

2008-6

a) Although these are differences, we are comparing the means of two populations of these differences, so we will use a 2 sample t-test for the difference of the means:

Define μ_{DIFFM} = population mean of the post-pre test differences in score for students who attended the magnet school.

μ_{DIFFO} = population mean of the post-pre test difference in score for students who attended their original school.

$$H_0: \mu_{\text{DIFFM}} = \mu_{\text{DIFFO}} \text{ (or } \mu_{\text{DIFFM}} - \mu_{\text{DIFFO}} = 0)$$

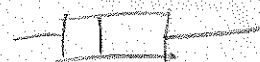
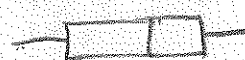
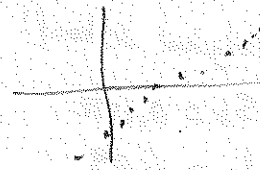
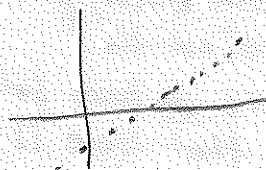
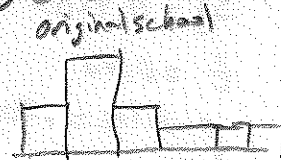
$$H_A: \mu_{\text{DIFFM}} > \mu_{\text{DIFFO}} \text{ (or } \mu_{\text{DIFFM}} - \mu_{\text{DIFFO}} > 0)$$

CONDITIONS:

- Groups indep: There is no reason to believe that the two groups of students do not vary independently of each other.

- SRS: The problem states that the students in each group were randomly selected.

- Nearly Normal: Checking histograms of distributions of the differences:



Some possible skew in original school differences. Investigate further w/ Normal Prob Plot & Boxplot

Skew is not severe, and no outliers so both distributions are Nearly Normal.

- n < 100 pop: The sample sizes of 8 and 13 are less than 100 of the respective populations.

On my TI-84 I performed a 2 Sample TTest using the sample differences for list data, and $\mu_1 > \mu_2$ with no pooling. The result was

$$t = 2.487$$

$$p\text{-value} = .0178$$

$$df = 8.689$$

with $\alpha = .05$, p-value of .0178 is low so we reject H_0 . There is sufficient evidence to conclude that the mean pre-to-post improvement in scores is higher for students at the magnet school compared to students at their original school.

$$b)(i) \hat{\text{post score}}_{\text{magnet}} = 73.27 + 0.1811 (\text{pre score}_{\text{magnet}})$$

The slope of 0.1811 indicates that for every 1 point increase in score on the pretest there will be, on average, an increase of 0.1811 points in score on the posttest (for students at the magnet school).

$$(ii) \hat{\text{post score}}_{\text{original}} = 9.24 + 0.9204 (\text{pre score}_{\text{original}})$$

The slope of 0.9204 indicates that for every 1 point increase in score on the pretest there will be, on average, an increase of 0.9204 points in score on the posttest (for students at the original school.)

c) For inference for regression (slope of Least Squares Regression Line, LSRL) the t -statistic and standard error for the slope coefficient are given in the software output, along with corresponding p -values (and the problem states conditions for inference have been met):

(i) For the magnet school: with $\alpha = .05$, $p\text{-value} = 0.706$ is high so we fail to reject H_0 .
 $S_b = 0.4583$
 $t = 0.40$
 $p\text{-value} = .706$
 We do not have evidence of an association (correlation) between pre and post test scores for the magnet school.

(ii) For the original school: with $\alpha = .05$, $p\text{-value} = 0.000$ is low so we reject H_0 .
 $S_b = 0.1572$
 $t = 6.09$
 $p\text{-value} = 0.000$
 We do have sufficient evidence to conclude that there is an association (correlation) between pre and post test for the original school.

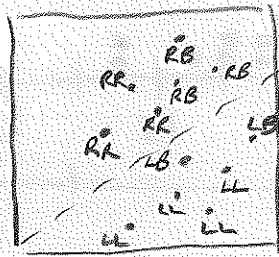
d) The slope of the LSRL for the magnet school is close to zero, meaning that all the students performed well on the posttest, regardless of how they performed on the pretest.

For the original school, the LSRL slope is close to one, meaning students who performed well on the posttest were mainly students who also scored well on the pretest.

It appears that the magnet school improves the performance of low pre-test score students more compared to the original school. At the original school, only students who scored high on the pretest scored high on the post test.

2008b#6

(a) Labelling each point on scatterplot w/ dominant foot & non-dominant foot:



The upper cluster represents patients who are right foot dominant.
The lower cluster represents patients who are left foot dominant.

(b) There is a positive, fairly strong, linear relationship between amount of dominant and non-dominant foot swelling. Also, swelling seems to be greater in the dominant foot.

(c) Because these are matched pairs so we should conduct a paired-t test on the mean of the differences.

Define μ_d = mean of the differences in swelling (dominant - non-dominant).

Then $H_0: \mu_d = 0$

$H_a: \mu_d \neq 0$

Conditions: paired data ✓
pairs indep. of each other. This is a random sample ✓
Sample representative of population? Yes ✓
Differences Nearly Normal? ✓

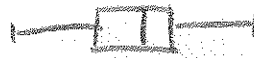
Histogram of
diffs.



Normal probability plot



Boxplot



Visualizations of the differences show an approximately Normal distribution with no obvious outliers.

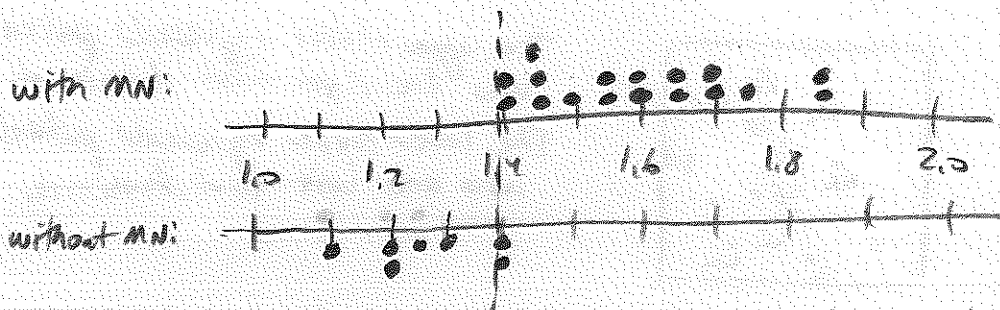
On TI-84 conduct a T-Test using the differences from the table, and $\mu_0 = 0$

$$t = 12.6817, \quad p\text{-value} = 3.8140 \times 10^{-7}$$

with $\alpha = .05$, p -value of 3.8140×10^{-7} is low so we reject H_0 .

There is strong evidence of a difference between the mean swelling in dominant and non-dominant foot (in this population).

(d) Divide the data into two groups: w/MN & w/o MN and plot the swelling values, \bar{x} and s .



These back-to-back dotplots suggest that a value of 1.4 would be a good criterion value. For values of 1.4 or higher, the patient is diagnosed with Morton's neuroma.

2010b #6

- (a) The slope of 0.165 means that for every additional 1 square foot in size the price of the house increases, on average, \$165.
- (b) This positive residual means that this house's actual price is \$49,000 higher than the price the LSRL model would predict for a home of this size.
- (c) "Use the residuals". We could find the average residual for each group:

$$\text{w/pool: } \frac{6 + 49 - 18 + 42 + 1 + 50 - 23 + 42}{8} = 18,625$$

$$\text{w/o pool: } \frac{13 + 26 - 45 + 22 + 10 - 46 - 57 + 1 - 2 - 69 + 23 + 44 - 19 + 26 - 58 - 5247}{17} = -8,824$$

The model underestimates pool houses by \$18,625 and overestimates non pool houses by \$8,824. So, on average, a pool increases the price by \$27,449.

(d) No. The value difference of 0 is contained in the interval, which is equivalent to failing to reject H_0 of $\mu_d = 0$ with $\alpha = 0.05$.

(e) The difference in predicted price depends upon the size of the house (because the slopes are not equal). But we can compare calculated values at particular sizes:

$$\text{at low size end (1500 sqft): w/pool price} = -11,602 + 0.166(1500) = 237,398$$

$$\text{w/o pool price} = -27,382 + .160(1500) = 212,618$$

$$\text{price diff} = \$24,780$$

$$\text{at mid size (2250 sqft): w/pool price} = -11,602 + .166(2250) = 361,893$$

$$\text{w/o pool price} = -27,382 + .160(2250) = 332,618$$

$$\text{price diff} = \$29,275$$

$$\text{at high size end (3000 sqft): w/pool price} = -11,602 + .166(3000) = 486,398$$

$$\text{w/o pool price} = -27,382 + .160(3000) = 452,618$$

$$\text{price diff} = \$33,780$$

The price differences between the models is similar to the residuals calculated price differences.

2009b #6

(a) Yes. Both groups need to receive both a pill and an injection. For group 1 pill is real and injection is placebo. For group 2 injection is real and pill is placebo. To be double-blind not only the subjects, but those administering the treatments must not know which is which, so the researchers must prepare the pills and injections so that those interacting with patients cannot identify real vs placebo pills or injections.

(b) This would be a difference of 2 proportions Z Interval
Let P_A = population proportion receiving pill^A who survive at least 15 yrs.
 P_B = population proportion receiving injection^B who survive at least 15 yrs.

$$\hat{p}_A = \frac{38}{154} = 0.2468 \quad \hat{p}_B = \frac{16}{164} = 0.0976$$

Using my TI-84 I conducted a Z-proportion Int

with $x_1: 38$

$n_1: 154$

$x_2: 16$

$n_2: 164$

and C-level: 0.95

The resulting confidence interval is:

$$(0.06735, 0.23104)$$

I am 95% confident that the difference in percentage of all patients surviving at least 15 yrs. (A-B) is between 6.73% and 23.10%.

This suggests that treatment A has a higher 15-yr survival rate.

$$(c) \frac{\hat{p}_A}{\hat{p}_B} = \frac{0.2468}{0.0976} = \boxed{2.53}$$

(d) 95% CI for $\ln\left(\frac{p_A}{p_B}\right)$ is (0.3868, 1.4690)

endpoint values:

$$\text{for } 0.3868: e^{0.3868} = 1.4723$$

$$\text{for } 1.4690: e^{1.4690} = 4.3449$$

The procedure given in the problem suggests that the corresponding 95% confidence interval for $\frac{p_A}{p_B}$ is

$$(1.4723, 4.3449)$$

This means we are 95% confident that people are between 1.47 and 4.34 times more likely to survive 15 years or more with treatment A than with treatment B.

(e) The values are larger (survival rates) so higher numbers have more practical meaning ~~for~~ to people.

2004 #6

- (a) If μ = mean reduction in (population) level of cholesterol, we are given $\mu_x = 24$ mg/dl $\sigma_x = 15$ mg/dl. We would build a confidence interval using a t -statistic.

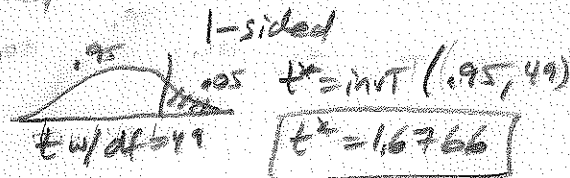
Conditions: Sample representative of population: we are told random selection ✓
sampling distn approx normal: can assume because $n > 25$. ✓
Samples independent: likely due to random sample. ✓

On TI.84 TInterval w/ $\bar{x} = 24$ $(19.737, 28.263)$
 $\sigma_x = 15$
 $n = 50$
 $C-level = .95$

We are 95% confident that the mean reduction in cholesterol for the new drug in the population is between 19.74 and 28.26 mg/dl.

- (b) The confidence interval is a two-sided test corresponding to a significance level of $\alpha = .05$. The new test is one-sided. If it were two-sided the p -value of .033 would be doubled to .066 which would no longer be significant.

- (c) For a 95% confidence interval



$$L = \bar{x} - t^* \frac{s}{\sqrt{n}}$$

$$L = 24 - 1.6766 \frac{15}{\sqrt{50}}$$

$$L = \boxed{20.4447}$$

- (d) Yes. 20 mg/dl (the decision point) would now be below the lower bound (L) of the new confidence interval.

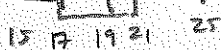
2013#6

(a) Enter data to display boxplots:

Western Pacific



Eastern Pacific



Western numbers of typhoons are higher for all but one year.

The median for Western, 31, is higher than all years of Eastern storms, 25.

3/4 of the Western years have higher number of storms than any Eastern year.

There is much more variability in Western year number of storms.

(b) Western storm numbers show a decreasing trend,
Eastern storm numbers are relatively consistent over time.

(c) Western data for 4 years ending in 2010: 28, 27, 28, 18

$$\text{Moving avg}_{\text{West}, 2010} = \frac{28+27+28+18}{4} = \boxed{25.25 \text{ storms}}$$

(d) (add a dot at 2010, 25, 25)

(e) (i) The moving average more clearly shows the trend over time by reducing the yearly variation. This reinforces the earlier finding that the number of storms is generally decreasing in the Western Pacific, but is fairly constant (with a slight upward trend from 2005 to 2010) in the Eastern Pacific.

(ii) The moving averages hide year-to-year variability. In the original yearly plots, there is more variability from year-to-year in the Western Pacific than the Eastern (at least for the period from 2002 to 2006.)