

Solutions version 1

Midterm 3 - STAT 301 Section 501/502 version 1 Spring 2012

Name:

UIN:

Signature:

1. Do not open this test until told to do so.
2. Turn in your exam with your answers circled when you are done with the exam. You should not take the exam with you.
3. This is a closed book examination. You may use three cheat sheets, that you have brought with you, and the tables. You should have no other printed or written material with you on the exam.
4. You have 50 minutes to work on this exam. There are 6 questions.
5. You may use a calculator but not a phone during the exam.
6. If you are unsure of what a question is asking for, do not hesitate to ask the instructor or teaching assistant for clarification.
7. Good Luck!!!

(1) Your aim is to design an experiment and collect data and construct a 95% confidence interval for the mean. You know that the standard deviation of a population is $\sigma = 3$.

(i) Your colleagues ask you to briefly explain what the confidence interval tells you about the mean. Explain this in one sentence. [1]

The confidence interval is a range of values where you are likely (95% sure) to find the unknown population mean.

(ii) What sample size should you use such that the margin of error is at most 0.5? [1]

Solve the equation

$$0.5 = \frac{1.96 \times 3}{\sqrt{n}}$$

$$n = \left(\frac{1.96 \times 3}{0.5} \right)^2 = 138.3 \text{ need a sample size of at least } 139.$$

(iii) Suppose you were incorrectly told that the standard deviation was $\sigma = 3$, and in fact the standard deviation of the population is $\sigma = 5$. What effect would this have on the margin of error of the CI for the sample size calculated in (ii). [1]

The margin of error would be larger than 0.5
→ larger the s.d. the wider the CI in order to accommodate the mean.

(2) Fred wants to use statistics to 'prove' his research. As he learnt in his statistics class, he stated his research hypothesis as the alternative and the opposite of that as his null hypothesis. After collecting the data, the test at the 5% level and is unable to reject the null. Disappointed with this result, he goes and collects another sample does the test at the 5% level and again is unable to reject the null. He does this 5 times and on the fifth attempt rejects the null. He then writes a paper saying that there is evidence to reject the null hypothesis, because his p-value was less than 5%.

Are you convinced by his results? Comment on the validity of the result? [2]

No! By using a 5% significance level, if the null is true then 5 times out of a hundred you will falsely reject the null. Therefore, if the null is true and you do the experiment several times at some point you will falsely reject the null. Therefore repeating the experiment until you reject, then claiming one rejection proves the result is ridiculous.

(3) During your statistics class you heard a lot about t-distributions with this and that degrees of freedom. The main reason we use a t-distribution instead of a normal distribution is to account for the extra variability in estimating the standard deviation rather than using the true standard deviation. Suppose the sample size is very large (say over 100), in general

does it matter whether you use a normal distribution or t-distribution, give a reason for your answer? [2]

No it doesn't. As the sample size grows the estimate of the standard ~~error~~ deviation gets better. This means that as the number of degrees of freedom grow the t-distribution gets closer to the normal distribution.

- (4) A tomato packing company advertises that the mean weight of their tomato boxes is 300g (it is also known that standard deviation of the weights is 10g). It is okay if the mean weight is greater than what they claim but it is an offense if the mean is less than what is claimed (and risks closing the company down). Therefore, every year an inspector takes a sample of 30 boxes of tomatoes to check their claims.

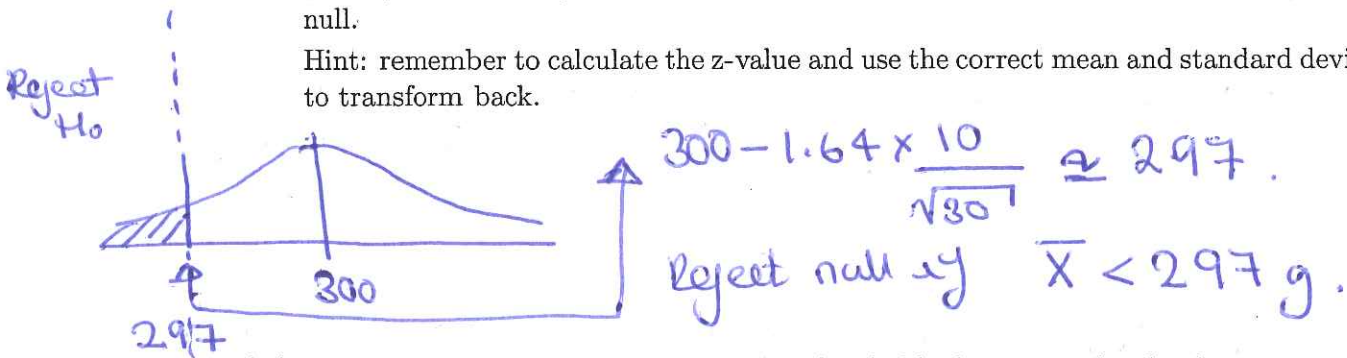
(Assume in this question that the standard deviation is known, therefore there is no need to use anything but the normal distribution).

- (i) State the null and the alternative for this scenario. [2]

$$H_0: \mu \geq 300g \quad H_1: \mu < 300g$$

- (ii) The inspector does the test at the 5% level. Based on this information, give the cut off point (critical region) such that if the average weight is less than this value we reject the null. [2]

Hint: remember to calculate the z-value and use the correct mean and standard deviation to transform back.

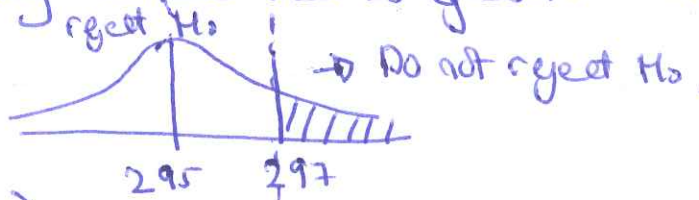


- (iii) If the sample average is less than this threshold, then we make the decision to accept the alternative. Will this decision be right or wrong? Explain your answer. [1]

Either (a) the true mean $\mu < 300g$ or (b) we have falsely rejected null. We can never be sure whether we are right or wrong.

- (iv) What is the probability that the sample mean will be greater than the threshold given in part (ii) if the true mean packing weight of tomato boxes is 300g? If you do not have a number for part (ii), then use the number 290 instead. [1]

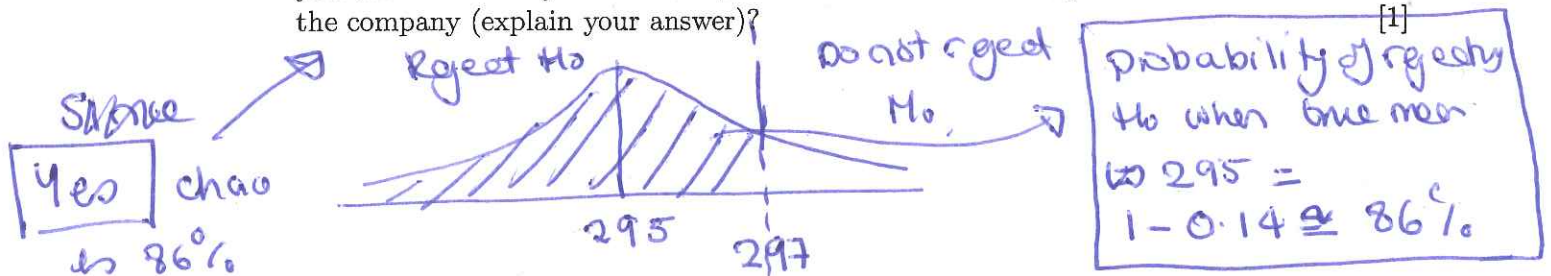
calculate probability sample mean is greater than 297, if $\mu = 295$.



$$= P\left(Z > \frac{297 - 295}{\frac{10}{\sqrt{30}}}\right) \approx 0.14$$

Standard error

- (v) Suppose that an employee has a grudge against the company. He recalibrates the machine such that its mean weight of tomatoes is now $\mu = 295g$. Using your answer in (ii), do you think it is likely that the inspector will detect the change in mean and close down the company (explain your answer)? [1]



- (5) A new feed has been developed that is believed to cause a greater weight gain in calves weight over conventional feeds. To statistically test these claims a farmer random sampled 20 sets of twins (calves). To one twin he gave conventional feed and the other calf he gave the new feed. The farmer did not know what test to do so he did both a paired t-test and an independent sample t-test. The statcrunch outputs of both tests are summarised below (X corresponds to the weight on conventional feed and Y to the weight on the new feed).

- (i) By selecting the correct output (see Figures 1 and 2), do we reject the test at the 5% level? [1]

Paired t-test. $p\text{-value} = 0.7\% < 5\%$ reject null.

- (ii) Using the output construct a 95% CI for the difference in the population means. [1]

$$[-0.75 \pm t_{19}(0.025) \times 0.28] = [-0.75 \pm 2.09 \times 0.28] = [-1.32, -0.164]$$

- (iii) The manufacturer of the new feed claims that the mean weight increase over the conventional feed is at least 1.5 pounds, is there any evidence to substantiate these claims (evidence here being at the 5% level)? [1]

Hint: Part (ii) can help you answer this question.

Since -1.5 ^{or less} does not lie in the above interval, there is no evidence ^{to support claim} that the mean weight gain is greater than 1.5 pounds.

- (iv) Using the plot in Figure 3 explain why there is a difference in the results of the tests for the paired t-test and the independent sample t-test. [1]

There is a large amount of variability within the weights of the calves. This means the standard deviation in the independent sample t-test is large. This 'masks' any difference there may be between the feeds. However the paired t-test (where we consider the differences) removes this variability. Therefore we can 'uncover' ~~the~~ or clearly see the difference.

Hypothesis test results:

μ_1 : mean of X

μ_2 : mean of Y

$\mu_1 - \mu_2$: mean difference

$H_0 : \mu_1 - \mu_2 = 0$

$H_A : \mu_1 - \mu_2 < 0$

(without pooled variances)

large standard error

Difference	Sample Mean	Std. Err.	DF	T-Stat	P-value
$\mu_1 - \mu_2$	-0.74636257	0.54786366	28.843489	-1.3623145	0.0918

Figure 1: Independent sample t-test

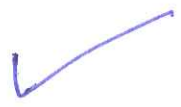
Hypothesis test results:

$\mu_1 - \mu_2$: mean of the paired difference between X and Y

$H_0 : \mu_1 - \mu_2 = 0$

$H_A : \mu_1 - \mu_2 < 0$

small standard error



Difference	Sample Diff.	Std. Err.	DF	T-Stat	P-value
X - Y	-0.74636257	0.2778224	19	-2.6864736	0.0073

Figure 2: Paired t-test

large amount of variability in the weights

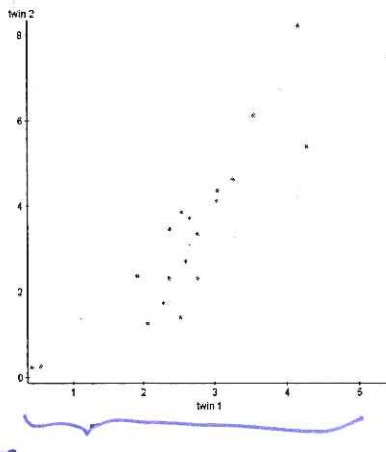


Figure 3:

- (6) A survey was conducted to see whether there is a difference in the preference of fruit juices of those living on the east and west coast. 1000 people were surveyed and the 2 by 2 table summarising the data is given below.

	East	West	Subtotal
Orange	100	300	400
Apple	500	100	600
Subtotal	600	400	1000

- (i) Fill in the missing cells in the table above.

[1]

- (ii) If there is no association between the location of a person and their fruit juice preference, on average how many of the West coast people in the survey would you expect to like orange juice and how many would you expect to like Apple juice? [1]

No associations \rightarrow proportions stay the same
 Expected number to like orange juice on the west coast = $\frac{400}{1000} \times 400 = 160$
 " " apple juice on the west coast = $\frac{600}{1000} \times 400 = 240$

- (iii) Without doing a formal test, by using the table and your answer to (ii), do you believe there is an association, give a reason for your answer? [1]

comparing what we expect to see [160 and 240] with what we do observe [300, 100] the differences between 160 and 300 and 240 and 100 are rather large. This suggests there is an ~~assoc~~ association between drink preference and coast.
 To test this statistically we use a χ^2 -test for association.