

CS306: Data Analysis and Visualization

Lab 7: Report

Rajdeep Pinge 201401103

Aditya Joglekar 201401086

Objective:

To study how t-test and confidence intervals can be used to compare data, and make informed decisions about the given problem.

Experiment 1:

Data analysis techniques are widely used in biological sciences. In this experiment, your goal is to analyze beak sizes of the same bird species that are found on two different islands. The files ‘island1.txt’ and ‘island2.txt’ provide the beak sizes from two samples of the bird species.

1. The goal is to compare the average beak sizes of birds found on the two islands, from a statistical viewpoint. You should use CIs for your analysis. Based on this, which of the two islands has birds with bigger beak sizes?

It is important to realize that we have been given samples of the beak sizes from the two islands. These samples may or may not represent the population. Thus, it would be naive to simply calculate the average bird size of the two samples. This would only be a comparison between the samples. What we are interested in is “inferring” about the population.

We realize that this is where ‘inferential statistics’ comes in. In the last lab we had used the t-test to determine if the sample statistic can be used to infer about the population statistic.

In this lab, we use the idea of Confidence Interval instead. In t-test, we try to estimate what is the likelihood of getting an extreme value of the test statistic by chance assuming that the null

hypothesis is True. This gives us an idea of how representative our sample statistic is, with respect to the population.

Here, when we define a Confidence Interval, we do it with a certain confidence level. For e.g when we say that I am 95% confident about our interval, what we are saying is that, **if we do the sampling(from the population) a large number of times and if we find Confidence Interval for each one, then 95% of the times actual population mean will lie in that CI.** Thus, we basically get an estimate of the range in which the population mean will lie most of the time.

Notice, that these two approaches are complementary to each other. One tells us how likely the sample statistic is and one gives us a range for the population parameter to lie in.

Confidence interval observations for Q1:

We set **alpha** to the standard value of **0.05**. This corresponds to the area on the pdf of the t-score which will be considered to be associated with the critical t- values. Since, the t-distribution is symmetric we need two critical t-values on either sides. Thus, area allocated to each side is $0.05/2 = 0.0025$.

This definition also points to the fact that the t-test and the CI definition are closely tied. The t-critical which decides which t-scores are statistically valid, also determine the ‘width’ of the CI.

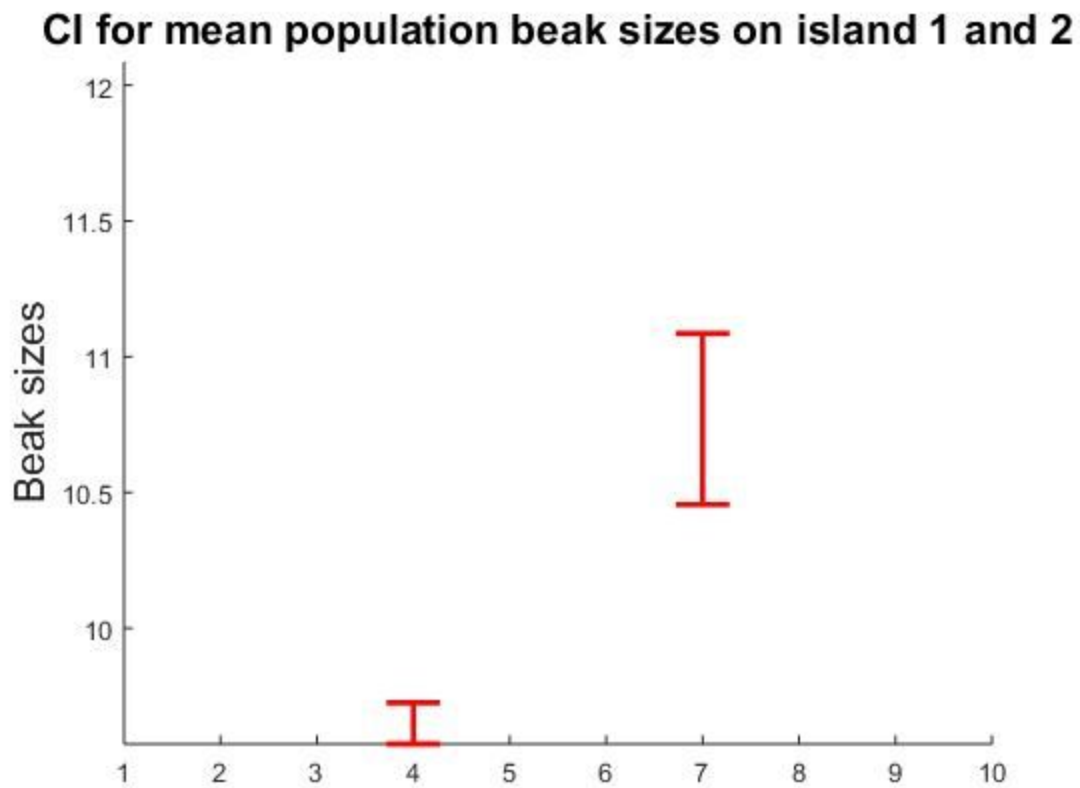
From, the derivation we get an interval in which we say that the true population mean lies more often than not provided we do the experiment a large number of times.

We calculate the CI for the given samples. The results are as follows,

CI- island 1 -> 9.5736 - 9.7255

CI- island 2 -> 10.4552 - 11.0853

The above will be clearer in the following plot.



Clearly, the two confidence intervals are NOT overlapping. This means that if we were to repeatedly sample the beak sizes on the two islands, 95% of the time, the population mean beak size for both the islands will be inside the respective CI's and hence there would be significant difference between the two means because the intervals are not overlapping.

Thus we, we can claim with statistical significance, that the mean beak size of the birds is large on island 2 than on island 1.

Approach 2 : Calculating CI in terms of difference of means

This is just another way of putting across the same idea. We instead look at the difference between the sample means and define a CI for population mean accordingly.

The calculations give the following results:

$$\text{CI_diff} = -1.4462 \quad -0.7953$$

Here also, the alpha value is 0.05. This means that 95% of the time the sampling experiment is conducted, the difference in population mean of the beak sizes will be between in the above interval.

Thus, we are very confident and can say with statistical significance that the difference is “not zero”. Moreover, since the difference is negative, it again implies that the mean beak size is larger on the second island.

Please note that this inference is subject to our choice of the significance level, the standard deviation of the data and the sample sizes taken. This is because of the way CI has been defined.

In fact, we will address one of these factors in part 2.

2. Since the sample sizes are different for the two islands, the conclusions in the previous part may be biased. Suggest and implement a strategy to verify this aspect.

From, the formula it is clear, that CI is affected by the sample size of the population. Ideally, we would want the sample sizes to be similar so as to eliminate the effect of this parameter during the comparison.

In this case unfortunately, island 1 data has 751 samples while island 2 data has only 43 samples. Hence, the CI may be biased as result. Note that larger sample sizes encourage a narrower CI.

We need to investigate if this sample size difference will affect our conclusion in Q1. Basically, we need to examine if the overlap was avoided because the larger sample size of data 1 narrowed it unfairly.

One way to tackle this problem will be to **take random sample from the larger data set so as to match the size with the smaller dataset.** However, the moment we introduce sampling, we introduce uncertainty. To make sure that the sampling does not skew our conclusions, we need to perform the sampling a large number of times[to account for the variations produced by sampling] and calculate the CI intervals.

So, we do the above for a large number of iterations [say 10000 times]. For every iteration, we sample from the larger dataset and calculate the CI. The hope is that now since the sample sizes are the same, the CI will be more indicative of the true nature of the dataset. Note that since we do it a large number of times, it should take care of the variations because of sampling.

We do two things:

a. Calculate the mean Confidence interval:

For every iteration we get a CI. We take the mean of these CI boundary values. This basically gives the interval where the population mean will lie with the given significance level considering the variations because of sampling.

The mean CI(island 1) obtained: 9.3280 9.9734

This means that the average CI for island 1 over a large number of iterations is the above for sample size equal to the second island. Thus, more often than not, **the population mean for island 1 will still be an interval which is much lesser than the interval of island 2 and there would be no overlapping of the CIs. This clearly, shows that the size was not a factor which was biasing the results and that birds in island 2 indeed have a larger mean beak size.**

b. Calculate number of times overlap happens:

One could argue that in **a. part** above, the average CI might be deceptive. It is possible for the CI for an iteration to overlap even if the average doesn't. This is unlikely to happen because as we have done the experiment a large number of times. However, just to be sure. We count the number of times the calculated CI for sampled island 1 data overlaps with island 2 sample.

The results are as follows:

```
count_overlap =
```

2

Thus, from 10000 iterations, the overlap happens only twice. Thus, we can be reasonably confident of our conclusion above that size difference is not due to any bias.

Experiment 2:

Suppose 6 different images have been processed by 2 algorithms: algorithm1 and algorithm2. In each case a number called PSNR which quantifies image quality is computed. Note that a higher PSNR implies higher quality. Load data_lab7.mat. 'psnr1' and 'psnr2' are 6-dimensional vectors whose values denote the PSNR values due to processing from algorithm1 and algorithm 2, respectively.

1. Your goal is to use CIs to establish if the difference between the mean PSNR values due to algorithm1 and algorithm2, is statistically significant. Present your findings for two different values of significance levels: 0.05 and 0.01. Do they agree? If not, which significance value is more justifiable?

Hypothesis

Null hypothesis $H_0: \mu_1 - \mu_2 = 0$

Alternate hypothesis $H_a: \mu_1 - \mu_2 \neq 0$

Where μ_1 : mean of the PSNR values generated by algorithm 1

μ_2 : mean of the PSNR values generated by algorithm 2

There are two methods to find the confidence interval.

- A. Finding Confidence Interval of individual samples and then comparing.
- B. Finding Confidence Interval for difference in means of two populations.

We will perform both the methods.

The result can be interpreted based on the following image:

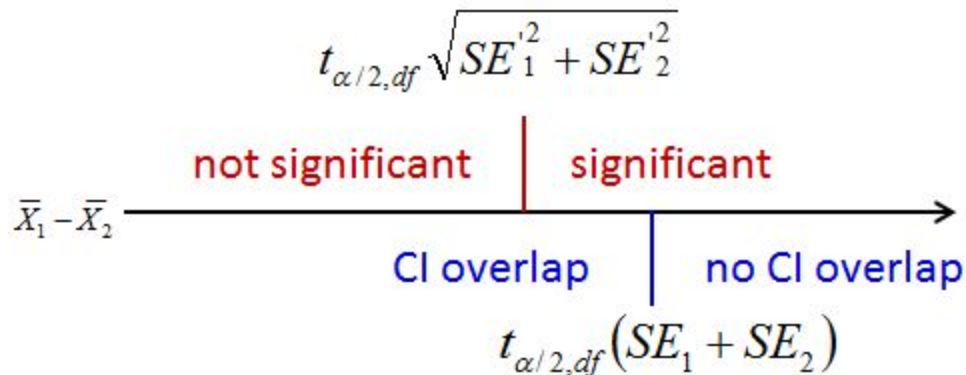


Fig.: Inference from CIs

A. Finding Confidence Interval of individual samples and then comparing.

a. Significance level $\alpha = 0.05$

Here we find the Confidence interval of given samples individually. This gives us two CIs for when significance level α is taken as 0.05.

Therefore $\alpha/2 = 0.025$. Therefore cumulative probability in two tailed t distribution = 0.975.

Degree of freedom $df = n_1 - 1 + n_2 - 1 = 10$

Therefore using t score calculator, t score = 2.571

For this t-score we can find the confidence interval using the following formula.

$$CI \text{ for } \mu: \bar{X} \pm t_{\alpha/2, df} * (s / \sqrt{n})$$

Where \bar{X} is the mean of the sample,

s : is the standard deviation of the sample

n : number of sample points

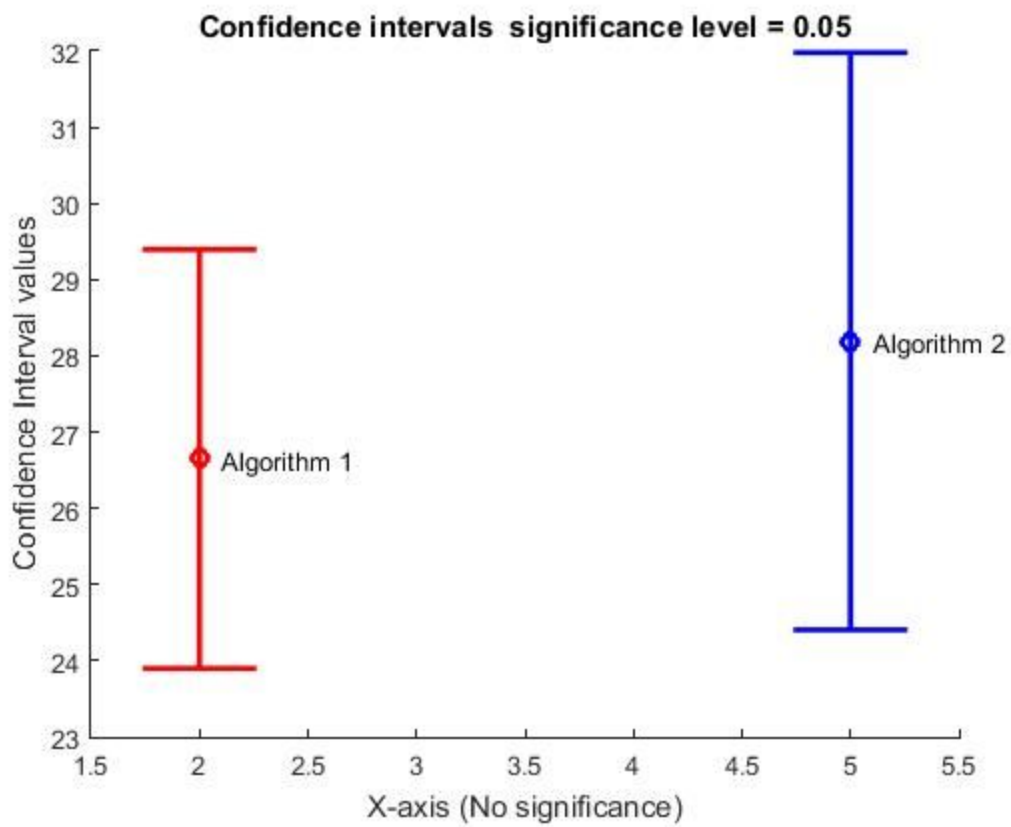
df: degree of freedom

From this, we get an interval for both the samples. The following are the endpoints of the interval.

Algorithm 1: 23.9006 - 29.3960

Algorithm 2: 24.4048 - 31.9784

Below is the graph showing the two CIs.



The above image shows that the two intervals are overlapping. Observing the Fig.: Inference from CIs, We cannot comment whether the difference between the two samples is significant or not.

b. Significance level $\alpha = 0.01$

$\alpha/2 = 0.005$. Therefore cumulative probability in two tailed t distribution = 0.995.

Degree of freedom $df = n_1 - 1 + n_2 - 1 = 10$

Therefore using t score calculator, t score = 4.032

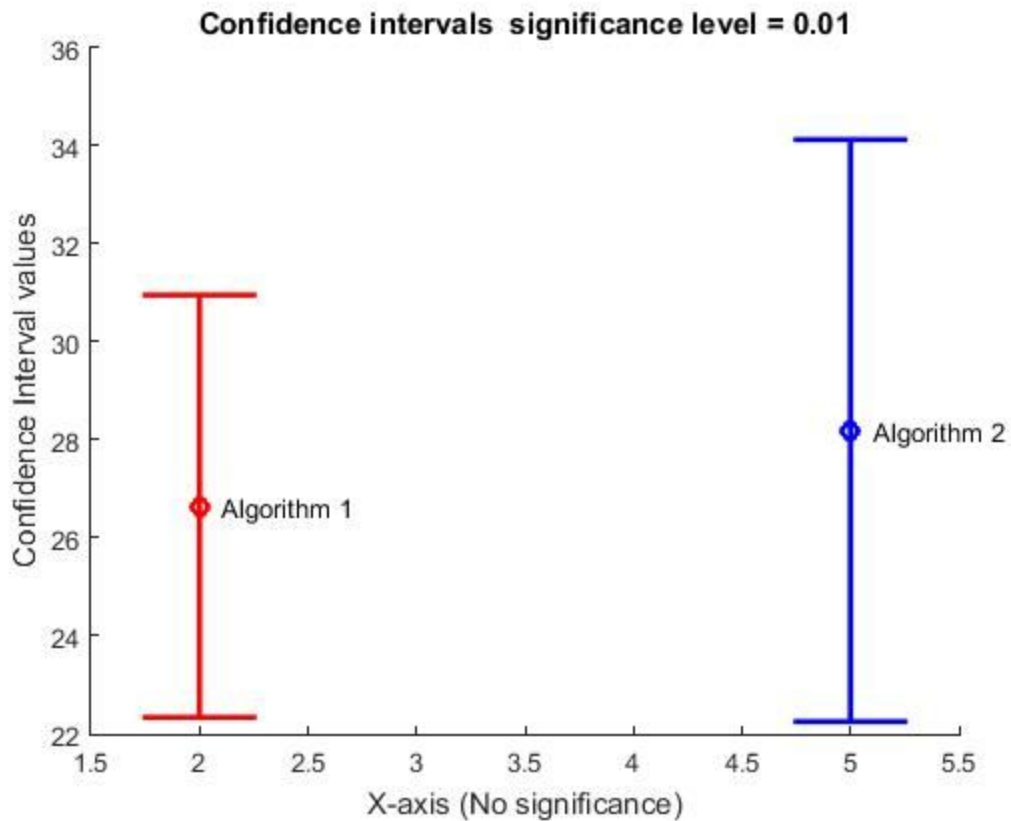
From the above formula of CI for μ , we can find two CI corresponding to two samples

The following are the endpoints of the interval.

Algorithm 1: 22.3392 - 30.9574

Algorithm 2: 22.2530 - 34.1303

Below is the graph showing the two CIs.



The above two graphs also overlap again indicating that we cannot comment on if the given samples have significant difference or not.

We can move to the second method to find out if we can categorically state any observation.

B. Finding Confidence Interval for difference in means of two populations.

We use the following formula to find CI for difference in means of two populations.

$$\text{CI for } \mu_1 - \mu_2 : \bar{X}_1 - \bar{X}_2 \pm t_{\alpha/2, n_1+n_2-2} \sqrt{\frac{s_1^2(n_1-1) + s_2^2(n_2-1)}{(n_1-1) + (n_2-1)}} \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}$$

Note here that we don't have the mean values of the two populations. We are just given two samples. Also the populations are same and the difference is the algorithm which is used to judge the quality of the images.

a. Significance level $\alpha = 0.05$

$\alpha/2 = 0.025$. Therefore cumulative probability in two tailed t distribution = 0.975.

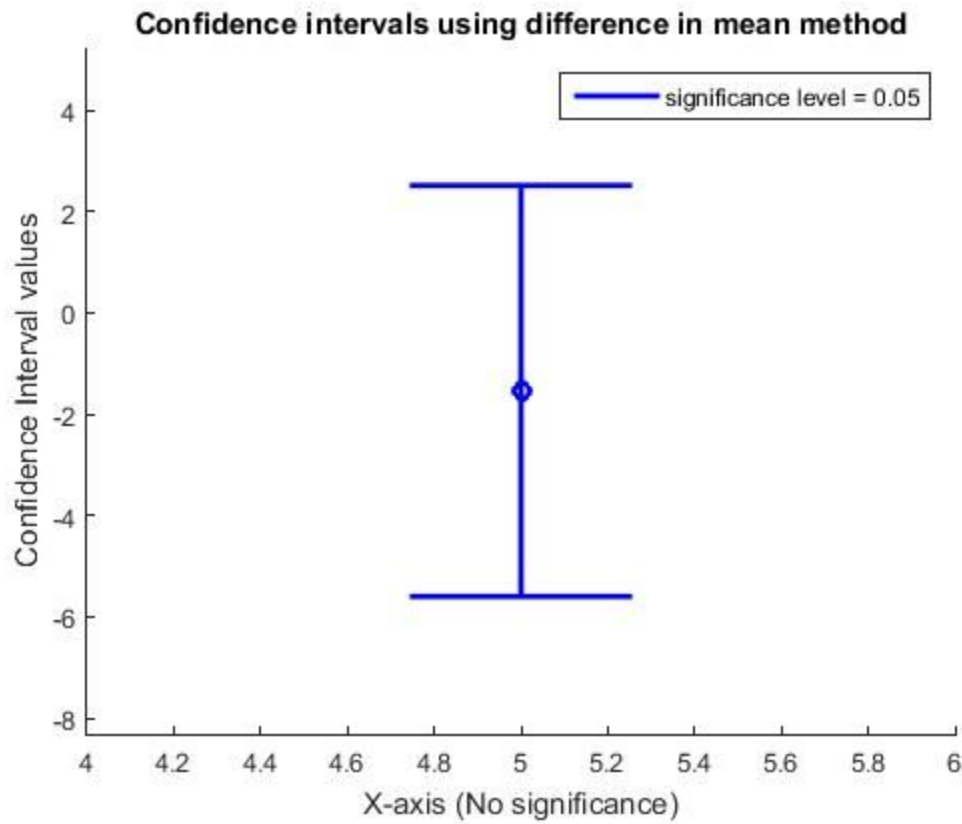
Degree of freedom $df = n_1 - 1 + n_2 - 1 = 10$

Therefore using t score calculator, t score = 2.228

From this, we get an interval for the combination. The following are the endpoints of the interval.

CI endpoints: -5.5977 and 2.5111

Below is the graph showing the CI for significance level: 0.05



In the above figure, it is clear that 0 lies in the confidence interval. This shows that the two means are very close and may be equal in some cases. By observing the Fig.: Inference from CIs, It can be concluded that there is no significant difference in the given two samples.

This means that there is no significant difference in the two algorithms used for finding quality of images.

b. Significance level $\alpha = 0.01$

$\alpha/2 = 0.005$. Therefore cumulative probability in two tailed t distribution = 0.995.

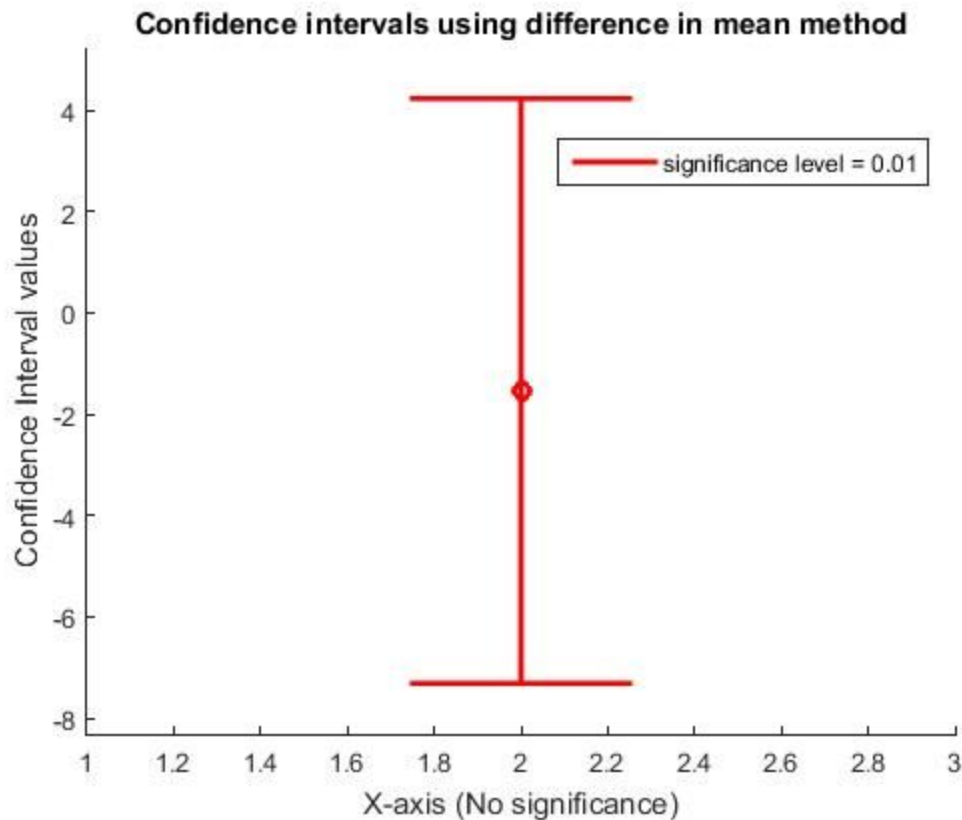
Degree of freedom $df = n_1 - 1 + n_2 - 1 = 10$

Therefore using t score calculator, t score = 3.169

From this, we get an interval for the combination. The following are the endpoints of the interval.

CI endpoints: -7.3101 and 4.2235

Below is the graph showing the CI for significance level: 0.01



In this figure also, it is clear that 0 lies in the confidence interval. This shows that the two means are very close and may be equal in some cases. It can be concluded that there is no significant difference in the given two samples.

This means that there is no significant difference in the two algorithms used for finding quality of images.

Hence by second method, we can conclude that the result for both the significance levels, 0.05 and 0.01 are same and statistically there is no significant difference in the two algorithms according to the given sample.

2. Do the conclusions in previous part agree with those made from a t-test carried out on the same data?

Since the first method in the above part does not allow us to comment on the statistical significance, t test must be performed to verify the results.

In the given experiment, we use two algorithms on the same underlying images to calculate PSNR values. The two samples are thus not independent. Hence, we will perform paired t-test.

Paired t-test

Mean of difference = -1.5433

Standard Deviation of difference = 1.6316

Standard deviation of first sample $s_1 = 2.6178$

Standard deviation of second sample $s_2 = 3.6078$

Degree of Freedom = $n - 1 = 5$

We assume that sampling distribution of the mean is normal due to Central Limit Theorem. Hence 1st assumption of t test is valid.

Here the standard deviation and consequently, the variance is of equal order ($s_2^2 < 2*s_1^2$). This satisfies the 2nd assumption for t test

Since the 6 images are different from each other as stated in the problem statement, the 3rd assumption stating the independence of the degree of freedom for t test is satisfied.

Effect size = -0.9459

Therefore | Effect Size | = 0.9459 (High)

This means that the t score is significant and we must do complete t test to get the correct statistical interpretation

t-score = -2.3169

P-value = 0.0342 (Using the t score calculator)

For Significance level: 0.05, $\alpha/2 = 0.025$

and for Significance level: 0.01, $\alpha/2 = 0.005$

The p-value obtained is greater than both, so the p value satisfied the null hypothesis that there is no statistically significant difference in the two samples which intern means that there is no statistically significant difference between the two algorithms used for rating.

Hence the conclusions made from CIs agree with the t test performed on the same data.

Code

Q1 Part 1:

```
% Author - Aditya Joglekar
% Date - 28-02-2017

% calculate the CI for the two datasets

clear;
close all;

load('island1.txt.txt');
load('island2.txt.txt');

% 0.025 0.975 % for alpha= 0.05

% for island 1
SEM1 = std(island1_txt)/sqrt(length(island1_txt));           % Standard
Error
ts1 = tinv([0.025 0.975],length(island1_txt)-1);           % T-Score, deg of
freedom, alpha= 0.05, returns 2 critical values upper and lower
CI1 = mean(island1_txt) + ts1*SEM1                           % takes care of both
upper and lower
%CI1_L= mean(island1_txt) - ts1*SEM1

% for island 2
```



```

SEM2 = std(island2_txt)/sqrt(length(island2_txt)); % Standard
Error

ts2 = tinv([0.025 0.975],length(island2_txt)-1); % T-Score, deg of
freedom, alpha= 0.05

CI2 = mean(island2_txt) + ts2*SEM2

%CI2_L= mean(island2_txt) - ts2*SEM2


axis([1 10 CI1(1) CI2(2)+1])

line([4-0.25,4+0.25],[CI1(2),CI1(2)],[0,0],'LineStyle','-','Color','r','LineWi
dth',2);

line([4-0.25,4+0.25],[CI1(1),CI1(1)],[0,0],'LineStyle','-','Color','r','LineWi
dth',2);

line([4,4],[CI1(1),CI1(2)],[0,0],'LineStyle','-','Color','r','LineWidth',2);

hold on

line([7-0.25,7+0.25],[CI2(2),CI2(2)],[0,0],'LineStyle','-','Color','r','LineWi
dth',2);

line([7-0.25,7+0.25],[CI2(1),CI2(1)],[0,0],'LineStyle','-','Color','r','LineWi
dth',2);

line([7,7],[CI2(1),CI2(2)],[0,0],'LineStyle','-','Color','r','LineWidth',2);

title('CI for mean population beak sizes on island 1 and 2','FontSize',16)

ylabel('Beak sizes','FontSize',16)

```

Q1 CI in terms of difference

```
% Author - Aditya Joglekar

% Date - 28-02-2017


% calculate CI in terms of difference of means

clear;

close all;


load('island1.txt.txt');
load('island2.txt.txt');


%hist(island2_txt)


% 0.025 0.975 % for alpha= 0.05


% for island 1

n1= length(island1_txt)
n2= length(island2_txt)


SEM_diff = sqrt( ( (std(island1_txt)^2*(n1-1) +
std(island2_txt)^2*(n2-1))/(n1+n2-2) ) *(1/n1 + 1/n2));           % Standard
Error

ts_diff = tinv([0.025 0.975],length(island1_txt)+length(island1_txt)-1);
% T-Score, deg of freedom, alpha= 0.05, returns 2 critical values upper and
lower

CI_diff = mean(island1_txt)- mean(island2_txt) + ts_diff*SEM_diff
% takes care of both upper and lower
```

Q1: Part 2

% Author - Aditya Joglekar

% Date - 28-02-2017

% Code randomly samples the larger dataset with size of smaller dataset for a large number
% of iterations, calculates CI each time and checks for overlap.

```
clear;
```

```
close all;
```

```
load('island1.txt.txt');
```

```
load('island2.txt.txt');
```

```
ite=1000;
```

```
% for island 2, no need to sample
```

```
SEM2 = std(island2_txt)/sqrt(length(island2_txt)); % Standard Error
```

```
ts2 = tinv([0.025 0.975],length(island2_txt)-1); % T-Score, deg of  
freedom, alpha= 0.05
```

```
CI2 = mean(island2_txt) + ts2*SEM2
```

```
samp_size= length(island2_txt);
```

```
sum_low=0;
```

```
sum_high=0;
```

```

count_nooverlap=0;

for i=1:ite

island1_samp = randsample(island1_txt,samp_size);

% for island 1

SEM_samp = std(island1_samp)/sqrt(length(island1_samp));           % Standard
Error

ts_samp = tinv([0.025  0.975],length(island1_samp)-1);           % T-Score, deg of
freedom, alpha= 0.05, returns 2 critical values upper and lower

CI_samp = mean(island1_samp) + ts_samp*SEM_samp;                   % takes care
of both upper and lower

%CI1_L= mean(island1_txt) - ts1*SEM1

sum_low= sum_low + CI_samp(1);
sum_high= sum_high + CI_samp(2);

if CI_samp(2)< CI2(1) | CI_samp(1) > CI2(2)
    count_nooverlap= count_nooverlap+1;
end

end

average_CI1_low= sum_low/ite
average_CI1_high= sum_high/ite

%CI2_L= mean(island2_txt) - ts2*SEM2

count_overlap= ite- count_nooverlap
count_overlap/ite

```

Q2.

Code for method 2 in part 1 and for part 2

```
% Author - Rajdeep Pinge
% Date - 28-02-2017

% Code to statistically compare two algorithms of image grading using
% CI method for difference of population mean and using paired t test

clear;
close all;

load('data_lab7');

% find mean of samples
mean_psnr1 = mean(psnr1);
mean_psnr2 = mean(psnr2);

% find standard deviation of samples
s1 = std(psnr1);
s2 = std(psnr2);
```

```

% find variance of samples

var1 = s1*s1;

var2 = s2*s2;


n1 = length(psnr1);
n2 = length(psnr2);


% degree of freedom

df = n1+n2-2;


% t scores using calculator

t_score_95 = 2.228;
t_score_99 = 3.169;


% find the limits of CI


% significance level = 0.05

CI_95_high = mean_psnr1 - mean_psnr2 + t_score_95 * sqrt( ((var1*(n1-1) +
var2*(n2-1))/df) * (1/n1 + 1/n2) );

CI_95_low = mean_psnr1 - mean_psnr2 - t_score_95 * sqrt( ((var1*(n1-1) +
var2*(n2-1))/df) * (1/n1 + 1/n2) );


% significance level = 0.01

CI_99_high = mean_psnr1 - mean_psnr2 + t_score_99 * sqrt( ((var1*(n1-1) +
var2*(n2-1))/df) * (1/n1 + 1/n2) );

CI_99_low = mean_psnr1 - mean_psnr2 - t_score_99 * sqrt( ((var1*(n1-1) +
var2*(n2-1))/df) * (1/n1 + 1/n2) );


% Plot CIs

axis([1 3 CI_99_low-1 CI_99_high+1])

```



```

% take difference of sample values
diff_psnr = psnr1 - psnr2;

% find mean of difference of sample values
mean_psnr = mean(diff_psnr);

% find standard deviation of difference of sample values
sd_psnr = std(diff_psnr);

% calculate degree of freedom for the paired test sample which is (n-1)
deg_of_freedom = length(diff_psnr) - 1;

% calculate effect size
effect_size = mean_psnr / sd_psnr

% find t value for paired t test
paired_t = effect_size * sqrt(deg_of_freedom + 1)

% p value as calculated on the t score calculator
pval = 0.0342

%% The p value is well within the significant level range

```

Code for method 1 in part 1

```

% Author - Rajdeep Pinge
% Date - 28-02-2017

```



```
% Code to statistically compare two algorithms of image grading using
% separate CI calculation for each sample

clear;

close all;

load('data_lab7');

% find mean of samples
mean_psnr1 = mean(psnr1);
mean_psnr2 = mean(psnr2);

% find standard deviation of samples
s1 = std(psnr1);
s2 = std(psnr2);

% find variance of samples
var1 = s1*s1;
var2 = s2*s2;

n1 = length(psnr1);
n2 = length(psnr2);

% degree of freedom
df1 = n1-1;
df2 = n2-1;

% t scores using calculator
```

```

t_score_95 = 2.571;
t_score_99 = 4.032;

% find the limits of CI for 1st sample
% significance level = 0.05
CI_95_high_1 = mean_psnr1 + t_score_95 * s1/sqrt(n1);
CI_95_low_1 = mean_psnr1 - t_score_95 * s1/sqrt(n1);
% significance level = 0.01
CI_99_high_1 = mean_psnr1 + t_score_99 * s1/sqrt(n1);
CI_99_low_1 = mean_psnr1 - t_score_99 * s1/sqrt(n1);

% find the limits of CI for 2nd sample
% significance level = 0.05
CI_95_high_2 = mean_psnr2 + t_score_95 * s2/sqrt(n2);
CI_95_low_2 = mean_psnr2 - t_score_95 * s2/sqrt(n2);
% significance level = 0.01
CI_99_high_2 = mean_psnr2 + t_score_99 * s2/sqrt(n2);
CI_99_low_2 = mean_psnr2 - t_score_99 * s2/sqrt(n2);

% plot CI of both for significance level = 0.01
line([2-0.25,2+0.25],[CI_99_high_1,CI_99_high_1],[0,0],'LineStyle','-','Color','r','LineWidth',2);
line([2-0.25,2+0.25],[CI_99_low_1,CI_99_low_1],[0,0],'LineStyle','-','Color','r','LineWidth',2);
line([2,2],[CI_99_low_1,CI_99_high_1],[0,0],'LineStyle','-','Color','r','LineWidth',2);
hold on

```

```

plot([2,2], mean_psnr1, 'ro', 'LineWidth',2)
text(2+0.1,mean_psnr1,'Algorithm 1')

line([5-0.25,5+0.25],[CI_99_high_2,CI_99_high_2],[0,0], 'LineStyle','-','Color',
'b','LineWidth',2);
line([5-0.25,5+0.25],[CI_99_low_2,CI_99_low_2],[0,0], 'LineStyle','-','Color','
b','LineWidth',2);
line([5,5],[CI_99_low_2,CI_99_high_2],[0,0], 'LineStyle','-','Color','b','LineW
idth',2);

hold on

plot([5,5], mean_psnr2, 'bo', 'LineWidth',2)
text(5+0.1,mean_psnr2,'Algorithm 2')

title('Confidence intervals  significance level = 0.01')
xlabel('X-axis (No significance)')
ylabel('Confidence Interval values')

% plot CI of both for significance level = 0.05

figure

line([2-0.25,2+0.25],[CI_95_high_1,CI_95_high_1],[0,0], 'LineStyle','-','Color',
'r','LineWidth',2);
line([2-0.25,2+0.25],[CI_95_low_1,CI_95_low_1],[0,0], 'LineStyle','-','Color','
r','LineWidth',2);
line([2,2],[CI_95_low_1,CI_95_high_1],[0,0], 'LineStyle','-','Color','r','LineW
idth',2);

hold on

plot([2,2], mean_psnr1, 'ro', 'LineWidth',2)
text(2+0.1,mean_psnr1,'Algorithm 1')

line([5-0.25,5+0.25],[CI_95_high_2,CI_95_high_2],[0,0], 'LineStyle','-','Color',
'b','LineWidth',2);

```

```
line([5-0.25,5+0.25],[CI_95_low_2,CI_95_low_2],[0,0], 'LineStyle', '-', 'Color', 'b', 'LineWidth', 2);

line([5,5],[CI_95_low_2,CI_95_high_2],[0,0], 'LineStyle', '-', 'Color', 'b', 'LineWidth', 2);

hold on

plot([5,5], mean_psnr2, 'bo', 'LineWidth', 2)

text(5+0.1, mean_psnr2, 'Algorithm 2')

title('Confidence intervals  significance level = 0.05')

xlabel('X-axis (No significance)')

ylabel('Confidence Interval values')
```