

## Samples and the Bootstrap

Based on a chapter by Chris Piech

Let's say you are the king of Bhutan and you want to know the average happiness of the people in your country. You can't ask every single person, but you could ask a random subsample. In this next section we will consider principled claims that you can make based on a subsample. Assume we randomly sample 200 Bhutanese and ask them about their happiness. Our data looks like this: 72, 85, ..., 71. You can also think of it as a collection of  $n = 200$  I.I.D. (independent, identically distributed) random variables  $X_1, X_2, \dots, X_n$ .

### Estimating Mean and Variance from Samples

We assume that the data we look at are IID from the same underlying distribution ( $F$ ) with a true mean ( $\mu$ ) and a true variance ( $\sigma^2$ ). Since we can't talk to everyone in Bhutan, we have to rely on our sample to estimate the mean and variance. From our sample we can calculate a sample mean ( $\bar{X}$ ) and a sample variance ( $S^2$ ). These are the best guesses that we can make about the true mean and true variance.

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

The first question to ask is: are those unbiased estimates? Yes. Unbiased, means that if we were to repeat this sampling process many times, the expected value of our estimates should be equal to the true values we are trying to estimate. We will prove that that is the case for  $\bar{X}$ . The proof for  $S^2$  is in the lecture slides.

$$\begin{aligned} E[\bar{X}] &= E\left[\sum_{i=1}^n \frac{X_i}{n}\right] = \frac{1}{n} E\left[\sum_{i=1}^n X_i\right] \\ &= \frac{1}{n} \sum_{i=1}^n E[X_i] = \frac{1}{n} \sum_{i=1}^n \mu = \frac{1}{n} n\mu = \mu \end{aligned}$$

The equation for sample mean seems like a reasonable way to calculate the expectation of the underlying distribution. The same could be said about sample variance except for the surprising  $(n-1)$  in the denominator of the equation. Why  $(n-1)$ ? That denominator is necessary to make sure that the  $E[S^2] = \sigma^2$ .

The intuition behind the proof is that sample variance calculates the distance of each sample to the sample mean, *not* the true mean. The sample mean itself varies, and we can show that its variance is also related to the true variance.

## Standard Error

Ok, you convinced me that our estimates for mean and variance are not biased. But now I want to know how much my sample mean might vary relative to the true mean.

$$\begin{aligned} \text{Var}(\bar{X}) &= \text{Var}\left(\sum_{i=1}^n \frac{X_i}{n}\right) = \left(\frac{1}{n}\right)^2 \text{Var}\left(\sum_{i=1}^n X_i\right) \\ &= \left(\frac{1}{n}\right)^2 \sum_{i=1}^n \text{Var}(X_i) = \left(\frac{1}{n}\right)^2 \sum_{i=1}^n \sigma^2 = \left(\frac{1}{n}\right)^2 n\sigma^2 = \frac{\sigma^2}{n} \\ &\approx \frac{S^2}{n} && \text{Since } S \text{ is an unbiased estimate} \\ \text{Std}(\bar{X}) &\approx \sqrt{\frac{S^2}{n}} && \text{Since Std is the square root of Var} \end{aligned}$$

That  $\text{Std}(\bar{X})$  term has a special name. It is called the **standard error** and its how you report uncertainty of estimates of means in scientific papers (and how you get error bars). Great! Now we can compute all these wonderful statistics for the Bhutanese people. But wait! You never told me how to calculate the  $\text{Std}(S^2)$ . True, that is outside the scope of CS109. You can find it on Wikipedia if you want.

Let's say we calculate the our sample of happiness has  $n = 200$  people. The sample mean is  $\bar{X} = 83$  (what is the unit here? happiness score?) and the sample variance is  $S^2 = 450$ . We can now calculate the standard error of our estimate of the mean to be 1.5. When we report our results we will say that the average happiness score in Bhutan is  $83 \pm 1.5$  with variance 450.

## Bootstrap

Bootstrap is a newly invented statistical technique for both understanding distributions of statistics and for calculating  $p$ -values (a  $p$ -value is a the probability that a scientific claim is incorrect). It was invented here at Stanford in 1979 when mathematicians were just starting to understand how computers, and computer simulations, could be used to better understand probabilities.

The first key insight is that: if we had access to the underlying distribution ( $F$ ) then answering almost any question we might have as to how accurate our statistics are becomes straightforward. For example, in the previous section we gave a formula for how you could calculate the sample variance from a sample of size  $n$ . We know that in expectation our sample variance is equal to the true variance. But what if we want to know the probability that the true variance is within a certain range of the number we calculated? That question might sound dry, but it is critical to evaluating scientific claims! If you knew the underlying distribution,  $F$ , you could simply repeat the experiment of drawing a sample of size  $n$  from  $F$ , calculate the sample variance from our new sample and test what portion fell within a certain range.

The next insight behind bootstrapping is that the best estimate that we can get for  $F$  is from our sample itself! The simplest way to estimate  $F$  (and the one we will use in this class) is to assume that  $P(X = k)$  is simply the fraction of times that  $k$  showed up in the sample. Note that this defines the probability mass function of our estimate  $\hat{F}$  of  $F$ .

```
def bootstrap(sample):
    n = number of elements in sample
    pmf = estimate the underlying pmf from the sample
    stats = []
    repeat 10,000 times:
        resample = draw n new samples from the pmf
        stat = calculate your stat on the resample
        stats.append(stat)
    stats can now be used to estimate the distribution of the stat
```

To calculate  $\text{Var}(S^2)$  we could calculate  $S_i^2$  for each resample  $i$  and after 10,000 iterations, we could calculate the sample variance of all the  $S_i^2$ s.

The bootstrap has strong theoretic grantees, and is accepted by the scientific community, when calculating any statistic. It breaks down when the underlying distribution has a "long tail" or if the samples are not I.I.D.