ . Now we take a random sample  $Y_1, \dots, Y_n$  and compute the sample mean in the usual way. Note that  $\sum_{i=1}^n Y_i$  is the number of successes in our sample. So  $\bar{Y}$  is the proportion of successes in the sample. It is a point estimator for  $\mu = p$ , the population proportion. So we denote it by  $\hat{p}$ . The expected value of  $\hat{p}$  is  $\mu = p$ . So the sample proportion  $\hat{p}$  is an unbiased estimator of  $p$ . The variance of  $Y$  is just  $p(1 - p)$  and so the variance of  $\hat{p}$  is  $p(1 - p)/n$ .

**Difference of the proportions for two populations:** Now suppose we have two populations and the experiments for the populations only have two outcomes. We let  $p_1$  be the proportion for population 1 and  $p_2$  the proportion for population 2. We want to estimate  $p_1 - p_2$ . We draw a random sample from each population and assume that the two samples are independent of each other. We let  $\hat{p}_1$  and  $\hat{p}_2$  be the sample proportions for the two random

samples. A natural estimator for  $p_1 - p_2$  is  $\hat{p}_1 - \hat{p}_2$ . The mean is  $p_1 - p_2$ . So it is an unbiased estimator. Since  $\hat{p}_1$  and  $\hat{p}_2$  are independent RV's, the variance of the estimator is

$$Var(\hat{p}_1 - \hat{p}_2) = Var(\hat{p}_1) + Var(\hat{p}_2) = \frac{p_1(1-p_1)}{n_1} + \frac{p_2(1-p_2)}{n_2} \quad (6)$$

An example of the above is the following. We divide our lab mice into a two groups (or populations). Group 1 gets the vaccine. Group 2 (the control group) does not. Both groups are exposed to the virus. For both groups success is getting the virus. So  $p_1$  is the probability a vaccinated mouse gets the virus,  $p_2$  the probability an unvaccinated mouse gets it. The parameter difference  $p_1 - p_2$  is a measure of the effectiveness of our vaccine.

Finally we consider the variance for a single population. Recall that the sample variance is defined by

$$S^2 = \frac{1}{n-1} \sum_{i=1}^n (Y_i - \bar{Y})^2 \quad (7)$$

The following proposition says that  $S^2$  is an unbiased estimator for  $\sigma^2$ .

**Proposition 2.** *Let  $S^2$  be the usual sample variance:*

$$S^2 = \frac{1}{n-1} \sum_{i=1}^n (Y_i - \bar{Y})^2 \quad (8)$$

*Then*

$$E[S^2] = \sigma^2 \quad (9)$$

*So  $S^2$  is an unbiased estimator for the population variance. The variance of  $S^2$  is*

$$var(S^2) = \frac{\mu_4}{n} - \frac{\sigma^4(n-3)}{n(n-1)} \quad (10)$$

*where  $\mu_4 = E[(Y - \mu)^4]$ .*

**Proof:**

$$\sum_{i=1}^n (Y_i - \bar{Y})^2 = \sum_{i=1}^n Y_i^2 - n(\bar{Y})^2 \quad (11)$$

$$(12)$$

We use the equation

$$\text{Var}(Y) = E[Y^2] - (E[Y])^2 \quad (13)$$

So  $E[Y_i^2] = \text{Var}(Y_i) + (E[Y_i])^2 = \sigma^2 + \mu^2$ . And

$$E[(\bar{Y})^2] = \text{Var}(\bar{Y}) + (E[\bar{Y}])^2 = \frac{\sigma^2}{n} + \mu^2 \quad (14)$$

So

$$\begin{aligned} E\left[\sum_{i=1}^n (Y_i - \bar{Y})^2\right] &= \sum_{i=1}^n E[Y_i^2] - nE[(\bar{Y})^2] \\ &= n(\sigma^2 + \mu^2) - n\left(\frac{\sigma^2}{n} + \mu^2\right) \\ &= (n-1)\sigma^2 \end{aligned}$$

Thus  $E[S^2] = \sigma^2$ .

The formula for the variance of  $S^2$  come from a rather long computation which we do not include.  $\square$

The estimator in homework problem 2 for the variance was

$$V = \frac{1}{n} \sum_{i=1}^n Y_i^2 - (\bar{Y})^2 \quad (15)$$

You should have found in the homework that

$$E[V] = \frac{n-1}{n} \sigma^2 \quad (16)$$

So  $V$  is a biased estimator of  $\sigma^2$ . Note that the bias goes to zero as  $n \rightarrow \infty$ .

**Remark:**  $S^2$  is an unbiased estimator for  $\sigma^2$ . Define  $S$  to be  $\sqrt{S^2}$ . This is a natural estimator for  $\sigma$ . However,  $E[S]$  is not equal to  $(E[S^2])^{1/2} = \sigma$ . So  $S$  is a biased estimator for  $\sigma$ .

## 8.4 The goodness of a point estimator

Throughout this section  $\hat{\theta}$  is an estimator for the population parameter  $\theta$ .

**Definition 4.** *The error of estimation  $\epsilon$  is defined to be*

$$\epsilon = |\hat{\theta} - \theta| \quad (17)$$

Note that the error of estimation is a random variable. Suppose that our estimator is unbiased. Recall that the standard error is just another name for the standard deviation of the estimator. The main point of this section is to see what we can say about the probability that the error of estimation is less than two times the standard error. We start with an upper bound that always holds.

**Theorem 1. (Chebyshev's theorem)** *For any RV  $Y$  with finite variance  $\sigma^2$  and mean  $\mu$ ,*

$$P(|Y - \mu| \geq \delta) \leq \frac{1}{\delta^2} \sigma^2 \quad (18)$$

We will use this with  $\delta = 2\sigma$ . Then it says

$$P(|Y - \mu| \geq 2\sigma) \leq \frac{1}{4} \quad (19)$$

This is the same as

$$P(|Y - \mu| < 2\sigma) \geq \frac{3}{4} \quad (20)$$

So for any random variable the probability the RV is within two standard deviations of the mean is at least 0.75. For most distributions, and especially those that occur as the distribution of an estimator, this probability is much closer to 1. A good rule of thumb is that this probability is around 0.95. For a normal distribution this probability is  $P(|Z| \geq 2) \approx 0.0544$  (where  $Z$  is standard normal). For a uniform distribution this probability is 1.

An example of a distribution for which the probability is essentially 0.75 is the following. Let  $Y$  be a discrete RV which only takes on the three values  $-1, 0, 1$ . We take

$$P(Y = 0) = 1 - p, P(Y = 1) = P(Y = -1) = p/2 \quad (21)$$

The mean is 0. The second moment is  $E[Y^2] = p$ . So  $\sigma = \sqrt{p}$ . If  $2\sqrt{p} < 1$ , then the probability that  $Y$  is more than  $2\sigma$  from the mean is just the probability that  $Y = \pm 1$  which is  $p$ . The condition  $2\sqrt{p} < 1$  is the same as  $p < 1/4$ . So if we take  $p$  just slightly smaller than  $1/4$ , then the probability that  $Y$  is more than  $2\sigma$  from the mean is almost  $1/4$ . This example shows that although the Chebyshev bound is usually far from the truth, there are distributions for which you cannot do any better.

Returning to estimators, we apply the above bound with  $Y = \hat{\theta}$  and  $\delta = 2\sigma_{\hat{\theta}}$ . Note that  $\sigma_{\hat{\theta}}$  is the standard error. Since the estimator is unbiased, the mean of  $\hat{\theta}$  is  $\theta$ . So we get

$$P(|\hat{\theta} - \theta| < 2\sigma_{\hat{\theta}}) \geq \frac{3}{4} \quad (22)$$

Thus for any unbiased estimator, with probability at least 0.75, the error of estimation is no more than twice the standard error. And for most estimators the probability is closer to 1, often around 0.95.

**Example:** Consider the mice and vaccine example from previous section. 100 mice are given vaccine and exposed to a virus. 15 of them get the virus. 50 control mice are not given the vaccine and exposed to the virus. 42 of them get the virus. Estimate  $p_2 - p_1$  and give a two standard error bound on the estimate.

We have  $\hat{p}_1 = 15/100 = 0.15$  and  $\hat{p}_2 = 42/50 = 0.84$ . So our estimate for  $p_2 - p_1$  is  $0.84 - 0.15 = 0.69$ . The standard error is given by  $\sqrt{p_1(1-p_1)/n_1 + p_2(1-p_2)/n_2}$ . We don't know  $p_1$  and  $p_2$ , so we have to approximate them with  $\hat{p}_1$  and  $\hat{p}_2$ . So the standard error is approximately

$$\sqrt{0.15(1-0.15)/100 + 0.84(1-0.84)/50} = 0.0629 \quad (23)$$

So 2 times the standard error is 0.126. So the true value of  $p_2 - p_1$  will be within 0.126 of 0.69 with probability of at least 0.75. The sample sizes are pretty large, so by the CLT, our estimator  $\hat{p}_2 - \hat{p}_1$  will be approximately normal. So the probability our estimate is within 0.126 of 0.69 should actually be about 0.95.

---

**End of lecture on Tues, 1/30**

---

So far we have only considered unbiased estimators for which the error of estimation is  $|\hat{\theta} - \theta| = |\hat{\theta} - E[\hat{\theta}]|$ . What if the estimator is biased? Can we get a similar statement? We start with another version of Chebyshev.

**Theorem 2. (Chebyshev's theorem - another version)** *For any RV  $Y$  with finite variance  $\sigma^2$  and any constant  $c$*

$$P(|Y - c| \geq \delta) \leq \frac{1}{\delta^2} E[(Y - c)^2] \quad (24)$$

Taking  $c = \mu$  gives the version we saw before. Now we apply this with  $Y = \hat{\theta}$ ,  $c = \theta$ , and  $\delta = 2\sqrt{MSE}$ . Recall that the MSE is  $E[(\hat{\theta} - \theta)^2]$ . So the above becomes

$$P(|\hat{\theta} - \theta| \geq 2\sqrt{MSE}) \leq \frac{1}{4} \quad (25)$$

Note that if the estimator is unbiased, then  $\sqrt{MSE}$  is the standard error and this becomes the same bound we saw before. As before this probability is typically much closer to 0 than 1/4. In many case it is about 0.05. So for a biased estimator, the error in estimation is no more than  $2\sqrt{MSE}$  with probability at least 75% and more typically 95%.

## 8.5 Confidence intervals

A confidence interval for  $\theta$  is an estimator for  $\theta$  which is not a point estimator. Instead the estimator gives an interval. We denote this interval by  $[\hat{\theta}_L, \hat{\theta}_U]$ .

**Definition 5.** *The confidence coefficient is*

$$P(\theta \in [\hat{\theta}_L, \hat{\theta}_U]) \quad (26)$$

The confidence coefficient is a measure of how reliable our confidence intervals is. The closer the confidence coefficient is to 1, the more reliable the confidence interval. Note that if we make the confidence interval larger, then the confidence coefficient will increase. But to say  $\mu$  belongs to a large confidence interval is less desirable than being able to say it belongs to a small confidence interval. So there is a trade-off. Typical choices for the confidence coefficient are 0.9, 0.95, 0.98 or maybe even 0.99. Confidence coefficients are often written in the form  $1 - \alpha$ . So  $\alpha$  is typically 0.1, 0.05, 0.02, 0.01.

We can also consider one-sided confidence intervals, i.e.,  $(-\infty, \hat{\theta}_U]$  or  $[\hat{\theta}_L, \infty)$ . For these the confidence coefficient would be  $P(\theta \leq \hat{\theta}_U)$  or  $P(\theta \geq \hat{\theta}_L)$ , respectively.

Given a particular confidence coefficient  $1 - \alpha$  we want to find the confidence interval, i.e.,  $\hat{\theta}_L$  and  $\hat{\theta}_U$ . Note that  $\hat{\theta}_L$  and  $\hat{\theta}_U$  must be functions of the



random sample and constants like the sample size, but they cannot depend on unknown population parameters like  $\theta$ . So we want to find estimators  $\hat{\theta}_L$  and  $\hat{\theta}_U$  so that

$$P(\theta \in [\hat{\theta}_L, \hat{\theta}_U]) = 1 - \alpha \quad (27)$$

This is the same as

$$P(\theta < \hat{\theta}_L) + P(\theta > \hat{\theta}_U) = \alpha \quad (28)$$

There is freedom here as to how much of  $\alpha$  goes into each of the two probabilities. The simplest thing to do is require

$$P(\theta < \hat{\theta}_L) = \alpha/2, \quad P(\theta > \hat{\theta}_U) = \alpha/2 \quad (29)$$

The difficulty here is that the distribution of the estimators  $\hat{\theta}_L$  and  $\hat{\theta}_U$  will typically depend on  $\theta$  which we do not know. To deal with this difficulty we introduce the idea of a *pivotal quantity*.

**Definition 6.** *A pivotal quantity (for the parameter  $\theta$ ) is a function of the random sample with two properties*

- *The pivotal quantity depends on the random sample and the unknown parameter  $\theta$ , but it does not depend on any other unknown parameters.*
- *The distribution of the pivotal quantity does not depend on  $\theta$ .*

Given a pivotal quantity  $U$ , the strategy for using it to get a confidence interval is the following. Since the distribution of  $U$  does not depend on  $\theta$ , we can find numbers  $a$  and  $b$  so that  $P(a \leq U \leq b) = 1 - \alpha$ . Since  $U$  depends on the random sample and  $\theta$ , the inequality  $a \leq U \leq b$  can be transformed into a statement that  $\theta$  belongs to some interval whose endpoints are functions of the random sample. We illustrate with two examples.

**Example:** We first do an unrealistic example to illustrate the idea. Suppose that the population has an exponential distribution. So  $Y$  has an exponential distribution. To be consistent with the notation in this section we let  $\theta$  be the parameter for the exponential. So the mean of  $Y$  is  $\theta$ . We do not know  $\theta$  and want a confidence interval for it. We take a random sample of size 1. So the random sample is just  $Y_1$ . We claim that  $Y_1/\theta$  is a pivotal quantity. It clearly meets the first condition. We need to show that its distribution does not depend on  $\theta$ . **SHOW THIS.**

Now we use this pivotal quantity to find a confidence interval for  $\theta$ . We take  $\alpha = 0.1$ . We start by finding an interval  $[a, b]$  so that the pivotal quantity is in  $[a, b]$  with probability 0.9. Let  $U = Y_1/\theta$  denote the pivotal quantity. So we want

$$P(a \leq U \leq b) = 0.9 \quad (30)$$

There are many  $a, b$  that satisfy this. We will ask that

$$P(U < a) = 0.05, \quad P(U > b) = 0.05 \quad (31)$$

The distribution of  $U$  is just an exponential with mean 1. So  $a = \text{qexp}(0.05, 1) = 0.05129$ . And  $b = \text{qexp}(0.95, 1) = 2.9957$ . Now we use the fact that the pivotal quantity depends on  $\theta$  but no other unknown parameter to turn the statement that  $a \leq U \leq b$  into a statement about  $\theta$ :

$$\{a \leq U \leq b\} = \{a \leq Y_1/\theta \leq b\} = \{1/b \leq \theta/Y_1 \leq 1/a\} \quad (32)$$

$$= \{Y_1/b \leq \theta \leq Y_1/a\} = \{\theta \in [Y_1/b, Y_1/a]\} \quad (33)$$

So

$$P(\theta \in [Y_1/b, Y_1/a]) = 0.9 \quad (34)$$

We have found a confidence interval for  $\theta$  with confidence coefficient 0.9. It is  $[Y_1/b, Y_1/a] = [0.3338Y_1, 19.50Y_1]$ . Note that this confidence interval is huge. But that is to be expected if we use a sample of size 1.

**Example(continued):** We continue to consider a  $Y$  that has an exponential distribution with mean  $\theta$ . But now we take a sample of size  $n$ . In a homework problem we worked out the pdf of  $\bar{Y}_n$  in this case. It has a gamma distribution with  $\alpha_{\bar{Y}} = n$  and  $\beta_{\bar{Y}} = \theta/n$ . We claim that  $U = \bar{Y}_n/\theta$  is a pivotal quantity. It clearly meets the first condition. We need to show that its distribution does not depend on  $\theta$ . To show this, recall that the mgf of a gamma distribution with parameters  $\alpha, \beta$  is

$$M(t) = \left( \frac{1}{1 - \beta t} \right)^\alpha \quad (35)$$

SHOW the distribution of  $U$  is a gamma distribution with  $\alpha = n$  and  $\beta = 1/n$ . We now look for  $a$  and  $b$  such that

$$P(U < a) = 0.05, \quad P(U > b) = 0.05 \quad (36)$$

The values of  $a$  and  $b$  will depend on  $n$ . To be concrete we work this out for sample sizes of  $n = 10, 100$ . Note that in R, the arguments for the gamma routine are a bit confusing. The default arguments are not  $\alpha, \beta$  but rather  $\alpha, 1/\beta$ . For example, to compute the value of the cdf at 0.05 for  $\alpha = 10$ ,  $\beta = 0.1$  we use `qgamma(0.05, 10, 10)`. You can force it to use  $\beta$  instead of  $1/\beta$  by `qgamma(0.05, 10, scale = 0.1)`. (scale is another name for the parameter  $\beta$ .)

$$\begin{aligned} a_{10} &= \text{qgamma}(0.05, 10, 10) = 0.5425 \\ b_{10} &= \text{qgamma}(0.95, 10, 10) = 1.5705 \\ a_{100} &= \text{qgamma}(0.05, 100, 100) = 0.8414 \\ b_{100} &= \text{qgamma}(0.95, 100, 100) = 1.170 \end{aligned}$$

So the 90% confidence intervals are  $[0.6367\overline{Y}_{10}, 1.843\overline{Y}_{10}]$  for  $n = 10$  and  $[0.8547\overline{Y}_{100}, 1.1885\overline{Y}_{100}]$  for  $n = 100$ . Note that as the sample size gets larger, the confidence interval gets smaller.

## 8.6 Large sample confidence intervals

In this section we consider confidence intervals corresponding to the point estimators we studied earlier :  $\hat{\mu}, \hat{\mu}_1 - \hat{\mu}_2, \hat{p}, \hat{p}_1 - \hat{p}_2$  under the assumption that the sample size(s) are large. If the sample size is large, then the distribution of all of these estimators is approximately normal. Recall that all four of these estimators are unbiased. Let  $\hat{\theta}$  be one of these four estimators and  $\theta$  the corresponding population parameter. So  $E[\hat{\theta}] = \theta$ . If we define

$$Z = \frac{\hat{\theta} - \theta}{\sigma_{\hat{\theta}}} \tag{37}$$

then the distribution of  $Z$  will be approximately the standard normal. This is one of the criteria for a pivotal quantity.  $Z$  does not quite meet the first criteria since it contains the unknown parameter  $\sigma_{\hat{\theta}}$  in addition to the unknown parameter  $\theta$ . But for large samples we can approximate  $\sigma_{\hat{\theta}}$  by a function of the random sample that does not involve any unknown parameters. We return to this point later.

Now suppose we want a confidence interval with confidence coefficient  $1 - \alpha$ . We start by looking for  $a$  and  $b$  so that

$$P(Z < a) = \alpha/2, \quad P(Z > b) = \alpha/2 \tag{38}$$

Since the distribution of  $Z$  is symmetric,  $a$  will just be  $-b$ . The usual notation is to write  $b$  as  $z_c$ . Here are some commonly used values of  $\alpha$  and the corresponding  $z_c$ :

$$\alpha = 0.1 \Rightarrow z_c = 1.645..., \quad (39)$$

$$\alpha = 0.05 \Rightarrow z_c = 1.960..., \quad (40)$$

$$\alpha = 0.01 \Rightarrow z_c = 2.576..., \quad (41)$$

$$(42)$$

So we now have  $P(-z_c \leq Z \leq z_c) = 1 - \alpha$ . Using the definition of  $Z$  this becomes

$$P(-z_c \leq \frac{\hat{\theta} - \theta}{\sigma_{\hat{\theta}}} \leq z_c) = 1 - \alpha \quad (43)$$

which we can rewrite as

$$P(\theta \in [\hat{\theta} - z_c \sigma_{\hat{\theta}}, \hat{\theta} + z_c \sigma_{\hat{\theta}}]) = 1 - \alpha \quad (44)$$

So the confidence interval is

$$[\hat{\theta} - z_c \sigma_{\hat{\theta}}, \hat{\theta} + z_c \sigma_{\hat{\theta}}] \quad (45)$$

This result applies to each of the four estimators  $\hat{\mu}, \hat{\mu}_1 - \hat{\mu}_2, \hat{p}, \hat{p}_1 - \hat{p}_2$ .

We now return to the issue of the dependence of  $\sigma_{\hat{\theta}}$  on unknown parameters. First suppose  $\theta = \mu$  and  $\hat{\theta} = \bar{Y}$ . For this case  $\sigma_{\hat{\theta}}$  is  $\sigma/\sqrt{n}$ , and  $\sigma$  is unknown. We approximate  $\sigma$  by the sample standard deviation  $S$ . So our confidence interval is

$$[\bar{Y} - z_c S/\sqrt{n}, \bar{Y} + z_c S/\sqrt{n}] \quad (46)$$

---

**End of lecture on Thurs, 2/1**

---

If  $\theta = \mu_1 - \mu_2$  and  $\hat{\theta} = \bar{Y}_1 - \bar{Y}_2$ , then we make the approximation

$$\sigma_{\hat{\theta}} = \sqrt{\sigma_1^2/n_1 + \sigma_2^2/n_2} \approx \sqrt{S_1^2/n_1 + S_2^2/n_2} \quad (47)$$

where  $S_1^2$  is the sample variance of the sample from population 1 and  $S_2^2$  is the sample variance of the sample from population 2. So our confidence interval is

$$[\bar{Y}_1 - \bar{Y}_2 - z_c \sqrt{S_1^2/n_1 + S_2^2/n_2}, \bar{Y}_1 - \bar{Y}_2 + z_c \sqrt{S_1^2/n_1 + S_2^2/n_2}] \quad (48)$$

If  $\theta = p$  and  $\hat{\theta} = \hat{p}$ , then we make the approximation

$$\sigma_{\hat{p}} = \sqrt{p(1-p)/n} \approx \sqrt{\hat{p}(1-\hat{p})/n} \quad (49)$$

So the confidence interval is

$$\hat{p} \pm z_c \sqrt{\hat{p}(1-\hat{p})/n} \quad (50)$$

Finally, for  $\theta = p_1 - p_2$  and  $\hat{\theta} = \hat{p}_1 - \hat{p}_2$ , we make the approximation

$$\sigma_{\hat{p}} = \sqrt{p_1(1-p_1)/n_1 + p_2(1-p_2)/n_2} \approx \sqrt{\hat{p}_1(1-\hat{p}_1)/n_1 + \hat{p}_2(1-\hat{p}_2)/n_2} \quad (51)$$

So the confidence interval is

$$\hat{p}_1 - \hat{p}_2 \pm z_c \sqrt{\hat{p}_1(1-\hat{p}_1)/n_1 + \hat{p}_2(1-\hat{p}_2)/n_2} \quad (52)$$

**Example:** We want to compare the average height of an adult female in Canada and the US. We randomly choose 1000 adult females in Canada (pop 1) and 1000 adult females in the US (pop 2). We find

$$\bar{Y}_1 = 161.0cm, \quad S_1 = 12.2cm \quad (53)$$

$$\bar{Y}_2 = 164.1cm, \quad S_2 = 13.1cm \quad (54)$$

$$(55)$$

Find a 95% confidence interval for  $\mu_1 - \mu_2$ .

$$\sigma_{\bar{Y}_1 - \bar{Y}_2}^2 = \frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2} \approx \frac{S_1^2}{n_1} + \frac{S_2^2}{n_2} = \frac{(12.2)^2}{1000} + \frac{(13.1)^2}{1000} = 0.320 \quad (56)$$

So  $\sigma_{\bar{Y}_1 - \bar{Y}_2} = 0.566$ . For 95% confidence interval  $z_c = 1.96$ . So confidence interval is

$$(161.0 - 164.1) \pm 1.96 * 0.566 = -3.1 \pm 1.109 \quad (57)$$

**Example:** (from the book) We have two brands of refrigerators. Both have a one year warranty. For brand A we find that in a sample of 50, 12 fail during the warranty period. For brand B we find that in a sample of 60, 12 fail during the warranty period. We define  $p_A, p_B$  to be the probability of a failure. Estimate  $p_A - p_B$  with a confidence interval with confidence coefficient of 98%.

$$\hat{p}_A = 12/50 = 0.24, \hat{p}_B = 12/60 = 0.20 \quad (58)$$

Critical  $z$  is 2.33. Confidence interval is

$$(0.24 - 0.20) \pm 2.33 \sqrt{\frac{(0.24)(0.76)}{50} + \frac{(0.20)(0.80)}{60}} \quad (59)$$

$$= 0.04 \pm 0.1851 = [-0.1451, 0.2251] \quad (60)$$

## 8.7 Choosing the sample size

If we have a sample of size  $n$  and we want a certain confidence coefficient  $1 - \alpha$ , then we have seen how to find the confidence interval  $[\hat{\theta}_L, \hat{\theta}_U]$  for several different estimates. Here we consider a different question. Suppose we have not drawn the sample yet and we want a given confidence coefficient and we want the confidence interval to be a certain size. Note that  $\hat{\theta}_L$  and  $\hat{\theta}_U$  depend on the sample size, and so we can use the given width for the confidence interval to solve for the sample size.

**Example** We continue the example comparing the average heights of adult females in the US and Canada. With sample sizes of 1000 we found a confidence interval of  $-3.1 \pm 1.109$ . So the width was 2.22. Suppose we want the width of a 95% confidence interval to be 1.0 and we are going to use samples of the same size. What should  $n$  be?

Answer: We had

$$S_1 = 12.2cm, \quad S_2 = 13.1cm \quad (61)$$

$$(62)$$

We continue to approximate

$$\sigma_{\bar{Y}_1 - \bar{Y}_2}^2 = \frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2} \approx \frac{S_1^2}{n} + \frac{S_2^2}{n} = \frac{(12.2)^2}{n} + \frac{(13.1)^2}{n} = \frac{320.45}{n} \quad (63)$$

So  $\sigma_{\bar{Y}_1 - \bar{Y}_2} = 17.9/\sqrt{n}$ . For 95% confidence interval  $z_c = 1.96$ . The width of the confidence interval will be  $2z_c\sigma_{\bar{Y}_1 - \bar{Y}_2} = 70.17/\sqrt{n}$ . So  $n = (70.17)^2 = 4925$ .

**Example** A marketing firm wants to test interest in a new product. They will ask the people sampled one question - would they buy the product or not. They want to estimate the fraction  $p$  of people who say yes with a 90% confidence interval and want the estimate to be within 2% of the true value, i.e., they want the width of the confidence interval to be 4%.

Answer: We have no idea what  $p$  really is, so how do we find  $\sigma_{\hat{p}}$ ? It is given by  $\sqrt{p(1-p)}/\sqrt{n}$ . The maximum of  $p(1-p)$  over  $0 \leq p \leq 1$  occurs when  $p = 1/2$ . So we have  $\sigma_{\hat{p}} \leq 1/(2\sqrt{n})$ . For a 90% confidence interval,  $z_c = 1.645$ . We want  $z_c/(2\sqrt{n}) = 0.02$ . This leads to  $n = 1692$ .

## 8.8 Small sample confidence intervals

In the section on large sample confidence intervals, we used the fact that the sample was large in two ways. First, a large sample size implies that  $\bar{Y}$  will be approximately normally distributed. Second, a large sample size allowed us to approximate the unknown population variance by the sample variance.

We now assume the sample size is not large. We assume the population is normal (or at least close to normal). We review the definition of the t-distribution. Let  $Z$  and  $W$  be independent RV's and suppose  $Z$  has a standard normal distribution and  $W$  has a  $\chi^2$  distribution with  $\nu$  degrees of freedom. Then

$$T = \frac{Z}{\sqrt{W/\nu}} \quad (64)$$

has a t-distribution with  $\nu$  degrees of freedom. If the population is normal, then  $\bar{Y}$  is normal with mean  $\mu$  and variance  $\sigma^2/n$ . And  $(n-1)S^2/\sigma^2$  has a  $\chi^2$  distribution with  $n-1$  degrees of freedom. So

$$T = \frac{\bar{Y} - \mu}{\sqrt{S^2/n}} \quad (65)$$

has a t-distribution with  $n-1$  degrees of freedom.

The statistic  $T$  depends on the random sample through  $\bar{Y}$  and  $S^2$ . It depends on only one unknown parameter, namely  $\mu$ . And the distribution

of  $T$  does not depend on  $\mu$ . So  $T$  is a pivotal quantity for  $\mu$ . So we can use it to construct a confidence interval for  $\mu$  with confidence coefficient  $1 - \alpha$ . Recall that the t-distribution is symmetric. So when we seek  $a$  and  $b$  such that  $P(a \leq T \leq b)$ , it is natural to take  $a = -b$ . Let  $t_{\alpha/2}$  be defined by

$$P(T > t_{\alpha/2}) = \alpha/2 \quad (66)$$

Note that  $t_{\alpha/2}$  depends on the number of degree of freedom,  $n - 1$ . So we have

$$P(-t_{\alpha/2} \leq T \leq t_{\alpha/2}) = \alpha \quad (67)$$

The next step in finding the confidence interval is to rearrange  $-t_{\alpha/2} \leq T \leq t_{\alpha/2}$  so that it gives a confidence interval for  $\mu$ . We find that  $-t_{\alpha/2} \leq T \leq t_{\alpha/2}$  is equivalent to

$$\mu \in [\bar{Y} - t_{\alpha/2}S/\sqrt{n}, \bar{Y} + t_{\alpha/2}S/\sqrt{n}] \quad (68)$$

So the confidence interval with confidence coefficient  $1 - \alpha$  is

$$\bar{Y} \pm t_{\alpha/2}S/\sqrt{n} \quad (69)$$

---

### End of lecture on Tues, 2/6

---

**Example:** A random sample of 16 Americans found that the amount of beef they ate (in lbs) in the past year was,

118, 115, 125, 110, 112, 130, 117, 112, 115, 120, 113, 118, 119, 122, 123, 126

You can enter this in R with

$$x = c(118, 115...126) \quad (70)$$

Then

$$\text{mean}(x) = 118.4, \quad \text{sd}(x) = 5.656, \quad \text{qt}(0.025, 15) = -2.131 \quad (71)$$

So we get a confidence interval of  $118.4 \pm 3.106$ . Note that if we had pretended 16 was a big sample we would have used 1.96 in place of 2.131 and gotten  $118.4 \pm 2.772$ .



We now consider a confidence interval for the difference of two population means  $\mu_1 - \mu_2$  when the sample sizes are not large. We assume that each of the two populations are normal (or close to normal) and that the samples drawn from the populations are independent of each other. We know that  $\bar{Y}_1$  is independent of  $S_1^2$  and that  $\bar{Y}_2$  is independent of  $S_2^2$ . We also know that  $(n_1 - 1)S_1^2/\sigma_1^2$  and  $(n_2 - 1)S_2^2/\sigma_2^2$  have  $\chi^2$  distributions. We want to replace

$$\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2} \rightarrow \frac{S_1^2}{n_1} + \frac{S_2^2}{n_2} \quad (72)$$

But ...

If we assume that the two populations have the same variance, then we can make progress. So we assume  $\sigma_1^2 = \sigma_2^2$ . We denote this common variance by  $\sigma^2$ . We estimate this variance by “pooling our two samples.” Define

$$S_p^2 = \frac{(n_1 - 1)S_1^2 + (n_2 - 1)S_2^2}{n_1 - 1 + n_2 - 1} \quad (73)$$

Note that  $(n_1 + n_2 - 2)S_p^2/\sigma^2$  is the sum of two independent  $\chi^2$  distributions with  $n_1 - 1$  and  $n_2 - 1$  d.f. So it has a  $\chi^2$  distribution with  $n_1 + n_2 - 2$  d.f. Also, it is independent of  $\bar{Y}_1 - \bar{Y}_2$ . After a bit of algebra this implies

$$T = \frac{\bar{Y}_1 - \bar{Y}_2 - (\mu_1 - \mu_2)}{S_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} \quad (74)$$

has a t-distribution with  $n_1 + n_2 - 2$  d.f. It is a pivotal quantity for  $\mu_1 - \mu_2$ . So we can use it to find a confidence interval for  $\mu_1 - \mu_2$ . We define  $t_{\alpha/2}$  as before. The event  $-t_{\alpha/2} \leq T \leq t_{\alpha/2}$  leads to the confidence interval

$$(\bar{Y}_1 - \bar{Y}_2) \pm t_{\alpha/2} S_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}} \quad (75)$$

**Example:** A drug is supposed to lower blood pressure. 10 subjects are divided into two groups. 6 subjects (group 1) are given the drug over some period of time and 4 subjects (group 2) are not. At the end of the trial their systolic blood pressure is measured.

$$\bar{Y}_1 = 117.5, \quad S_1 = 9.7, \quad (76)$$

$$\bar{Y}_2 = 126.8, \quad S_2 = 12.0 \quad (77)$$

Find a 95% confidence interval for  $\mu_1 - \mu_2$ . Does this study provide evidence that the drug lowers blood pressure?

Describe what population 1 and 2 are.

We assume that the variances of the the two populations are the same. The pooled estimator for this common variance is

$$S_p^2 = \frac{(n_1 - 1)S_1^2 + (n_2 - 1)S_2^2}{n_1 - 1 + n_2 - 1} \quad (78)$$

$$= \frac{(6 - 1)(9.7)^2 + (4 - 1)(12.0)^2}{6 - 1 + 4 - 1} = 112.8 \quad (79)$$

So  $S_p = 10.6$ . The number of degrees of freedom is  $6 + 4 - 2 = 8$ . The critical  $t_c$  is given by

$$t_c = -qt(0.025, 8) = 2.306 \quad (80)$$

The confidence interval works out to  $-9.3 \pm 15.77$ . We cannot conclude from this that  $\mu_1 - \mu_2 < 0$ . So this study does not provide evidence that the drug lowers blood pressure. (This does not mean that the study shows it does not lower blood pressure.)

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

In this section we consider a confidence interval for the population variance  $\sigma^2$ . We have already seen a point estimator for this parameter, namely, the sample variance

$$S^2 = \frac{1}{n - 1} \sum_{i=1}^n (Y_i - \bar{Y})^2 \quad (81)$$

We assume that the population is normal. Recall that under this assumption

$$U = \frac{(n - 1)S^2}{\sigma^2} \quad (82)$$

has a  $\chi^2$  distribution with  $n - 1$  degrees of freedom. Note also that  $U$  only involves one unknown parameter,  $\sigma^2$ . So  $U$  is a pivotal quantity. So to find a confidence interval for  $\sigma^2$  we first seek numbers  $\chi_L^2$  and  $\chi_U^2$  such that

$$P(\chi_L^2 \leq U \leq \chi_U^2) = 1 - \alpha \quad (83)$$

As we have encountered before, there are many solutions of this equation. One natural way to pick out one would be to find the solution that minimizes  $\chi_U^2 - \chi_L^2$ . This will lead to the smallest possible confidence interval. But there is no explicit way to compute this solution. We could find it by trial and error.

We take a simpler approach and choose  $\chi_u^2$  and  $\chi_L^2$  so that

$$P(U < \chi_L^2) = \alpha/2, \quad P(U > \chi_U^2) = \alpha/2 \quad (84)$$

Using the definition of  $U$ , eq. (83) implies

$$P((n-1)S^2/\chi_U^2 \leq \sigma^2 \leq (n-1)S^2/\chi_L^2) = 1 - \alpha \quad (85)$$

So our confidence interval for  $\sigma^2$  is

$$[(n-1)S^2/\chi_U^2, (n-1)S^2/\chi_L^2] \quad (86)$$

Often  $\chi_L^2$  is written as  $\chi_{\alpha/2, n-1}^2$  and  $\chi_U^2$  is written as  $\chi_{1-\alpha/2, n-1}^2$ .

**Example:** A population is approximately normally distributed. I want to estimate the population variance  $\sigma^2$ . I pick a random sample of size 10 and find a sample mean of 102.7 and a sample variance of 4.2. Find a 95% confidence interval for the population variance  $\sigma^2$ . We have  $\alpha = 0.05$ . So

$$\chi_L^2 = qchisq(0.025, 9) = 2.70, \quad \chi_U^2 = qchisq(0.975, 9) = 19.0, \quad (87)$$

So the confidence interval is

$$\left[ \frac{9 \times 4.2}{19.0}, \frac{9 \times 4.2}{2.70} \right] = [1.99, 14.0] \quad (88)$$

We can get a 95% confidence interval for  $\sigma$  by taking square roots  $[\sqrt{1.99}, \sqrt{14}] = [1.41, 3.74]$ .

---

**End of lecture on Thurs 2/8**

---