

Assignment 14

G HARSHA VARDHAN REDDY (CS21BTECH11017)

June 15, 2022

AI1110

Outline

1 Problem Statement

2 Definitions

3 Solution

Problem Statement

Papoulis Pillai Probability Random Variables and Stochastic Processes Exercise : 8-13

We plan a poll for the purpose of estimating the probability p of Republicans in a community. We wish our estimate to be within ± 0.02 of p . How large should our sample be if the confidence coefficient of the estimate is 0.95?

Definitions

Sample Proportion

If X is a binomial random variable, then $X \sim B(n, p)$ where n is the number of trials and p is the probability of a success. To form a sample proportion, take X , the random variable for the number of successes and divide it by n , the number of trials (or the sample size). The random variable P' is the sample proportion

$$P' = \frac{X}{n} \quad (1)$$

And
 p' = the estimated proportion of successes or point estimate for p

Binomial distribution

Let the random variable X be the sum of n Bernoulli random variables, i.e.

$$X = X_1 + X_2 + \dots + X_n \quad (2)$$

Assume $E(X_i) = p$. Then,

$$E(X) = E(X_1 + X_2 + \dots + X_n) \quad (3)$$

$$= E(X_1) + E(X_2) + \dots + E(X_n) \quad (4)$$

$$= p + p + \dots + p \quad (5)$$

$$= np \quad (6)$$

Binomial Distribution

The variance is given by:

$$E((X - np)^2) = E(X^2 - 2Xnp + n^2p^2) \quad (7)$$

$$= E(X^2) - 2npE(X) + E(n^2p^2) \quad (8)$$

$$= E(X^2) - n^2p^2 \quad (9)$$

Now,

$$E(X^2) = \sum_{k=0}^n k^2 \times \binom{n}{k} p^k q^{n-k} \quad (10)$$

$$= \sum_{k=1}^n npk \times \binom{n-1}{k-1} p^{k-1} q^{n-k} \quad (11)$$

Binomial Distribution

$$\begin{aligned} \sum_{k=1}^n npk \times \binom{n-1}{k-1} p^{k-1} q^{n-k} &= n(n-1)p^2 \times \sum_{k=2}^n \binom{n-2}{k-2} p^{k-2} q^{n-k} \\ &+ np \times \sum_{k=1}^n \binom{n-1}{k-1} p^{k-1} q^{n-k} \end{aligned} \quad (12)$$

$$= n(n-1)p^2(p+q)^{n-2} + np(p+q)^{n-1} \quad (13)$$

$$= n(n-1)p^2 + np \quad (14)$$

From 9,

$$\sigma^2 = (n(n-1)p^2 + np) - n^2p^2 \quad (15)$$

$$= np - np^2 = np(1-p) \quad (16)$$

$$\implies \sigma^2 = npq \quad (17)$$

Let the random variable $Y = \frac{X}{n}$. Then, Y is a binomial random variable that maps to $\frac{k}{n}$ when X maps to k . Therefore,

$$\mu_Y = \frac{\mu_X}{n} \quad (18)$$

$$= p \quad (19)$$

$$\sigma_Y^2 = \frac{\sigma_X^2}{n^2} \quad (20)$$

$$= \frac{pq}{n} \quad (21)$$

$$\implies \sigma_Y = \sqrt{\frac{pq}{n}} \quad (22)$$

Youth percentile or Z score

Z-score

Z-score indicates how much a given value differs from the standard deviation. The Z-score, or standard score, is the number of standard deviations a given data point lies above or below mean.

$$\implies Z_u = \frac{x - \mu}{\sigma} = \frac{p - p'}{\sigma_{p'}} \quad (23)$$

Where,

Z_u = Normal (Youth) percentile or Z score

x = Observed value

σ = Standard deviation

Confidence interval for a population proportion

A **Confidence Interval** is an estimate for an unknown parameter. It is governed by a number $\gamma = 1 - \delta$, which determines the accuracy of the estimation method. γ is called the **confidence coefficient**.

The confidence interval for a population proportion (p)

$$|p - p'| \leq \sigma Z_u \quad (24)$$

$$p' - \sigma Z_u \leq p \leq p' + \sigma Z_u \quad (25)$$

From (22)

$$\sigma_{p'} = \sqrt{\frac{(1 - p')(p')}{n}} \quad (26)$$

Therefore,

$$p' - Z_u \times \sqrt{\frac{(1 - p')(p')}{n}} \leq p \leq p' + Z_u \times \sqrt{\frac{(1 - p')(p')}{n}} \quad (27)$$

Solution

Given,

$$p \leq \pm 0.02 \text{ of } p' \quad (28)$$

$$\implies |p - p'| \leq 0.02 \text{ and,} \quad (29)$$

$$\text{Confidence Coefficient (CF)} = 0.95 \implies Z_u = 2 \quad (30)$$

From (27),(29) and (30),

$$\sqrt{\frac{(1 - p')(p')}{n}} \times 2 \leq 0.02 \quad (31)$$

As $n > 0$, From (31)

$$\frac{(1 - p')(p')}{n} \leq \left(\frac{1}{100}\right)^2 \quad (32)$$

$$\implies n \geq (1 - p')(p') \times 100^2 \quad (33)$$

$$(1 - p')(p') = p' - p'^2 \quad (34)$$

$$= -(p'^2 - p') = -\left(\left(p' - \frac{1}{2}\right)^2 - \frac{1}{4}\right) \quad (35)$$

$$= \frac{1}{4} - \left(p' - \frac{1}{2}\right)^2 \quad (36)$$

$$\implies (1 - p')(p') \leq \frac{1}{4} \quad (37)$$

From (33),(37)

$$n \geq \frac{100^2}{4} \quad (38)$$

$$\implies n \geq 2500 \quad (39)$$

Therefore, the size of sample(n) must be greater than equal to 2500.