

*October 24, 2017*

---

# CS 361: Probability & Statistics

Samples, populations, and  
intervals

---

# Samples, Urns, and Populations

---

# Populations and samples

---

Another set of problems in which point estimates arise is when we only have a subset of the data that we could have. If we want to know the average heights of students, we probably won't have the means to measure the height of every student. If we want to predict how the public is going to vote, we probably won't be able to afford to ask every single potential voter

The terminology we will use to talk about the set of data items we can actually measure is a **sample**. And the term for the entire dataset to which we may not have access is called the **population**

So the question becomes under what conditions are measurements of the sample likely to be good approximations of doing the same measurements on the whole population?

---

# Samples and populations

---

For the next few slides we will suppose that our sample is of size  $k$  and that our population is of size  $N_p$ . And we will suppose that  $k$  is much smaller than  $N_p$

How we obtain the sample determines the characteristics of the sample relative to the population. To model sampling, we will assume that each item of the sample is chosen independently and fairly from the population.

An classic way of thinking about this is by imagining that each potential data item in the population corresponds to a ticket and that we are drawing these tickets one by one from an urn, writing down what we draw, and then putting the ticket back in the urn

This is also known as **sampling with replacement**

---

# Samples and populations

---

With this model for sampling we will see whether the sample can tell us anything about the population. In particular, we will focus on the mean of the population and estimating it by looking at our sample

The **population mean**, which we can't actually observe, is indicated with the notation  $\text{popmean}(\{x\})$ . If we could look at the entire population we would calculate this directly using our formula for the mean of a dataset from Chapter 2

Instead, we have a sample of size  $k$ . The mean of our sample is called the **sample mean**

The **sample mean** is a random variable (it could be different depending on which  $k$  items we get when we randomly sample). We will write  $\chi^{(k)}$  for this random variable and determine its value for a given sample according to our formula for the mean of a dataset

---

# Samples and populations

---

Let's try and relate the sample mean to the population mean

$\chi^{(k)}$  is the random variable that outputs the sample mean of a sample of size  $k$

The value of the sample mean for a given sample is given by

$$\chi^{(k)} = \frac{1}{k}(X_1 + X_2 + \dots + X_k)$$

If we take expectations we get

$$E[\chi^{(k)}] = \frac{1}{k}(E[X^{(1)}] + E[X^{(1)}] + \dots + E[X^{(1)}])$$

Or

$$E[\chi^{(k)}] = E[X^{(1)}]$$

---

# Samples and populations

---

But we can reason about the expected value of  $X^{(1)}$  on the basis of the way we've defined how we are sampling

$X^{(k)}$  is the random variable that outputs the sample mean of a sample of size  $k$

Using the definition of expectation—that we just sum over every possible value of a random variable and multiply by its probability—we have

$$E[X^{(1)}] = \sum_{i \in 1, \dots, N_p} x_i p(i)$$

We've assumed we draw fairly from the urn, so this can be rewritten

$$E[X^{(1)}] = \sum_{i \in 1, \dots, N_p} x_i \frac{1}{N_p}$$

---

# Samples and populations

---

$$E[X^{(1)}] = \sum_{i \in 1, \dots, N_p} x_i \frac{1}{N_p}$$

Rewriting

$$E[X^{(1)}] = \frac{\sum_{i \in 1, \dots, N_p} x_i}{N_p}$$

Which is the formula for the mean of the population, i.e.  $E[X^{(1)}] = \text{popmean}(\{x\})$

We showed earlier that  $E[X^{(k)}] = E[X^{(1)}]$  thus we've shown the following

$$E[X^{(k)}] = \text{popmean}(\{x\})$$

---

# Estimators

---

$$E[X^{(k)}] = \text{popmean}(\{x\})$$

This was a nice result. In general when we find a random variable that has as its expectation something we are trying to estimate, we call the random variable an **unbiased estimator** of that quantity

Something we will note but not show here is that the computing the variance of the sample dataset with the variance formula applied to our sample

$$\frac{\sum_{i=1}^k (x_i - X^{(k)})^2}{k}$$

gives us a random variable, but this random variable is not an unbiased estimator of the population variance. I.e. the expected value of the above does not equal the population variance

---

# Estimators

---

The quantity given by

$$\frac{\sum_{i=1}^k (x_i - \bar{x}^{(k)})^2}{k - 1}$$

however, is an unbiased estimate.

In other words, the above quantity is a random variable (it is different depending on exactly which sample of  $k$  items we happen to draw) with expectation equal to the population variance

When we refer to the **sample variance** or the **sample standard deviation** we will be referring to this quantity or, respectively, its square root

---

# Sample mean - variance

---

So we know that the expected value for the sample mean is equal to the population mean, which means that for a large  $k$  (sample size) the sample mean will be a good estimate of the population mean. In order to get a sense of what constitutes a large  $k$  we should look at the variance / standard deviation of the random variable corresponding to sample mean

If we do a little algebra we find that the variance and standard deviation of the random variable  $X^{(k)}$  are. (Note this is not the variance of the sample or the population)

$$\begin{aligned}\text{var}[X^{(k)}] &= \frac{\text{popstd}(\{x\})^2}{k} \\ \text{std}(\{X^{(k)}\}) &= \frac{\text{popstd}(\{x\})}{\sqrt{k}}\end{aligned}$$

where popstd is the the standard deviation of the population, which we cannot measure

---

# Sample mean - variance

---

$$\text{std}(\{X^{(k)}\}) = \frac{\text{popstd}(\{x\})}{\sqrt{k}}$$

The standard deviation of the estimate of the mean is sometimes called **the standard error** of our estimate. Note that we cannot calculate the above exactly. Though it still tell us some interesting things

A larger sample size reduces the error of our estimate of the popmean but that this reliability grows less than linearly in the sample size

It also says that if we are trying to estimate the mean of a population which itself has a large amount of variance, we will need a larger number of samples

Note that our error does not depend on the size of the population. Whether the population is 1000 or 10,000,000 this is our error

---

# Distribution of the sample mean

---

The sample mean is a random variable. In particular it is a random variable that we get by summing together a bunch of independent and identically distributed random variables—our samples. From the central limit theorem, then, we can expect that our estimate—the sample mean—has a normal distribution.

That is, if we took a sample of size  $k$ , recorded the sample mean, then did this again and again the numbers we write down would be normal data

We have just shown that the mean of this normal random variable is  $\text{popmean}(\{x\})$  and the standard deviation is  $\frac{\text{popstd}(\{x\})}{\sqrt{k}}$

---

# Urn model, samples, and populations

---

If we have managed to sample the population with a procedure that corresponds to this urn model — independent samples taken fairly from the population—then we see that estimating the mean is fairly straightforward

If our samples are not independent or aren't taken fairly from the population, however, determining what our sample tells us about the population is more difficult.

If I want to estimate the average height of people in Champaign and I do my sampling at the local daycare or at the university's basketball practice, my sample mean is not going to be a good estimate of the population mean, no matter how large it is!

# Interval estimates

---

# Interval estimates

---

So far we have dealt with point estimation or point inference, where we look at a dataset and get out a number estimating some quantity of interest

It can also be desirable to give a range for a quantity we are trying to estimate, along with the probability that the true value of the quantity lies somewhere in the range

If we have an interval that contains the true value with some probability we call the interval and probability a **confidence interval**

---

# Confidence intervals for samples

---

Suppose we have a population and are sampling from that population

Recall that the sample mean  $\bar{x}^{(k)}$  is a normal random variable with expectation  $\text{popmean}$  and standard deviation  $\frac{\text{popstd}(\{x\})}{\sqrt{k}}$

Given what we know about normal random variables. This tells us that for about 68% of the samples we will have

$$\text{popmean} - \frac{\text{popstd}}{\sqrt{k}} \leq \bar{x}^{(k)} \leq \text{popmean} + \frac{\text{popstd}}{\sqrt{k}}$$

# Confidence intervals for samples

It turns out that the variance for our estimate of popsd can be assumed to be small for reasonable sized  $k$ . So we can replace popsd in the interval with the sample standard deviation

And say that the true value of the population mean is bounded by

Sample standard deviation


$$X^{(k)} - \frac{sd(\{x\})}{\sqrt{k}} \leq \text{popmean} \leq X^{(k)} + \frac{sd(\{x\})}{\sqrt{k}}$$

in 68% of samples. Or in 95% of samples we have

$$X^{(k)} - 2 \frac{sd(\{x\})}{\sqrt{k}} \leq \text{popmean} \leq X^{(k)} + 2 \frac{sd(\{x\})}{\sqrt{k}}$$

---

# So what?

---

So what does all that mean? It means we have a recipe for how we can report an interval when estimating the population mean as opposed to what we learned earlier today in the context of point estimates where we just reported a single number, the sample mean

First we calculate the sample mean  $\bar{x}^{(k)}$

Then we calculate the sample standard deviation using our unbiased estimator from before

$$\sqrt{\frac{\sum_{i=1}^k (x_i - \bar{x}^{(k)})^2}{k - 1}}$$

Then we form our standard error by dividing the sample standard deviation by the square root of k

Then we report the sample mean plus or minus the standard error for a 68% confidence interval or  $2 \times \text{stderr}$  for a 95% confidence interval

---

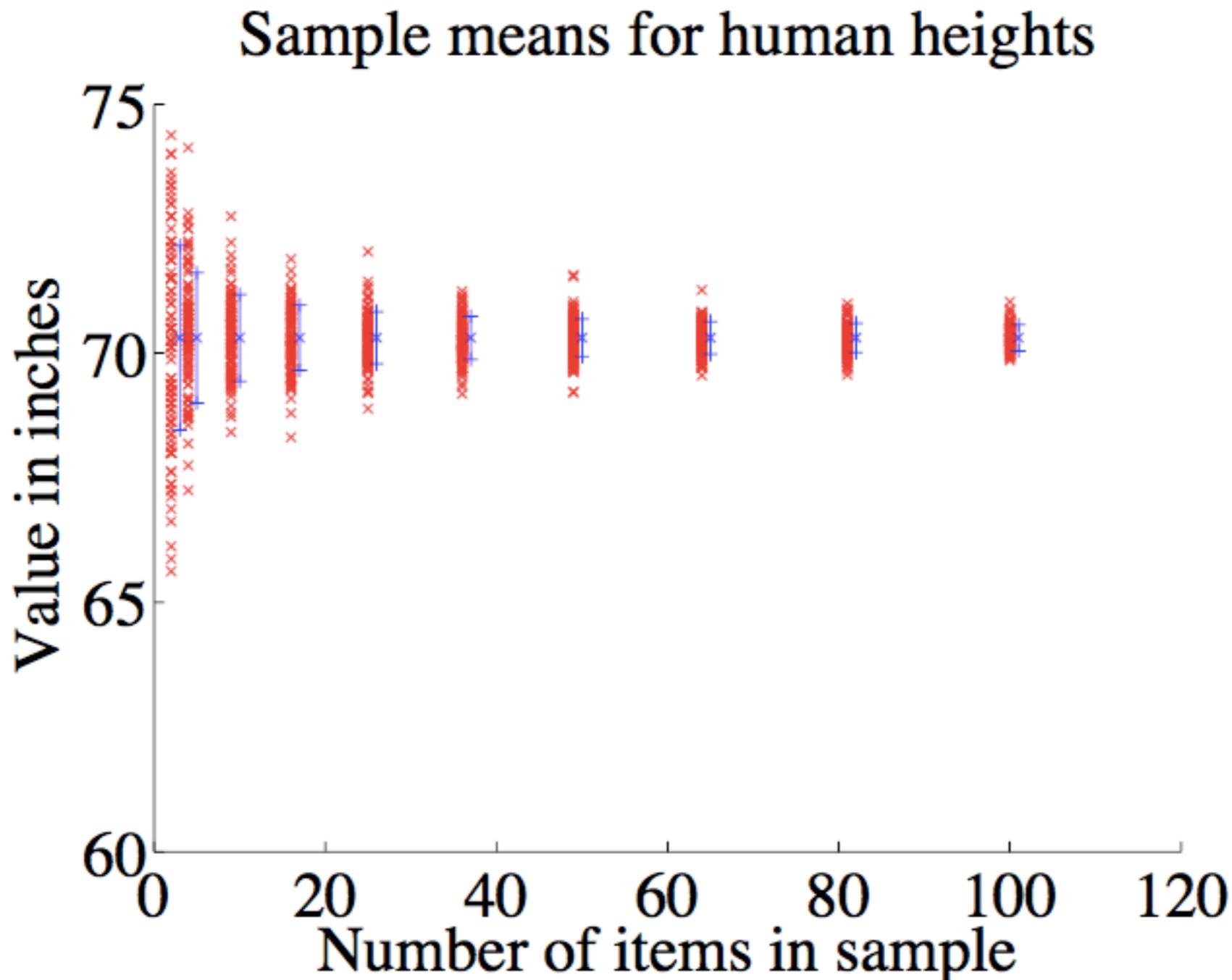
# Example

---

Suppose we have 100 heights in inches of students at UIUC and the sample mean is 66 inches and the sample variance is 16.

Our sample standard deviation is thus 4. So our standard error is  $4/\sqrt{100}$  or 0.4. Thus with 95% confidence the value for the mean of heights of all students at UIUC is between 65.2 and 66.8 inches

# Simulation



251 human heights in the population and various sample estimates of the mean

popmean and 1 popsd standard error bars in blue. Estimates from simulating the sampling process in red

---

# Bayesian confidence intervals

---

If we are doing parameter estimation we might want to evaluate the probability that the parameter lies within an interval and in many cases we just compute the following integral

$$\int_a^b p(\theta|\mathcal{D}) d\theta$$

---

# Example

---

We have a coin with unknown probability of coming up heads. Starting with a uniform Beta prior, suppose we flip the coin 10 times and observe 7 heads. What is the probability that theta is between 0.5 and 0.8?

Our posterior is given by

$$p(\theta|\mathcal{D}) = \frac{\Gamma(12)}{\Gamma(8)\Gamma(4)} \theta^7 (1 - \theta)^3$$

And we have

$$\int_{0.5}^{0.8} p(\theta|\mathcal{D}) d\theta \approx 0.73$$

---

# Example

---

We have a coin with unknown probability of coming up heads. Starting with a uniform Beta prior, suppose we flip the coin 10 times and observe 7 heads. Construct an interval  $[a,b]$  such that  $P(\theta \leq a|\mathcal{D}) \approx 0.05$  and  $P(\theta \geq b|\mathcal{D}) \approx 0.05$ . This is the 90% confidence interval for theta

We have to get some software to solve this for us, but we want to solve for a in

$$\int_{-\infty}^a P(\theta|\mathcal{D}) d\theta \approx 0.05 \quad \text{And b in} \quad \int_b^{\infty} P(\theta|\mathcal{D}) d\theta \approx 0.05$$

Doing so yields the interval  $[0.435, 0.865]$

---

# Bayesian confidence intervals

---

If the specified level of confidence we are looking for is  $1-2u$  for  $0 \leq u \leq 0.5$  we must find an  $a$  and  $b$  such that  $P(\{\theta \in [a, b]\} | \mathcal{D}) = 1 - 2u$

For any value on the RHS the  $a$  and  $b$  aren't unique. The reason we are writing it as  $1-2u$  instead of just  $p$  is that we would like to find an  $a$  and  $b$  such that

$$P(\{\theta \leq a\}) = \int_{-\infty}^a P(\theta | \mathcal{D}) d\theta = u$$

and

$$P(\{\theta \geq b\} | \mathcal{D}) = \int_b^{\infty} P(\theta | \mathcal{D}) d\theta = u.$$

Which are valid solutions since

$$P(\{\theta \in [a, b]\} | \mathcal{D}) = 1 - P(\{\theta \leq a\} | \mathcal{D}) - P(\{\theta \geq b\} | \mathcal{D})$$