

LOGIC - HYPOTHESIS TESTING

Null hypothesis, H_0 - NOTHING GOING ON

Alternative " , H_A - SOMETHING EXTRA GOING ON

E.g.

$$H_0 : P(T) = \frac{1}{2} \quad (\text{FAIR COIN}) \quad T = \text{Tails}$$

$$H_A : P(T) \neq \quad (\text{UNFAIR COIN})$$

TEST STATISTIC -

Measures how far away the data are from the expected if H_0 were true

The most common test statistic is the z-statistic.

$$z = \frac{\text{observed} - \text{expected}}{SE}$$

E.g. COIN TOS

$$\text{expected} = 10 \cdot \frac{1}{2} = 5$$

$$SE = \sqrt{10} \cdot \sqrt{\frac{1}{2} \cdot \frac{1}{2}} = 1.58$$

$$z = \frac{7 - 5}{1.58} = 1.27$$

P-Values as Measures of Evidence

large values of $|Z|$ are evidence against H_0

The strength of the evidence is the

p-value (observed significance level)

So, values of p por debajo del

5% se suelen considerar "statistically significant"

$$Z = \frac{\text{observed} - \text{expected}}{SE}$$

- If p-value is larger than 5% we will not reject the null hypothesis.

T-Test

▷ We use it when the sample is small.

$$S = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2}$$

$$\bar{x} \pm t_{n-1} SE$$

Statistical Significant vs Importance

TWO SAMPLE Z-TEST

Rating:

$$P_1 \rightarrow n=1000 \text{ \& } 55\%$$

$$P_2 \rightarrow n=1500 \text{ \& } 58\%$$

$$H_0 \rightarrow P_1 = P_2 \text{ (nothing is really coming up)}$$

$$z = \frac{\text{observed difference} - \text{expected difference}}{\text{SE of difference}}$$

$$\hat{P}_1 = 55\% \quad P_1 - P_2 = 0$$

$$\hat{P}_2 = 58\%$$

$$z = \frac{(\hat{P}_2 - \hat{P}_1) - (P_2 - P_1)}{\sqrt{(\text{SE}(P_1))^2 + (\text{SE}(P_2))^2}} = \frac{(0.58 - 0.55) - 0}{0.0202} = 1.48$$

\swarrow M. nomos
Tobla

$$\sqrt{\sqrt{\frac{P_1(1-P_1)}{n_1}} + \sqrt{\frac{P_2(1-P_2)}{n_2}}}$$

p-value = 2 (7%)
14% > 5% so we cannot
reject null hypothesis

Also confidence interval $(\hat{P}_1 - \hat{P}_2) \pm z \text{SE}(\hat{P}_1 - \hat{P}_2)$

$$\text{in case of } 95\% \quad z=2 \quad [-1\%, 7\%]$$

$$\left. \begin{array}{l} 0.33 \cdot 1000 = 330 \\ 0.58 \cdot 1500 = 870 \end{array} \right\} \rightarrow 1420 \text{ out of } 2500$$

proportion of $\frac{1420}{2500} = 0.568$

$$SE(\hat{p}_2 - \hat{p}_1) = \sqrt{\frac{0.568 \cdot (1 - 0.568)}{1000} + \frac{0.568(1 - 0.568)}{1500}} = 0.02022$$

which gives the se given

TO COMPARE TWO MEANS

$$SE(\bar{x}_2 - \bar{x}_1) = \sqrt{(SE(\bar{x}_1))^2 + (SE(\bar{x}_2))^2}$$

$$SE(\bar{x}_2) = \frac{\sigma_2}{\sqrt{n_2}} \quad \text{estimated} \quad \frac{s_2}{\sqrt{n_2}} \quad t\text{-test}$$

Paired Difference Test

paired t-test

$$t = \frac{\bar{d} - 0}{SE(\bar{d})}$$

$$SE(\bar{d}) = \frac{\sigma}{\sqrt{n}} \rightarrow SD = 0.55$$

$$t = \frac{1.4 - 0}{0.55/\sqrt{5}} = 5.69$$

L

COMPUTER SIMULATIONS IN

PLACE OF CALCULATIONS

Recall - confidence interval

$$\boxed{\bar{x} \pm z SE(\bar{x})}$$

If we instead of \bar{x} are interested in other estimator: for example $\hat{\theta}$ for a parameter θ and the normal approximation is not valid for that estimator $\hat{\theta}$

En estas situaciones las simulaciones pueden ser usadas para estimar

LAW OF LARGE NUMBERS

TO APPROXIMATE QUANTITIES OF INTEREST

Monte Carlo Method

For the explanation we will use the average height of people living in USA

Sample $n = 100$

We are interested in a parameter θ of a population that we estimate with $\hat{\theta}$

Our statistic (estimator) $\hat{\theta}$ is the average so

$$\hat{\theta} = \text{average of sample} = \frac{1}{n} \sum_{i=1}^n x_i$$

Monte Carlo Method or Simulation:

Approximation to a fixed quantity θ by the average of independent random variables with expected value θ . The larger the sample the smaller the SE of the statistic

$$SE(\hat{\theta}) = \sqrt{E(\hat{\theta} - E(\theta))^2} \quad E \rightarrow \text{variance}$$

Example:

1. 1000 samples of 100 observations
2. Compute $\hat{\theta}$ for the samples $\hat{\theta}_1 \dots \hat{\theta}_{1000}$
3. Compute Standard Dev

$$s(\hat{\theta}_1 \dots \hat{\theta}_{1000}) = \sqrt{\frac{1}{1000-1} \sum_{i=1}^{1000} (\hat{\theta}_i - \text{average}(\hat{\theta}_i))^2}$$

$$\bar{\hat{\theta}} = \text{average}(\hat{\theta})$$

Summary

- o Random Sampling
- o The more samples the less SE
- o We can compute the SE with it

Plug-In / Bootstrap principle

The bootstrap principle uses the plug-in principle and the monte carlo method to approximate quantities such as $SE(\hat{\theta})$

1. Draw sample $X_1 \dots X_n$ to compute $\hat{\theta}$
2. Repeat B times (e.g. $B=1000$) to get $\hat{\theta}_1 \dots \hat{\theta}_B$
3. If we have only 1 sample then the bootstrap simulates from the sample since we don't have access to the population.

Basically it interacts with the sample like "if the sample was the population"

Non-parametric bootstrap

Sometimes we can know or suppose characteristics of the data. For example we might know that it follows a normal distribution but not its SE or mean

BOOTSTRAP CONFIDENCE INTERVALS

$$\hat{\theta} \pm z_{\alpha/2} SE(\hat{\theta})$$

We can estimate the sampling distribution of $\hat{\theta}$ not just $SE(\hat{\theta})$

Making a histogram of the bootstrap copies

Also we can do $\hat{\theta} - \theta$

bootstrap pivotal interval

$$[2\hat{\theta} - \hat{\theta}_{(1-\alpha/2)}^*, 2\hat{\theta} - \hat{\theta}_{(\alpha/2)}^*]$$

Bootstrapping in regression

We have data $(X_1, Y_1), \dots, (X_n, Y_n)$ from simple linear regression model

$$Y_i = a + bX_i + e_i$$

From the data we can compute estimates \hat{a}, \hat{b} .

And Compute residuals:

$$\tilde{e}_i = Y_i - \hat{a} - \hat{b} X_i$$

* Remember: residuals are the difference between the observed values and the values predicted by the model/regression line

$$\text{Residual} = \text{observed} - \text{predicted}$$

Steps

1. Compute residuals $\tilde{e}_i = Y_i - \hat{a} - \hat{b} X_i$
2. Resample from those residuals to get $e_1^* \dots e_n^*$
3. Compute the bootstrap responses $Y_i^* = \hat{a} + \hat{b} X_i + e_i^*$