

22S:105
Statistical Methods and Computing

Two independent-sample t-tests

Lecture 17
 Mar. 20, 2006

Kate Cowles
 374 SH, 335-0727
 kcowles@stat.uiowa.edu

Two independent sample problems

- Goal of inference:
 - to compare the characteristics of two different populations
 - to compare responses to two different "treatments"

- Examples of two-independent-sample problems:
 - A medical researcher is interested in the effect on blood pressure of added dietary calcium. She conducts a randomized comparative experiment in which one group of subjects receives a calcium supplement and a control group gets a placebo.
 - A climatologist wishes to test whether seeding with silver nitrate affects the amount of rainfall produced from clouds. He randomly selects 26 clouds to seed, and measures the rain output by each of them as well as the rain output by 26 other randomly selected unseeded clouds.

Which study design?

The following situations require inference about a population mean or means. Identify which type of problem each one is:

- one-sample
- paired-sample
- two independent samples

1. To check a new method of chemical analysis, a chemist gets a reference specimen of known concentration from the National Institute of Standards and Technology. She then makes 20 measurements of the concentration of this specimen using the new method and checks for bias by comparing the mean of her 20 measurements with the known concentration.
2. Another chemist is checking the same new method. He has no reference specimen, but a familiar analytic method is available. He wants to know if the new and old methods agree. He takes a specimen of unknown concentration and measures the concentration 10 times with the new method and 10 times with the old method.

The set-up for two-independent-sample t-tests

		Group 1	Group 2
Population	Mean	μ_1	μ_2
	Standard deviation	σ_1	σ_2
Sample	Mean	\bar{x}_1	\bar{x}_2
	Standard deviation	s_1	s_2
	Sample size	n_1	n_2

Comparing means from two different populations

Assumptions:

- We have two simple random samples, from two distinct populations.
 - The samples are independent.
 - * The selection of one sample has no influence on the selection of the other. In particular, there is no matching.
 - The sizes of the two samples need not be the same.
 - We measure the same variable for both samples.
- The populations are normally distributed.

Example: Cloud seeding

We wish to use our sample data to test whether the population mean of rainfall produced per cloud is the same for unseed clouds as for seeded clouds. We will use a two-sided test assuming that we don't know in advance in what direction a difference is likely to go.

$$H_0 : \mu_u = \mu_s \text{ or } \mu_u - \mu_s = 0$$

$$H_a : \mu_u \neq \mu_s \text{ or } \mu_u - \mu_s \neq 0$$

We will conduct our test at $\alpha = .05$.

Thus the quantity we really want to estimate is the difference between the two population means

$$\mu_u - \mu_s$$

As usual, we will use the corresponding sample statistics

$$\bar{x}_u - \bar{x}_s$$

as our best guess of the unknown population value of interest.

Now we need to standardize $\bar{x}_u - \bar{x}_s$ in order to find out whether it is different enough from 0 to provide strong evidence against H_0 .

That is, we need to compute

$$\frac{(\bar{x}_u - \bar{x}_s) - (\mu_{u0} - \mu_{s0})}{\text{standard error of } (\bar{x}_u - \bar{x}_s)}$$

Suppose:

- we knew the standard deviations σ_u and σ_s in the populations of rainfall amounts from unseeded and seeded clouds
-

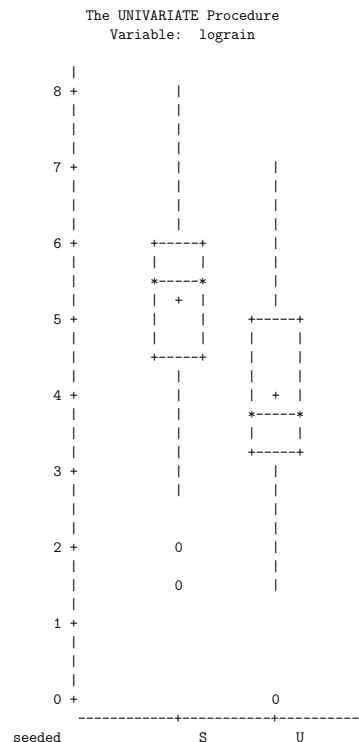
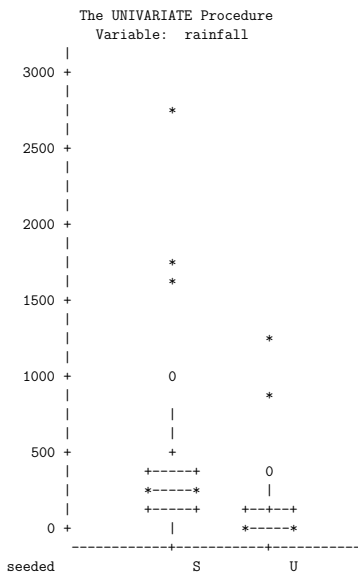
$$\sigma_u = \sigma_s = \sigma \text{ (some known value)}$$

Then the standard error of $(\bar{x}_s - \bar{x}_u)$ would be

$$\sqrt{\frac{\sigma^2}{n_u} + \frac{\sigma^2}{n_s}}$$

And the z statistic is

$$z = \frac{(\bar{x}_u - \bar{x}_s) - (\mu_{u0} - \mu_{s0})}{\sqrt{\frac{\sigma^2}{n_u} + \frac{\sigma^2}{n_s}}}$$



Example:

We will log-transform the rainfall amounts to get more symmetrical distributions of sample values. Suppose we knew that the population standard deviation of log-transformed rainfall amount was 1.5 log acre-feet for both seeded and unseeded clouds.

$$\sigma_u = \sigma_s = \sigma = 1.5$$

ulation mean of log rainfall is different in seeded clouds from unseeded clouds.

The results of measuring and logtransforming the rainfall are as follows:

Variable: LOGRAIN

SEEDED?	N	Mean	Std Dev	Std Error
S	26	5.13418678	1.59951361	0.31369043
U	26	3.99040563	1.64184748	0.32199278

The z statistic is

$$\begin{aligned} z &= \frac{3.99 - 5.13 - 0}{\sqrt{\frac{1.5^2}{26} + \frac{1.5^2}{26}}} \\ &= \frac{-1.14}{0.416} \\ &= -2.74 \end{aligned}$$

The cutoff values of z for a two-sided hypothesis test are -1.96 and 1.96. Because -2.74 is farther from 0 than either of these cutoffs, we can reject the null hypothesis and conclude that the pop-

If we do *not* know σ in the two populations, we need to use the data to estimate the standard error of $(\bar{x}_u - \bar{x}_s)$.

This can be done under two different assumptions:

1. The standard deviations in the two populations are known to be equal.
2. The standard deviations in the two populations are *not* known or assumed to be equal.

Two-sample t-test when variances are assumed to be equal

We must estimate the common variance σ^2 using the *pooled variance* s_p^2 from the two samples:

$$s_p^2 = \frac{(n_u - 1)s_u^2 + (n_s - 1)s_s^2}{n_u + n_s - 2}$$

Then we compute the t-statistic by substituting s_p^2 for σ^2 in the formula for the z-statistic.

$$t = \frac{(\bar{x}_u - \bar{x}_s) - (\mu_{u0} - \mu_{s0})}{\sqrt{\frac{s_p^2}{n_u} + \frac{s_p^2}{n_s}}}$$

When we cannot assume that the variances in the two populations are equal

When we do not assume that the standard deviations in the two populations are equal; i.e. when we think

$$\sigma_m \neq \sigma_f$$

then we estimate

- σ_u^2 with s_u^2
- σ_s^2 with s_s^2

For the cloud-seeding example,

$$\begin{aligned} s_p^2 &= 1.62 \\ t &= \frac{3.99 - 5.13 - 0}{\sqrt{\frac{1.62^2}{25} + \frac{1.62^2}{25}}} \\ &= \frac{-1.144}{0.458} \\ &= -2.5 \end{aligned}$$

We would compare this to the .025 cutoff for a t distribution with $n_u + n_s - 2 = 50$ degrees of freedom.

According to Table C, this cutoff would be 2.009.

Because the value we obtained is farther from 0 than this, we can reject the null hypothesis.

Then we compute the t statistic as

$$t = \frac{(\bar{x}_u - \bar{x}_s) - (\mu_{u0} - \mu_{s0})}{\sqrt{\frac{s_u^2}{n_u} + \frac{s_s^2}{n_s}}}$$

For the cloud seeding example, this gives

$$\begin{aligned} t &= \frac{3.99 - 5.13 - 0}{\sqrt{\frac{1.64^2}{25} + \frac{1.60^2}{25}}} \\ &= \frac{-1.144}{0.458} \\ &= -2.54 \end{aligned}$$

This form of t statistic does *not* come from an exact t distribution. Statistical software uses an approximation to compute the p-value in this case. Generally, the p-value is very close to that obtained under the assumption of equality of variance.

SAS for two-independent-sample t tests

```
data cloud ;
infile 'cloud.dat' ;
input rainfall seeded $ ;
lograin = log(rainfall) ;
run ;
```

```
proc ttest ;
class seeded ;
var lograin ;
run ;
```

TTEST PROCEDURE Statistics

Variable	seeded	N	Lower CL		Upper CL	
			Mean	Mean	Mean	Mean
lograin	S	26	4.4881	5.1342	5.7802	
lograin	U	26	3.3272	3.9904	4.6536	
lograin	Diff (1-2)		0.2409	1.1438	2.0467	

Statistics

Variable	seeded	Lower CL		Upper CL	
		Std Dev	Std Dev	Std Dev	Std Err
lograin	S	1.2544	1.5995	2.208	0.3137
lograin	U	1.2876	1.6418	2.2664	0.322
lograin	Diff (1-2)	1.3562	1.6208	2.0148	0.4495

Statistics

Variable	seeded	Minimum	Maximum
lograin	S	1.411	7.9178
lograin	U	0	7.0922
lograin	Diff (1-2)		

T-Tests

Variable	Method	Variances	DF	t Value	Pr > t
lograin	Pooled	Equal	50	2.54	0.0141
lograin	Satterthwaite	Unequal	50	2.54	0.0141

Equality of Variances

Variable	Method	Num DF	Den DF	F Value	Pr > F
lograin	Folded F	25	25	1.05	0.8971