

CSE 544

**PROBABILITY AND STATISTICS
FOR DATA SCIENCE**

ASSIGNMENT 5

Team members:

Name	SBU ID
Venkata Satwik Chilukuri	113263262
Akshay Somayaji	113322316
Shardul Churi	114403283

1)

* Given 10 samples,

$$D = \{2.78, 0.84, 1.88, 2.23, 1.99, 0.04, 2.65, 0.74, 1.19, 2.34\}$$

Null hypothesis (H_0): The samples are from the uniform $(0, 3)$ distribution.

Alternate hypothesis: The samples are not from the uniform $(0, 3)$ distribution.

Sorting the given data,

$$D = \{0.04, 0.74, 0.84, 1.19, 1.88, 1.99, 2.23, 2.57, 2.65, 2.78\}$$

Given,

$\hat{F}_x^-(x) = \text{ecdf to the left of } x$

$\hat{F}_x^+(x) = \text{ecdf to the right of } x$

$f_y(x) = \text{cdf of uniform}(0, 3) \text{ at } x$

So, $f_y(x) = \frac{x-0}{3-0} \quad \left(\text{since, cdf of uniform}(0, 3) = \frac{x-0}{3-0} \right)$

x	$f_y(x)$	$\hat{F}_x(x)$	$\hat{F}_{x^*}(x)$	$ \hat{F}_x(x) - f_y(x) $	$ \hat{F}_{x^*}(x) - f_y(x) $
0.04	0.0133	0	$\frac{1}{10}$	0.0133	0.0867
0.74	0.2466	$\frac{1}{10}$	$\frac{2}{10}$	0.1466	0.0466
0.84	0.28	$\frac{2}{10}$	$\frac{3}{10}$	0.08	0.02
1.19	0.3966	$\frac{1}{10}$	$\frac{4}{10}$	0.0966	0.0034
1.88	0.626	$\frac{4}{10}$	$\frac{5}{10}$	0.226	0.126
1.99	0.663	$\frac{5}{10}$	$\frac{6}{10}$	0.163	0.063
2.23	0.743	$\frac{6}{10}$	$\frac{7}{10}$	0.143	0.043
2.57	0.856	$\frac{7}{10}$	$\frac{8}{10}$	0.156	0.056
2.65	0.883	$\frac{8}{10}$	$\frac{9}{10}$	0.683	0.017
2.78	0.926	$\frac{9}{10}$	$\frac{10}{10}$	0.026	0.074

Maximum value = 0.226

Let ecdf represent \hat{F}_D

Now,

$$d = \max_{\alpha} |\hat{F}_D(\alpha) - F_x(\alpha)|$$

$$\approx 0.226$$

Given threshold is 0.25

Here, d is less than given threshold so Accept the null hypothesis (H_0)

The given samples are from the uniform(0,3) distribution

Given data,

$$x = \{5\} \quad , \quad y = \{2, 3\}$$

So,

$$N = |x| + |y| = 3$$

T_{observed} is the sampled mean difference of x and y .

i.e.,

$$\begin{aligned} T_{\text{obs}} &= |\bar{x} - \bar{y}| = |5 - \left(\frac{2+3}{2}\right)| \\ &= |5 - 4.5| \\ &= 0.5 \end{aligned}$$

In the table below, we will list the permutations for each permutation we have listed the sample

values of x_i, y_i, T_i :

i	x_i	y_i	$T_i = \bar{x}_i - \bar{y}_i $
1	5	{2, 3}	0.5
2	5	{3, 2}	0.5
3	2	{5, 3}	4
4	2	{3, 5}	4
5	3	{2, 5}	3.5
6	3	{5, 2}	3.5

Now,

we know the Indicator R.V for an Event E

$$I(E) = \begin{cases} 0 & \text{if } E \text{ does not occur} \\ 1 & \text{if } E \text{ occurs} \end{cases}$$

So,

$$T_{\text{obs}} = 0.5$$

$T_i = \bar{x}_i - \bar{y}_i $	$I(T_i > T_{\text{obs}})$
0.5	0
0.5	0
4	1
4	1
3.5	1
3.5	1

Now,

$$\text{P-value} = \frac{4}{6} = \frac{2}{3} = 0.666$$

So

if the threshold is 0.05 and the p-value here is not less than 0.05. So we can't reject the null hypothesis.

Thus x and y are from same distribution.

QUESTION 3.(a) INDEPENDENCE TEST TO SAVE YOUR CASINO

lets assume the dealers and outcome to be independent of each other.

H_0 : outcome of the table independent of the dealer.

H_1 : outcome of the table not independent of the dealer

The OBSERVED OUTCOME.

	A	B	C	Total
Win	48	54	19	121
Draw	7	5	4	16
Lose	55	50	25	130
Total	110	109	48	267

Considering the total values, we calculate the expected values.

	A	B	C	Total
Win	49.85	49.4	21.75	121
Draw	6.59	6.53	2.88	16
Lose	53.56	53.07	23.37	130
Total	110	109	48	267

Now, calculating χ^2_{obs}

$$= \sum_{\text{rows}=i}^n \sum_{\text{cols}=j}^n \frac{(\bar{E}_{ij} - O_{ij})^2}{\bar{E}_{ij}}$$

$$= \frac{(49.85 - 48)^2}{49.85} + \frac{(49.4 - 54)^2}{49.4} + \frac{(21.75 - 19)^2}{21.75}$$

$$+ \frac{(6.59 - 7)^2}{6.59} + \frac{(6.53 - 5)^2}{6.53} + \frac{(2.88 - 4)^2}{2.88}$$

$$+ \frac{(53.56 - 55)^2}{53.56} + \frac{(53.07 - 50)^2}{53.07} + \frac{(23.37 - 25)^2}{23.37}$$

$$= 0.0689 + 0.428 + 0.348$$

$$+ 0.024 + 0.358 + 0.314$$

$$+ 0.038 + 0.188 + 0.106$$

$$= \underline{\underline{1.8733}}$$

Calculating df : degree of freedom.

$$= (\# \text{ of rows} - 1) (\# \text{ of columns} - 1)$$

$$= (4 - 1)(3 - 1)(3 - 1)$$

$$= 2 \times 2$$

$$= 4$$

Calculating p value for 1.8733 and $df = 4$

$$= 1 - 0.241$$

[using tables]

$$= \underline{\underline{0.759}}$$

As $p > 0.05 [\alpha]$, we can't reject H_0 .

∴ we conclude that the outcome of the table is independent of the dealer.

(b) To find: Pearson Correlation Coefficient

DAY A B C

1 48 54 19

2 40 48 40

3 58 51 35

4 53 47 41

5 65 62 38

6 25 35 32

7 52 70 32

8 34 20 37

9 30 25 37

10 45 40 15

$$\bar{A} = \frac{\sum A_i}{n} = \frac{450}{10}$$

$$= 45$$

$$\bar{B} = \frac{\sum B_i}{n} = \frac{452}{10}$$

$$= 45.2$$

$$\bar{C} = \frac{\sum C_i}{n} = \frac{326}{10}$$

$$= 32.6$$

$$\hat{S}_{A,B} = \sum_{i=1}^n \{ (A_i - \bar{A})(B_i - \bar{B}) \}$$

$$\sqrt{(\sum (A_i - \bar{A})^2) \cdot (\sum (B_i - \bar{B})^2)}$$

$$= \frac{1396}{\sqrt{1462 \times 2193.6}} = \frac{1396}{1790.82} = 0.779$$

$$\hat{S}_{B,C} = \sum_{i=1}^n \{ (B_i - \bar{B})(C_i - \bar{C}) \}$$

$$\sqrt{(\sum (B_i - \bar{B})^2) \cdot (\sum (C_i - \bar{C})^2)}$$

$$= \frac{-96.2}{\sqrt{2193.6 \times 694.4}} = \frac{-96.2}{1234.19} = -0.078$$

$$\hat{\rho}_{A,C} = \sum_{i=1}^n \{ (A_i - \bar{A})(C_i - \bar{C}) \}$$

$$\sqrt{(\sum (A_i - \bar{A})^2)(\sum (C_i - \bar{C})^2)}$$

$$\sqrt{\frac{1962}{22} \times \frac{694.4}{22}} = \frac{1007.58}{0.022} = 0.022$$

$$\hat{\rho}_{A,B} = 0.779 \quad \hat{\rho}_{B,C} = -0.048 \quad \hat{\rho}_{A,C} = 0.022$$

There seems to be a +ve correlation between Dealers (A and B)

However there seems to be almost no correlation between (A and C) or (B and C). With a value of +0.779, the correlation between Dealer A and Dealer B appears to be quite strong with value of A increasing or decreasing in accordance with value of B.

Hence, we can say that Dealer C may not be loyal, and could be the cause of the money loss.

- 4) Refer to q4.py, and the eCDF graph q4 KS eCDF.png

5) Given, $O_1 = \{x_1, x_2, \dots, x_n\} \stackrel{iid}{\sim} \text{Normal}(\mu_1, \sigma_1^2)$

$O_2 = \{y_1, y_2, \dots, y_m\} \stackrel{iid}{\sim} \text{Normal}(\mu_2, \sigma_2^2)$

x 's \perp y 's, n and m are large.

s_x and s_y are the sample std deviations of O_1 & O_2 .

→ Also, $H_0: \mu_1 > \mu_2$ and $H_1: \mu_1 \leq \mu_2$ with $\delta > 0$ as the critical value

a) To calculate: Type-1 & Type-2 errors for the unpaired T-test using O_1 & O_2

→ Let's define $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$, $\bar{y} = \frac{1}{m} \sum_{i=1}^m y_i$

$$\Rightarrow \bar{O} = \bar{x} - \bar{y} \quad (\text{i})$$

→ Now, since n & m are assumed to be large, we have

$\bar{x} \sim \text{Normal}(\mu_1, \frac{\sigma_1^2}{n})$ (via CLT)

$\bar{y} \sim \text{Normal}(\mu_2, \frac{\sigma_2^2}{m})$

$$\Rightarrow \bar{O} \sim \text{Normal}(\mu_1 - \mu_2, \frac{\sigma_1^2}{n} + \frac{\sigma_2^2}{m}) \quad (\text{ii})$$

($\because \bar{x} - \bar{y}$ follows weighted sum of \perp Normals)

Now, T-statistic is given by

$$T = \frac{\bar{O}}{\sqrt{\frac{s_x^2}{n} + \frac{s_y^2}{m}}} \quad \text{for an unpaired T-test}$$

i) For Type-I error, we have

$$\Pr(\text{Type-I error}) = \Pr(\text{False Positive}) \\ = \Pr(\text{Reject } H_0 \mid H_0 \text{ is true})$$

For a one-tailed T-test with $H_0: \mu_1 \geq \mu_2$

& $H_1: \mu_1 < \mu_2$, we reject H_0 when $T < -\delta$

$$\Rightarrow \Pr(\text{Type I error}) = \Pr(T < -\delta \mid \mu_1 > \mu_2)$$

$$= \Pr\left(\frac{\bar{D}}{\sqrt{\frac{s_x^2}{n} + \frac{s_y^2}{m}}} < -\delta \mid \mu_1 > \mu_2\right)$$

$$= \Pr\left(\bar{D} < -\delta \sqrt{\frac{s_x^2}{n} + \frac{s_y^2}{m}} \mid \mu_1 > \mu_2\right)$$

$$= F_{\bar{D}}\left(-\delta \sqrt{\frac{s_x^2}{n} + \frac{s_y^2}{m}}\right) \quad (\text{CDF of dist}^n \bar{D} \text{ at point } D = -\delta \sqrt{\frac{s_x^2}{n} + \frac{s_y^2}{m}}) \\ \text{(iii)}$$

→ Now, since final result is in terms of \bar{D} , let's transform \bar{D} into standard normal $Z(0, 1)$.

From transformation property, we have at any point D , $Z = \frac{D - \mu_{\bar{D}}}{\sigma_{\bar{D}}} \quad (\text{iv})$

$$\text{From (ii), } \mu_{\bar{D}} = \mu_1 - \mu_2, \quad \sigma_{\bar{D}} = \sqrt{\frac{\sigma_1^2}{n} + \frac{\sigma_2^2}{m}}$$

→ Now since m & n are large, we can plugin true, unknown std deviations σ_1 & σ_2 with the MLE estimator for std deviation of normal dist n , which is the sample std deviations s_x & s_y respectively, via Asymptotic Normality.

$$\Rightarrow (\text{iv}) \text{ becomes } Z = \frac{D - (\mu_1 - \mu_2)}{\sqrt{\frac{s_x^2}{n} + \frac{s_y^2}{m}}}$$

$$\text{At } D = -\delta \sqrt{\frac{s_x^2}{n} + \frac{s_y^2}{m}}, \quad Z = -\delta \sqrt{\frac{s_x^2}{n} + \frac{s_y^2}{m}} - (\mu_1 - \mu_2)$$

$$\sqrt{\frac{s_x^2}{n} + \frac{s_y^2}{m}}$$

$$\therefore Z = -\delta - \frac{\mu_1 - \mu_2}{\sqrt{\frac{s_x^2}{n} + \frac{s_y^2}{m}}} \quad (\text{v})$$

Applying this transformation, (iii) becomes

$$\Pr(\text{Type I error}) = F_Z(-\delta - \frac{\mu_1 - \mu_2}{\sqrt{\frac{s_x^2}{n} + \frac{s_y^2}{m}}})$$

$$\therefore \Pr(\text{Type I error}) = \Phi\left(-\delta - \frac{\mu_1 - \mu_2}{\sqrt{\frac{s_x^2}{n} + \frac{s_y^2}{m}}}\right)$$

ii) Type 2 error probability can be written as

$$\begin{aligned} \Pr(\text{Type 2 error}) &= \Pr(\text{False Negative}) \\ &= \Pr(\text{Accept } H_0 \mid H_0 \text{ is false}) \end{aligned}$$

For a one-tailed T -test with $H_0: \mu_1 \geq \mu_2$

& $H_1: \mu_1 \leq \mu_2$, we accept H_0 when $T \geq -\delta$

$$\begin{aligned} \Rightarrow \Pr(\text{Type 2 error}) &= \Pr(T \geq -\delta \mid \mu_1 \leq \mu_2) \\ &= 1 - \Pr(T \leq -\delta \mid \mu_1 \leq \mu_2) \\ &= 1 - \Pr\left(\bar{D} \leq -\delta \mid \mu_1 \leq \mu_2\right) \\ &= 1 - \Pr\left(\bar{D} \leq -\delta \sqrt{\frac{s_x^2}{n} + \frac{s_y^2}{m}} \mid \mu_1 \leq \mu_2\right) \\ &= 1 - F_{\bar{D}}\left(-\delta \sqrt{\frac{s_x^2}{n} + \frac{s_y^2}{m}}\right) \quad (\text{vi}) \end{aligned}$$

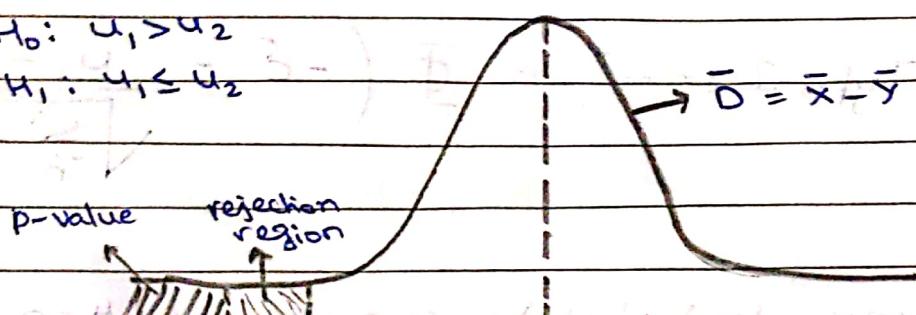
Using Z -transformation properties as in part i) and using results of (iii), (iv) & (v), we get

$$\Pr(\text{Type 2 error}) = 1 - \Phi\left(-\delta - \frac{\mu_1 - \mu_2}{\sqrt{\frac{s_x^2}{n} + \frac{s_y^2}{m}}}\right)$$

b) For a 1-tailed T-test, we can represent the region of p-value as below -

$$H_0: \mu_1 > \mu_2$$

$$H_1: \mu_1 \leq \mu_2$$



Now, p-value = $\Pr(\text{finding more extreme statistic than } t_0 \mid H_0 \text{ is true})$

$$\text{Let } t_0 = \frac{\bar{x} - \bar{y}}{\sqrt{\frac{s_x^2}{n} + \frac{s_y^2}{m}}}$$

$$\Rightarrow \text{p-value} = \Pr(\bar{D} < \bar{x} - \bar{y})$$

$$= F_{\bar{D}}(\bar{x} - \bar{y})$$

Using (iv), we can transform into std normal by

$$\sqrt{\frac{s_x^2}{n} + \frac{s_y^2}{m}}$$

$$\text{p-value} = F_z\left(\frac{\bar{x} - \bar{y} - (\mu_1 - \mu_2)}{\sqrt{\frac{s_x^2}{n} + \frac{s_y^2}{m}}}\right)$$

$$\boxed{\text{p-value} = \Phi\left(\frac{\bar{x} - \bar{y} - (\mu_1 - \mu_2)}{\sqrt{\frac{s_x^2}{n} + \frac{s_y^2}{m}}}\right)}$$

(a) Given, $X \sim \text{Normal}(1.5, 1)$ and $Y \sim \text{Normal}(1, 1)$

i) Z-Test for 20 sample x_i & y_i -

$\sigma_{x_1} = 1$, $\sigma_{y_1} = 1$, $\bar{x}_1 = 1.5987$, $\bar{y}_1 = 1.0625$, $n = m = 20$

sample pool std deviation = $\sqrt{\frac{\sigma_{x_1}^2}{n} + \frac{\sigma_{y_1}^2}{m}} = \sqrt{\frac{1}{20} + \frac{1}{20}} = \sqrt{0.1} = 0.3162$

Since we're assuming null hypothesis, the hypothesized difference between sample means = 0

$$\Rightarrow Z = \frac{\bar{x}_1 - \bar{y}_1}{\text{sample pool std dev}} = \frac{1.5987 - 1.0625}{0.3162}$$

$$\Rightarrow Z_{20} = 1.6957 \quad (\text{i})$$

$$Z \text{ p-value}_{20} = 2 * (1 - \Phi(1.6957))$$

$$= 2 * (1 - \Phi(1.6957))$$

$$= 0.0899 \quad (\text{ii})$$

Given, critical threshold = 1.962 = c

since $|1.6957| < c$, we fail to reject the null hypothesis, and accept that the mean values of X and Y are equal.

ii) Unpaired T-test for 20 samples of x_i & y_i -

$$\Rightarrow S_{x_1}^2 = 1.1772, S_{y_1}^2 = 1.1412$$

now, $T = \frac{\bar{x}_1 - \bar{y}_1}{\sqrt{\frac{S_{x_1}^2}{n} + \frac{S_{y_1}^2}{m}}} = \frac{1.5987 - 1.0625}{\sqrt{\frac{1.1772}{20} + \frac{1.1412}{20}}}$

$$\Rightarrow T_{20} = 1.5750 \quad (\text{iii})$$

$$T \text{ p-value}_{20} = 2 * (1 - F_{T_{DF}}(1.5750))$$

Now, degrees of freedom $df = n+m-2 = 38$

$$\Rightarrow T \text{ p-value} = 2 * (1 - F_{38}(1.5750)) = 0.1236 \quad (\text{iv})$$

Notes : $S_{x_1}^2 = \sum_{i=1}^{20} (x_i - \bar{x}_1)^2 / n-1 = 1.1772$

$$S_{y_1}^2 = \sum_{i=1}^{20} (y_i - \bar{y}_1)^2 / m-1 = 1.1412$$

Given, critical threshold $c = 2.086$
 since $|T| = 1.5750 < c$, we fail to reject
 the null hypothesis, and accept that the
 mean values of X and Y are equal.

b) i) Z-test for 1000 sample $X_2 \& Y_2$

$$\sigma_{X_2} = 1, \sigma_{Y_2} = 1, \bar{X}_2 = 1.4613, \bar{Y}_2 = 0.9835$$

$m = n = 1000$, sample pool std dev = 0.0447

$$\Rightarrow Z_{1000} = \frac{\bar{X}_2 - \bar{Y}_2}{\text{std dev}} = \frac{1.4613 - 0.9835}{0.0447} = 10.6826 \quad (\text{v})$$

$$Z \text{ p-value}_{1000} = 2(1 - \Phi(10.6826)) \approx 0.0000 \quad (\text{vi})$$

Given, critical value $c = 1.962$

since $|Z| = 10.6826 > c$, we reject the
 null hypothesis and accept the alternate one
 that the means of X and Y are not equal

iii) T-test for 1000 sample $X_2 \& Y_2$

$$S_{X_2}^2 = 1.0299, S_{Y_2}^2 = 0.9586$$

$$\Rightarrow T_{1000} = \frac{\bar{X}_2 - \bar{Y}_2}{\sqrt{\frac{1.0299}{1000} + \frac{0.9586}{1000}}} = 10.7134 \quad (\text{vii})$$

$$T \text{ p-value}_{1000} = 2(1 - F_{T_{998}}(10.7134)) \approx 0.000 \quad (\text{viii})$$

Given, critical value $c = 1.962$

Since $|T| = 10.7134 > c$, we reject the
 null hypothesis and accept the alternate
 one that the means of X and Y are
 not equal

Observations -

1. We got a True Positive result in part b), as using both Z and T tests, we were successfully able to reject the null hypothesis, that is, infer that the means of X and Y are not same, given we already have the ground truth that mean of X is 1.5 and mean of Y is 1. We can hence claim that both Z and T tests work well with large data samples
2. There is not much of a significant advantage of using Z test on smaller samples, as both tests on the 20-sample data failed to reject the null hypothesis (even though ground truth says otherwise), thereby resulting in a false negative. However, the p-value is much closer to $\alpha=0.05$ in the case of the Z-test (0.0899) than in the T-test (0.1236)
3. We can see that the p-value is significantly lesser than $\alpha=0.05$ in part b) with 1000 samples, for both Z and T tests. Therefore, we can support our rejection of the null hypothesis with extremely high confidence, since $0.0000 \lll 0.05$.