Lorenzo  Beltrame          08/02/2022          MAT : 12137286

<u>Statistic take home exam</u>

$Z$ test :          $X_1, \dots, X_n \sim N(\mu_x, 1)$   i.i.d.

$Y_1, \dots, Y_n \sim N(\mu_y, 1)$   i.i.d. $\}$ indipendently from each other

<u>def</u>    $\bar{X}_n = \dfrac{1}{n} \displaystyle\sum_{i=1}^{n} X_i$  ,        $\bar{Y}_n = \dfrac{1}{n} \displaystyle\sum_{i=1}^{n} Y_i$

<u>def</u>   test Statistic    $S = \dfrac{(\bar{Y}_n - \bar{X}_n)}{\sqrt{\frac{2}{n}}}$

<u>Valid p-values</u> :   $\cdot \ P_A := 1 - \Phi(S)$    for   $H_0 : \mu_y = \mu_x$

against   $H_1^A : \mu_y > \mu_x$

$\cdot \ P_B := 1 - \Phi(-S)$   for   $H_0 : \mu_y = \mu_x$

against   $H_1^B : \mu_y < \mu_x$

<u>Obj.</u>   test   $H_0 : \mu_x = \mu_y$   against   $H_1 : \mu_y \neq \mu_x$

$\begin{cases} \text{if} \ \bar{Y}_n > \bar{X}_n \ \text{report} \ P_A \\ \text{if} \ \bar{Y}_n < \bar{X}_n \ \text{report} \ P_B \end{cases}$          [1]

**a)** <u>Is this a valid p-value?</u>

We can write [1] as:

$$p(x) = \begin{cases} p_A = 1 - \Phi(S(x)) & \text{if } \bar{Y}_n > \bar{x}_n \\ \\ p_B = 1 - \Phi(-S(x)) & \text{if } \bar{Y}_n < \bar{x}_n \end{cases}$$

$$= \begin{cases} 1 - \Phi\left[\dfrac{(\bar{Y}_n - \bar{x}_n)}{\sqrt{\frac{2}{n}}}\right] & \text{if } \bar{Y}_n < \bar{x}_n \\ \\ 1 - \Phi\left[\dfrac{\bar{x}_n - \bar{Y}_n}{\sqrt{2/n}}\right] & \text{if } \bar{Y}_n < \bar{x}_n \end{cases}$$

<span style="color:magenta">Since we are changing the sign we can use the abs. value</span>

$$p(x) = 1 - \Phi\left[\left|\bar{Y}_n - \bar{x}_n\right| \cdot \sqrt{\frac{n}{2}}\right]$$
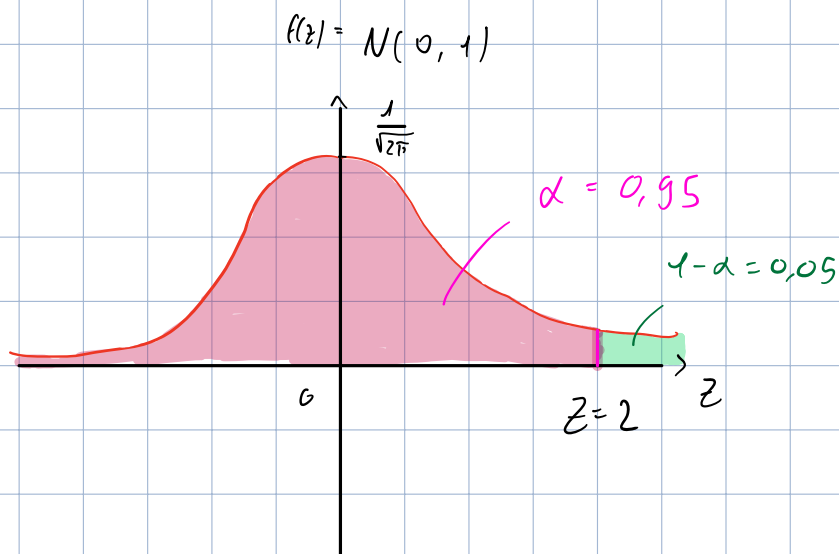
We know that a valid p value satisfyies:

$$\alpha = P_{H_0}\left(p(x) < \alpha\right)$$

$$= P_{H_0}\left[1 - \Phi\left(\left|\bar{Y}_n - \bar{x}_n\right|\sqrt{\frac{n}{2}}\right) < \alpha\right] =$$

$$= P_{H_0}\left[\Phi\left(\left|\bar{Y}_n - \bar{x}_n\right|\sqrt{\frac{n}{2}}\right) > 1 - \alpha\right] =$$

$$= P_{H_0}\left[\left|\bar{Y}_n - \bar{x}_n\right|\sqrt{\frac{n}{2}} > \Phi^{-1}(1-\alpha)\right] =$$

$f(z) = N(0, 1)$

$\frac{1}{\sqrt{2\pi}}$

$\alpha = 0.95$

$1 - \alpha = 0.05$

$0$

$Z = 2$

$z$

this is the valid p-value for test A
$$P_{H_0}\left(1 - \Phi(s) > \alpha\right) = \alpha$$

$$\Rightarrow P_{H_0}\left[\left(\bar{Y}_n - \bar{X}_n\right)\sqrt{\frac{n}{2}} > \Phi^{-1}(1-\alpha)\right] +$$

$$+ P_{H_0}\left[\left(\bar{Y}_n - \bar{X}_n\right)\sqrt{\frac{n}{2}} < \Phi^{-1}(1-\alpha)\right] =$$

this is the valid p-value for test B
$$P_{H_0}\left(1 - \Phi(-s) > \alpha\right) = \alpha$$

Since $P_A$ and $P_B$ are valid P-values

$$= \alpha + \alpha = 2\alpha$$

We got $\alpha = 2\alpha$ which is valid only for $\alpha = 0$.

$$\Rightarrow P(x) = 1 - \Phi\left[\left|\bar{Y}_n - \bar{X}_n\right|\sqrt{\frac{n}{2}}\right] \text{ is not a valid pvalue.}$$

**b)** <u>Correct the p-value:</u>

I can correct the p value considering

$$p(x) = 2 \left[ 1 - \Phi\left( |\bar{Y}_n - \bar{X}_n| \sqrt{\frac{n}{2}} \right) \right]$$

And we can verify it:

$$\alpha = P_{H_0}\left( p(x) < \alpha \right)$$

$$= P_{H_0}\left( 2\left[ 1 - \Phi\left( |\bar{Y}_n - \bar{X}_n| \sqrt{\frac{n}{2}} \right) \right] < \alpha \right)$$

$$= P_{H_0}\left( \Phi\left( |\bar{Y}_n - \bar{X}_n| \sqrt{\frac{n}{2}} \right) > 1 - \frac{\alpha}{2} \right) =$$

this is the valid p-value for test A
$$P_{H_0}(1 - \Phi(s) > \beta) = \beta$$

$$= P_{H_0}\left( (\bar{Y}_n - \bar{X}_n) \sqrt{\frac{n}{2}} > \Phi^{-1}\left( 1 - \frac{\alpha}{2} \right) \right) +$$

$$+ P_{H_0}\left( |\bar{X}_n - \bar{Y}_n| \sqrt{\frac{n}{2}} > \Phi^{-1}\left( 1 - \frac{\alpha}{2} \right) \right)$$

this is the valid p-value for test B
$$P(1 - \Phi(-s) > \beta) = \beta$$

I can substitute $\frac{\alpha}{2} = \beta$

$$= \beta + \beta = \alpha$$

This proves that $p(x) = 2\left[ 1 - \Phi\left( |\bar{Y}_n - \bar{X}_n| \sqrt{\frac{n}{2}} \right) \right]$

is a valid p-value.

c) $\Delta_\mu := (\mu_y - \mu_x) \sqrt{\frac{2}{n}} > 0$

See Exercise_1_C.r file in the submission folder to run the code.

# Exercise 2

See the file: Exercise_2.r

It prints all the required graphs

When we order our data the median is the point in the $(n//2)+1$ position.

Since we want to compute which is the maximum local sensitivity I want to achieve it I change K entries of my data.

This can be visualize as moving the middle value to 0 or to the max value. This means that the new middle point is one of the adjacent values.

The function get_indices provides the indices of the potential new middle points

For each potential new middle point, the local sensitivity is computed. Subsequently we store the maximum.

## Exercise 3

We have the population graph $G = (V, E)$

and we consider an estimation

$$\tau = \sum_{e \in E} y_e$$

Using the Horowitz Thompson approach:

$$\hat{\tilde{\tau}} = \sum_{e \in E^*} \frac{y_e}{\tilde{\pi}_e^{(1)}}$$

<span style="color:magenta">probability that the edge is sampled</span>

$$\tilde{\pi}_e^{(1)} = P(e \in E^*) \quad [1]$$

We already computed the vertex pair inclusion probability

$$\tilde{\pi}_{\{i,j\}}^{(2)} := P(\text{"vertex pair } \{i,j\} \text{ is sampled"})$$

for all the $\{i,j\} \in V^{(2)}$.

Consider $\hat{\tilde{\tau}} = \sum_{e \in E^*} \frac{y_e}{\tilde{\pi}_e^{(2)}}$ as the estimator for $\tau$

Prove that it is in general unbiased

**pf.** to prove that it is not unbiassed let's consider its value in expectation:

$$\mathbb{E}\left(\hat{\tau}\right) = \mathbb{E}\left(\sum_{e \in E^*}' \frac{Y_e}{\widetilde{\Pi}_e^{(2)}}\right) =$$

$$= \mathbb{E}\left(\sum_{e \in E}' \frac{Y_e}{\widetilde{\Pi}_e^{(2)}} \mathbb{1}_{E^*}\right) =$$

*I can consider the sum of all the $e \in E$ and add the indicator function*

$$\mathbb{1}_{E^*} = \begin{cases} 0 & e \notin E^* \\ 1 & \text{else} \end{cases}$$

*in expectation we have that this is equal to $\mathbb{P}(e \in E^*)$*

$$= \sum_{e \in E}' \frac{Y_e}{\widetilde{\Pi}_e^{(2)}} \mathbb{E}\left(\mathbb{1}_{E^*}\right)$$

*linearity*

$$= \sum_{e \in E}' \frac{Y_e}{\widetilde{\Pi}_e^{(1)}} \mathbb{P}\left(e \in E^*\right)$$

$$\overset{[1]}{=} \sum_{e \in E}' \frac{Y_e}{\widetilde{\Pi}_e^{(1)}} \widetilde{\Pi}_e^{(1)}$$

To have an unbiassed estimator we require:

$$\tau := \sum_{e \in E}' Y_e = \sum_{e \in E}' \frac{Y_e}{\widehat{\Pi}_e^{(2)}} \widetilde{\Pi}_e^{(1)} := \widetilde{\tau}$$

which is true only for $\boxed{\widetilde{\Pi}_e^{(1)} = \widehat{\Pi}_e^{(2)}}$

Let's consider an example where $\widetilde{\pi}_e^{(1)} \neq \widetilde{\pi}_e^{(2)}$.

I choose to consider the ==unlabeled star sampling==:

- $$\widetilde{\pi}_e^{(1)} = 1 - \binom{N_v - 2}{n}$$

- $$\widetilde{\pi}_e^{(2)} = \frac{n(n-1)}{N_v(N_v - 1)}$$

where $N_v$ is the number of vertex.

$$1 - \binom{N_v - 2}{n} \neq \frac{n(n-1)}{N_v(N_v - 1)}$$

We can see that, in a general setting:

$$\widetilde{\pi}_e^{(1)} \neq \widetilde{\pi}_e^{(2)}$$

▨