

## Lecture 7

### 7.1 Sampling, sample statistics and sampling distributions

- Let's think about our journey through the world of inferential statistics so far. Our goal the entire time has been to come up with the best possible way to make inferences about populations from samples. To do this, we have:
  - recast social phenomena as *experiments* that yield observations for analysis.
    - \* we called the outcomes of these experiments *events*
    - \* and defined the *sample space* of an experiment as the set of all the events that are possible
    - \* and considered carefully how we assign probabilities to events of interest.
  - we then defined a *random variable* as a function mapping a sample space to the real numbers
    - \* and then discussed the ways we describe the probability distribution of a random variable...
    - \* and the probability distribution of a *function* of random variables.
- In this context, observed social phenomena (election results, outbreaks of war, passage of legislation) can all be considered realizations of random variables.
  - Let's be a little more clear about this. The phenomena that social scientists study are random events.
  - This may sound odd: in common usage we say something is "random" when it cannot be anticipated.
  - Social scientists use this term differently. When we say an event is random, we mean that it is probabilistic rather than deterministic.

- For a random event, we attempt to specify the causal processes that alter the chance that it occurs. But we cannot specify the causal process that guarantee the event will occur. If we could, we would be studying a deterministic event.
  - \* An election is a good example of a random event. The best we can do is develop a theory that specifies factors important to determining election winners and then hypothesize that these factors change the odds of a particular result.
- Thus, saying that the variables in our theories are random variables amounts to saying that we expect that the values of these variables that we observe are draws from from an associated probability distribution.
- Finally, we arrived at a very powerful result. If the probability distribution is an accurate representation of the population frequency distribution, then the *expected value* of a random variable is the population mean,  $\mu$ , where we define expected value as

$$E(Y) \equiv \sum_y yp(y) \text{ in the discrete case and}$$

$$E(Y) \equiv \int_y yf(y) dy \text{ in the continuous case.}$$

- Now it's time to put all this theory to work helping us undertake the fundamental challenge we face in statistics: making inferences from samples to populations.
- The relevance of what we've learned previously to this task is that estimates are almost always functions of the  $n$  random observations that appear in a sample. Therefore, they are:
  - outcomes of experiments that are themselves random variables with their own probability distributions.
- If we can be very specific about the process giving rise to the sample, we can develop an **estimator** to make inferences from the sample to the population.
- An **estimator** is a *rule*, often expressed as a formula, that tells us how to calculate an estimate from a sample.
- As an example, consider the following as an estimator for the population mean,  $\mu$  :

- Draw a random sample of  $n$  observations,  $y_1, y_2, \dots, y_n$ , from the population and employ the observed sample mean

$$\bar{Y} \equiv \frac{y_1 + y_2 + \dots + y_n}{n} = \frac{1}{n} \sum_i y_i$$

- As an estimator for  $\mu$ ,  $\bar{Y}$  seems pretty intuitive. But let's be precise about just how good this estimate is for  $\mu$ .

- To do this, we first think of each of the  $n$  draws that gave rise to the sample as a realization of a random variable.
- An example: we want to know the mean income of the American population,  $\mu$ .
- To do so, I draw a sample of the American population and ask each person in the sample one question: what is your annual income?
- Denote the response given by the first person in my sample, draw number 1, as  $y_1$ .
- Now here's the key concept:  $y_1$  is one of literally millions of responses I could have observed in my first survey participant. Therefore it is a realization of the random variable  $Y_1$ .
- Now do this again: pick participant number 2 and ask him or her the same question. My observation  $y_2$  is again one of millions of responses I could have observed for respondent number 2. Therefore it is a realization of the random variable  $Y_2$ .
- So if we think about it this way, a sample of  $n$  observations is the realization of how many random variables?  $n$ : from  $Y_1, Y_2, \dots, Y_n$ .
- OK, now let's think again about  $\bar{Y}$ .

$$\bar{Y} = \frac{1}{n} \sum_i Y_i,$$

- In this context, our sample of  $n$  observations is just one possible realization of  $\bar{Y}$ .
- That is,  $\bar{Y}$  is a *function* of random variables, and therefore is *itself* a random variable.

- Note use of capital letters here, we are talking about theoretical, not observed, quantities. In keeping with our convention, little  $\bar{y}$  is a realization of the random variable  $\bar{Y}$ .
- Because it is a random variable,  $\bar{Y}$  has a probability distribution. To specify it, we will stipulate some simple but powerful assumptions that hold whenever we have a **random sample**. (Recall how we defined random sample: each of the  $\binom{N}{n}$  different possible samples has an equal probability of being drawn.)
- This generates the canonical case, in which we have a function of the random variables  $Y_1, Y_2, \dots, Y_n$  observed in a random sample from a population of interest.
  - When we have a random sample from a large enough population, we can safely assume that the RVs  $Y_1, Y_2, \dots, Y_n$  are independent and identically distributed (“i.i.d.”).
    - \* To recall:  $Y_1, Y_2, \dots, Y_n$  are said to be **independent** iff

$$F(y_1, y_2, \dots, y_n) = F_1(y_1)F_2(y_2) \cdot \dots \cdot F_n(y_n) \text{ for every } n\text{-tuple } (y_1, y_2, \dots, y_n), \text{ and}$$

- Now we introduce another concept:  $Y_1, Y_2, \dots, Y_n$  are said to be **identically distributed** iff

$$F_1(y_1) = F_2(y_2) = \dots = F_n(y_n) = F(y) \text{ for } y_1, y_2, \dots, y_n.$$

- In this context,  $\bar{Y} = \frac{1}{n} \sum_i Y_i$  is a **sample statistic**, which we define as a function of the observable random variables in a sample and known constants.
- We know (from the handout last time) that  $E(\bar{Y}) = \mu$ , [note that this is due to the identity assumption] and so we can be assured that, on average,  $\bar{Y}$  should equal  $\mu$ . So now let’s ask, how good of an estimate is it?
  - A straightforward measure of “goodness” would be how far off we can expect any realization of  $\bar{Y}$  to be from the population mean,  $\mu$ .
  - And since the mean of  $\bar{Y}$  is  $\mu$ , this quantity is just the standard deviation of  $\bar{Y}$ , or  $\sigma_{\bar{Y}}$ . But THIS of course is just  $\sqrt{\text{VAR}(\bar{Y})}$ , which we showed last time to be  $\sqrt{\frac{\sigma^2}{n}} = \frac{\sigma}{\sqrt{n}}$ . [recall that we need BOTH identity and independence to achieve this result.]

- Well, because sample statistics are themselves random variables, they have probability distributions (recall: table, graph, formula). We have a special term for the probability distributions of sample statistics: **sampling distributions**. The sampling distribution of a sample statistic is a theoretical model for the possible values of the statistic we would expect to observe through repeated sampling.
- The expected value and the variance of a sample statistic are important properties of the statistic's sampling distribution. For example, here are graphs depicting two sampling distributions of  $\bar{Y}$ :  $\bar{Y}$  and  $\bar{Y}'$ :

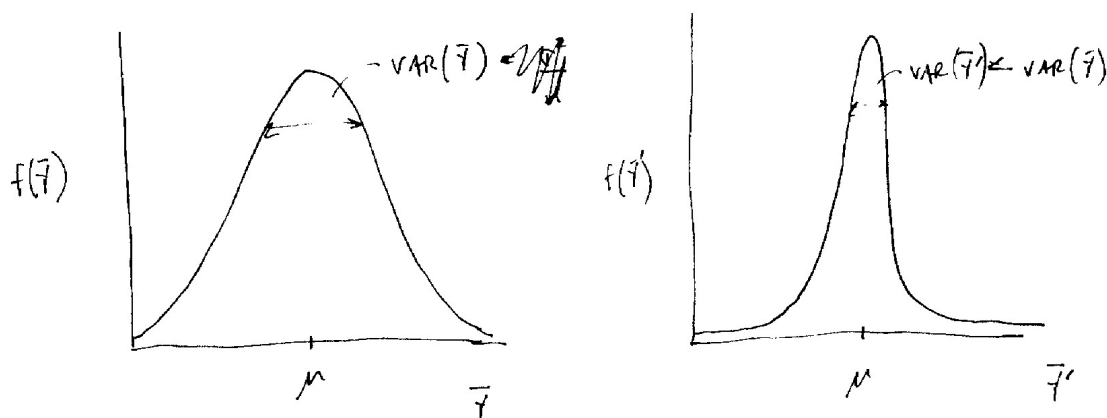


Figure 1: Variance of  $\bar{y}$

- $\bar{Y}$  and  $\bar{Y}'$  have the same expected value:  $\mu$ . But clearly  $\bar{Y}'$ , which has a smaller variance, is a better estimate of  $\mu$  than  $\bar{Y}$ . It's closer, on average, to  $\mu$  than  $\bar{Y}$ .
- So:
  - we've got  $\bar{Y}$ , an estimate for  $\mu$  that, on average, equals  $\mu$ .
  - we now also know how good an estimate this is: on average, it is  $\frac{\sigma}{\sqrt{n}}$  units away from  $\mu$ .
- But now an additional question arises: what does the distribution of the random variable  $Y$  look like? Often, we'll be drawing samples from populations about whose frequency

distribution we have no idea. E.g., I want to estimate  $\mu$ , the mean number of potatoes eaten by the average American per month. Does the distribution of  $Y$  look like this? Or this? Or this?

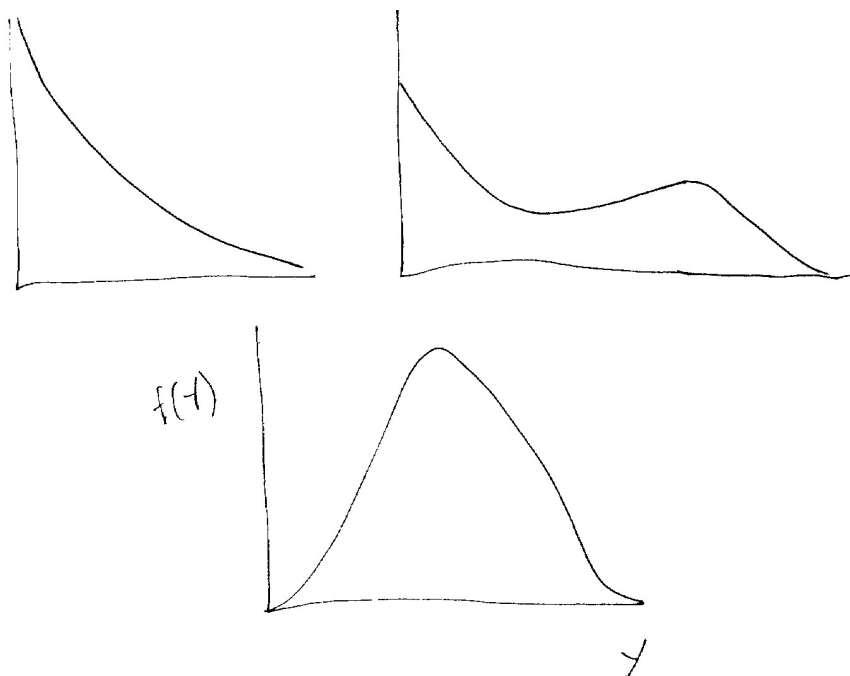


Figure 2: Distributions

- If possible, we'd really like to avoid making assumptions about the shape of the population distribution of  $Y$  in order to describe the distribution of  $\bar{Y}$ .
- And, it turns out, we can...

## 7.2 The Central Limit Theorem

- ...because if my sample size is large enough, I don't need any assumptions about the distribution of  $Y$  to describe the sampling distribution of  $\bar{Y}$ . For it turns out that if  $Y_1 \dots Y_n$  are i.i.d, then:
  - $\bar{Y}$  has a sampling distribution that is approximately Normal
  - as the sample size becomes large.
  - This is the **central limit theorem**.

- Before digging into the math (and we won't dig that far), let's have a look at a few distributions that convey the intuition (go over handout: "Probability Distributions of  $Y$  and Simulated Sampling Distributions of  $\bar{Y}$ ")
- The CLT is more formally stated as follows:
- Let  $Y_1, Y_2, \dots, Y_n$  be i.i.d. random variables with  $E(Y_i) = \mu$  and  $VAR(Y_i) = \sigma^2$ . Define

$$U_n \equiv \frac{\bar{Y} - \mu}{\sigma/\sqrt{n}}.$$

That is, we standardize  $\bar{Y}$  by (1) subtracting its hypothesized mean ( $\mu$ ) and then (2) dividing this difference by  $\bar{Y}$ 's standard deviation ( $\frac{\sigma}{\sqrt{n}}$ ). Then the CDF of  $U_n$  converges in probability to (where "converges in probability to" means "as  $n$  becomes large it is distributed as") the standard normal CDF. That is,

$$\lim_{n \rightarrow \infty} F_{U_n}(u) = \lim_{n \rightarrow \infty} P(U_n \leq u) = P(Z \leq u) = \int_{-\infty}^u \frac{1}{\sqrt{2\pi}} e^{-t^2/2} dt \text{ for all } u.$$

- Take a moment to appreciate the power of the CLT. It yields the sampling distribution of  $\bar{Y}$  *without requiring any assumptions about the probability distribution of  $Y$* . We will not cover the proof of the CLT in this class, as it requires knowledge of moment generating functions, which itself requires knowledge of Taylor series expansions. Those equipped with these tools who want to satisfy themselves with a proof may see section 7.4 of the text.

### 7.3 Estimation

- With the sampling distributions yielded by the CLT in hand, we can develop **estimators** of population parameters. An estimator is a rule—often expressed as a formula—that tells us how to calculate an estimate of a population parameter.
  - Two kinds of estimates that we'll focus on here include
    - \* **point estimates**—in which a single value, or point, is given as the estimate of the parameter

\* **interval estimates**-in which two values are used to construct an interval that we believe contains/traps/encloses the parameter of interest.

- So the sample mean,  $\bar{Y} = \frac{1}{n} \sum_i Y_i$ , is one possible point estimator of the population mean  $\mu$ .
- But there are lots of others. For example, consider the estimator  $\bar{Y}_B = \frac{1}{n} \sum_i (Y_i + 1)$ .
- Intuitively, we know this is a worse estimator than  $\bar{Y}$ . But can we put some meat on this intuition? We do this by specifying two desirable criteria for evaluating a potential estimator:
  - 1. unbiasedness
  - 2. (relatively) small variance, which is also known as **efficiency** or **precision**.
- First, some terminology. We typically write the population parameter for which we seek a point estimate as  $\theta$ , and a proposed estimator for this parameter as  $\hat{\theta}$ .
- Now, here is why we've spent a fair amount of time learning about the math of expectations: we can use these tools to show whether an estimator is unbiased and to determine how precise it is.
- In fact, we define an estimator as unbiased if—in expectation—it is equal to the parameter it claims to estimate. That is, an estimator is unbiased if the expected value of its distribution is the parameter. Formally,
  - $\hat{\theta}$  is an **unbiased estimator** for  $\theta$  if  $E(\hat{\theta}) = \theta$ .
  - If  $E(\hat{\theta}) \neq \theta$ , then we say  $\hat{\theta}$  is **biased**.
  - The **bias** of a point estimator is given by  $B(\hat{\theta}) = E(\hat{\theta}) - \theta$ .
- E.g. So one way to say that the estimator  $\bar{Y}_B$  isn't a good estimator is to show that it is a



biased estimator for  $\mu$ . We do that by showing that  $E(\overline{Y_B}) \neq \mu$ , and thus that  $B(\overline{Y_B}) \neq 0$ .

$$\begin{aligned}
 E(\overline{Y_B}) &= E \left[ \frac{1}{n} \sum_i (Y_i + 1) \right] \\
 &= \frac{1}{n} \left\{ \left[ \sum_i E(Y_i) \right] + nE(1) \right\} \\
 &= \frac{1}{n} \{ [n\mu] + n \} \\
 &= \mu + 1 \neq \mu, \text{ and } B(\overline{Y_B}) = 1.
 \end{aligned}$$

- So much for the first criterion. Our second criterion is that we'd like our estimator to be as close as possible to  $\theta$  in repeated sampling. In mathematical terms, we of course want the variance of the sampling distribution of our estimator to be as small as possible. Recalling the definition of the variance of a random variable, we write the variance of an estimator  $\hat{\theta}$  as  $VAR(\hat{\theta}) = E \left\{ [\hat{\theta} - E(\hat{\theta})]^2 \right\}$ . In an ideal world, we wish for this to be as small as possible.
- Sometimes we face a tradeoff between reducing an estimator's bias and reducing its variance. One way to evaluate the tradeoff is to minimize an estimator's **mean square error (MSE)**, which is defined as the expected value of the square of the distance between the estimator and the parameter

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

- It can be shown that the MSE is equal to the sum of an estimator's variance plus the square of its bias:

$$MSE(\hat{\theta}) = VAR(\hat{\theta}) + [B(\hat{\theta})]^2.$$

– Ask: Why does this measure include the square of its bias?

## 7.4 An example

- For example, let us say that we wish estimate the parameter  $\mu_1 - \mu_2$ , the difference in means of two different populations drawn from independent samples. Is the intuitive estimator,

$\bar{Y}_1 - \bar{Y}_2$ , unbiased? Let's see:

$$\begin{aligned} E(\bar{Y}_1 - \bar{Y}_2) &= E(\bar{Y}_1) - E(\bar{Y}_2) \\ &= \mu_1 - \mu_2. \end{aligned}$$

- Yes,  $\bar{Y}_1 - \bar{Y}_2$  is an unbiased estimator for  $\mu_1 - \mu_2$ . Now what is its variance? Well,

$$\begin{aligned} \text{VAR}(\bar{Y}_1 - \bar{Y}_2) &= \text{VAR}(\bar{Y}_1) + \text{VAR}(\bar{Y}_2) + 2\text{COV}(\bar{Y}_1, \bar{Y}_2) \\ &= \text{VAR}(\bar{Y}_1) + \text{VAR}(\bar{Y}_2) \quad [\bar{Y}_1, \bar{Y}_2 \text{ independent}] \\ &= \frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}, \quad [\text{variance of } \bar{Y}\text{-bar}] \end{aligned}$$

where  $\sigma_1^2$  and  $\sigma_2^2$  are the variances of populations 1 and 2 respectively.

- Note that we can talk about the standard error of an estimator  $\hat{\theta}$ , which we write  $\sigma_{\hat{\theta}}$ . It is just the square root of the variance of the estimator. The standard error of the estimator  $\bar{Y}_1 - \bar{Y}_2$  is thus  $\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$ .
- Another variant of the CLT (which does not require identicality, but instead requires some other conditions) tells us that the sampling distribution of a sum of independent random variables (like  $\bar{Y}_1$  and  $\bar{Y}_2$ ) approximates the Normal as  $n$  becomes large. And so

$$\bar{Y}_1 - \bar{Y}_2 \sim N\left(\mu_1 - \mu_2, \frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}\right) \text{ as } n_1, n_2 \rightarrow \infty.$$