

Supplemental 1

RATIOS

This lecture provides supplemental material concerning the handling of ratios in statistics. When you analyze, say with a t-test, a collection of ratios, the raw data are the ratios, and you are asking and answering questions about the individual ratios. When you are interested, however, in the population ratio, this is a mathematically incorrect way of asking and answering the question. We discuss a better way below.

Population Ratios

For the moment, let's consider traffic fatalities in 1995 and in 1996 in a country with two states. The data below have been stratified by the two states. Percentage increase, or decrease, in fatalities is computed as the difference in fatalities over the number of fatalities in 1995. This is the ratio minus one, so it is a quantity that is very closely related to the ratio. For starters, we are simply going to consider ratios for now.

	1995	1996	ratio
State 1	185	200	1.081
State 2	589	556	0.944
Totals	774	756	0.974

If you wanted to know the ratio for the country, however, you would not average the ratios 1.081 and 0.944 to report an average of 1.0125. Instead, you would total the fatalities for the two years and compute the ratio of the totals which gives 0.974. Even that statistic is slightly off. It is biased and, on average, is slightly away from where it should be. There are formulas, however, that, while complicated, are reasonable to compute.

Let

$$R = \frac{\mu_y}{\mu_x}$$

be a population ratio of interest. Let \bar{X} and \bar{Y} be the sample means, s_x^2 and s_y^2 be the sample variances, $\hat{\rho}$ be the sample correlation, and n be the sample size. A bias-corrected estimate of R is given by the following formula which uses all of the symbols just defined and can be calculated in Excel:

$$\hat{R} = \frac{\bar{Y}}{\bar{X}} - \left(\frac{1}{n}\right) \left(\frac{1}{(\bar{X})^2}\right) \left(\left(\frac{\bar{Y}}{\bar{X}}\right) s_x^2 - \hat{\rho} s_x s_y\right)$$

Its estimated variance is

$$\text{var}(\hat{R}) = \left(\frac{1}{n}\right) \left(\frac{1}{(\bar{X})^2}\right) \left(\left(\frac{\bar{Y}}{\bar{X}}\right)^2 s_X^2 + s_Y^2 - 2 \left(\frac{\bar{Y}}{\bar{X}}\right) \hat{\rho} s_X s_Y \right)$$

so that the standard error of the estimator \hat{R} is

$$S.E.(\hat{R}) = \sqrt{\text{var}(\hat{R})}$$

All of that may look complicated, but the formulas only involve quantities that are easy to calculate from the data.

Example: Did raising the speed limit increase fatalities?

Did the law allowing stated to raise the speed limit on American highways increase fatalities in car crashes? You did a paired analysis on these data treating the ratios in each state as independent data of interest. But what if we wanted an estimate of the percentage increase in fatalities in the collection of states that chose to increase the speed limit and we wanted to compare that estimate to the estimate from the collection of states that chose to retain the speed limit. Then we need the collection of estimates in the formula for estimating the collective ratio above. Let \mathbf{X} denote the 1995 fatalities and \mathbf{Y} denote the 1996 fatalities.

States that increased their speed limits	States that retained their speed limits
$\bar{X} = 941.5$ $\bar{Y} = 950.8$ $s_X^2 = 893,057$ $s_Y^2 = 934,842$ $n = 32$ $\hat{\rho} = 0.993$ $\hat{R} = 1.010$ S.E.(\hat{R}) = 0.021	$\bar{X} = 614.8$ $\bar{Y} = 612.6$ $s_X^2 = 218,486$ $s_Y^2 = 217,704$ $n = 19$ $\hat{\rho} = 0.997$ $\hat{R} = 0.996$ S.E.(\hat{R}) = 0.012

Based on these estimates and their standard errors, it does not seem like there is any difference between the states who chose to increase their speed limits and those that chose to retain them at 55 mph. We can do a formal, but approximate, test:

$$Z \approx \frac{\hat{R}_1 - \hat{R}_2}{\sqrt{(\text{S.E.}(\hat{R}_1))^2 + (\text{S.E.}(\hat{R}_2))^2}} = \frac{\hat{R}_1 - \hat{R}_2}{\sqrt{(\text{var}(\hat{R}_1)) + (\text{var}(\hat{R}_2))}}$$

is approximately normally distributed. It is only approximate, but we have no formal correction like when we use a t-distribution instead of a standard normal distribution, and so we simply use the normal approximation. Note that the variance of the sum or difference of two independent random variables is the sum of the variances - variances add, standard deviations and standard errors do not.

In our case,

$$Z \approx \frac{1.010 - 0.996}{\sqrt{(0.021)^2 + (0.012)^2}} = 0.55,$$

which is clearly an insignificant score.

The data follow. After that there is a technical explanation about why these formulas work. The technical explanation requires an understanding of calculus.

Increased speed limit				Retained speed limit			
State	1995 deaths	1996 deaths	percent change	State	1995 deaths	1996 deaths	percent change
AL	1114	1146	2.9	AK	87	81	-6.9
AZ	1025	994	-3	CT	317	310	-2.2
AR	631	615	-2.5	D.C.	58	62	6.9
CA	4192	3989	-4.8	HI	130	148	13.8
CO	645	617	-4.3	IN	960	984	2.5
DE	121	116	-4.1	KY	849	842	-0.8
FL	2805	2753	-1.9	LA	894	902	0.9
GA	1488	1573	5.7	ME	187	169	-9.6
ID	262	258	-1.5	MN	597	576	-3.5
IL	1586	1477	-6.9	NH	118	134	13.6
IS	527	465	-11.8	NJ	774	814	5.2
KS	442	490	10.9	NY	1679	1593	-5.1
MD	671	608	-9.4	NC	1448	1494	3.2
MA	444	417	-6.1	OR	574	526	-8.4
MI	1530	1505	-1.6	SC	881	930	5.6
MS	868	811	-6.6	VT	106	88	-17
MO	1109	1148	3.5	VA	900	877	-2.6
MT	215	200	-7	WV	377	348	-7.7
NE	254	293	15.4	WI	745	761	2.1
NV	313	348	11.2				
NM	485	485	0				
ND	74	85	14.9				
OH	1360	1391	2.3				
OK	669	772	15.4				
PA	1480	1469	-0.7				
RI	69	69	0				
SD	158	175	10.8				
TN	1259	1239	-1.6				
TX	3183	3742	17.6				
UT	325	321	-1.2				
WA	653	712	9				
WY	170	143	-15.9				

The Delta Method: Taylor Series Expansions and Approximation Methods

The reason why the formulas work is that we have expanded a function of one or more random variables in a Taylor series about the mean(s). Recall that the Taylor series expansion of a function of one variable about a specific point \mathbf{x}_0 is:

$$f(\mathbf{x}) = f(\mathbf{x}_0) + (\mathbf{x} - \mathbf{x}_0)f'(\mathbf{x}_0) + \frac{(\mathbf{x} - \mathbf{x}_0)^2}{2!}f''(\mathbf{x}_0) + \dots$$

For our purposes, we will only consider out to the linear or the quadratic terms. There is a more complicated, but similar, formula for a function of two variables. This section just describes functions of a single variable, so we will not derive the formula for ratios here.

If \mathbf{X} is our random variable and $\boldsymbol{\mu} = \mathbf{E}\mathbf{X}$ is its mean, then the first degree (linear) Taylor series approximation is

$$f(\mathbf{X}) \approx f(\boldsymbol{\mu}) + (\mathbf{X} - \boldsymbol{\mu})f'(\boldsymbol{\mu})$$

The mean is a good place to expand about because it is usually a good guess as to what a random data point will be and one generally expects a lot of data to be around the mean value. We use $\mathbf{E}f(\mathbf{X})$ to denote the mean of this function of the random variable \mathbf{X} . An initial approximation to this mean is

$$\mathbf{E}f(\mathbf{X}) \approx f(\boldsymbol{\mu}) + \mathbf{E}(\mathbf{X} - \boldsymbol{\mu})f'(\boldsymbol{\mu}) = f(\boldsymbol{\mu}) + \mathbf{0} = f(\boldsymbol{\mu})$$

since the expected value of a constant is just the constant and the expected value of \mathbf{X} is $\boldsymbol{\mu}$. Our approximation to the variance of $f(\mathbf{X})$ is

$$\text{var}(f(\mathbf{X})) \approx \text{var}(f(\boldsymbol{\mu}) + (\mathbf{X} - \boldsymbol{\mu})f'(\boldsymbol{\mu})) = \text{var}(\mathbf{X})(f'(\boldsymbol{\mu}))^2$$

A better approximation to the mean uses the second order (quadratic) Taylor series expansion of $f(\mathbf{X})$ and corrects for bias:

$$\mathbf{E}f(\mathbf{X}) \approx f(\boldsymbol{\mu}) + \mathbf{E}(\mathbf{X} - \boldsymbol{\mu})f'(\boldsymbol{\mu}) + \frac{\mathbf{E}(\mathbf{X} - \boldsymbol{\mu})^2}{2}f''(\boldsymbol{\mu}) = f(\boldsymbol{\mu}) + \mathbf{0} + f''(\boldsymbol{\mu})\text{var}(\mathbf{X})/2$$

Note that all of those formulas involve things that we do not know, such as the mean and the variance of the population. However, we use the "plug" in principle in statistics that says we should plug in our best guesses, namely the mean and variance of the data. That gives the following formulas:

$$\mathbf{E}f(\mathbf{X}) \approx f(\bar{X}) + f''(\bar{X})s_X^2/2$$

and

$$(\text{var})(f(\mathbf{X})) \approx s_X^2(f'(\bar{X}))^2$$

Example 1 An exponential random variable, \mathbf{X} , with mean $\mathbf{1}$ has variance $\mathbf{1}$. Use this information to approximate the mean and variance of $\ln(\mathbf{X})$.

Let $f(x) = \ln(x)$. Then $f'(x) = 1/x$ and $f''(x) = -1/x^2$. So that $E \ln(\mathbf{X}) \approx \ln(\mathbf{1}) - (\mathbf{1}/\mathbf{2}) = -\mathbf{1}/\mathbf{2}$ and $(var)(\ln(\mathbf{X})) \approx \mathbf{1}$. These are not great approximations because the exponential distribution and the log of it are both very spread out. Simulations suggest that the mean is about $-\mathbf{0.57}$ and the variance is about $\mathbf{1.7}$. However, if instead of a single member of the population, we were dealing with population averages, the approximation would improve. See the next example. ■

Example 2 Suppose we have a sample of size 100 from any distribution and the sample mean is $\mathbf{1}$ and the sample variance is $\mathbf{1}$. Approximate the mean and variance of $\ln(\bar{\mathbf{X}})$. The difference here is that the sample average, $\bar{\mathbf{X}}$ is approximately normal and is close to its mean.

Note that $\ln(\mathbf{1}) = \mathbf{0}$ but the $var(\bar{\mathbf{X}}) = \mathbf{1}/\mathbf{100}$ since the variance of the mean decreases proportional to the sample size. Thus, $E \ln(\bar{\mathbf{X}}) \approx \ln(\mathbf{1}) - (\mathbf{1}/\mathbf{200}) = -\mathbf{0.005}$ and $(var)(\ln(\bar{\mathbf{X}})) \approx \mathbf{1}/\mathbf{100} = \mathbf{0.01}$. Simulations suggest that these approximations are quite accurate. ■

For a derivation of the formula for ratios, see Rice: *Mathematical Statistics and Data Analysis*.

REFERENCES AND READINGS

- [1] John Rice. *Mathematical Statistics and Data Analysis*. Brooks/Cole, 3rd edition, 2006.

Exercises for Lecture 1

1. –

2. –