Term Test 1: Practice Problems: Solutions

- Problem 1. We consider a population of N units with bernoulli random variables $\{0,1\}$ data values. Suppose we choose a simple random sample without replacement of n units from this population. Let $p = \frac{1}{N} \sum_{i=1}^{N} y_i$ be the population proportion, and $\sigma^2 = \frac{1}{N} \sum_{i=1}^{N} (y_i p)^2$ denotes the population variance. Let $\widehat{p} = \frac{1}{n} \sum_{i=1}^{n} y_i$ and $s^2 = \frac{1}{n-1} \sum_{i=1}^{n} (y_i \widehat{p})^2$ be the sample proportion and the sample variance, respectively.
 - (a) Is \hat{p} unbiased estimator to p? Justify. $E(\hat{p}) = E(\frac{1}{n} \geq \hat{z}_{i}) = \frac{1}{n} \geq E(\hat{z}_{i}) = \frac{1}{n} \geq p = \frac{1}{n} (np) = p$ Yes \hat{p} is unbiased for \hat{p} .
 - (b) Show that σ^2 can be writen as $\sigma^2 = p(1-p)$. $S^2 = \frac{1}{N} \left[\sum_{i=1}^{N} y_i 2p \sum_{i=1}^{N} y_i + \sum_{i=1}^{N} y_i 2p \sum_{i=1}^{N} y_i 2p \sum_{i=1}^{N} y_i + \sum_{i=1}^{N} y_i + \sum_{i=1}^{N} y_i 2p \sum_{i=1}^{N} y_i + \sum_{i=1}^{N} y_i 2p \sum_{i=1}^{N} y_i + \sum_{i=1}^{N} y_i + \sum_{i=1}^{N} y_i 2p \sum_{i=1}^{N} y_i + \sum_{i=1}^{N} y_i 2p \sum_{i=1}^{N} y_i + \sum_{i=1}^{N} y_i +$
 - (c) Show that s^2 can be writen as $s^2 = \frac{n}{n-1}\widehat{p}(1-\widehat{p})$. $S = \frac{1}{n-1}\left[\sum_{n=1}^{\infty} y_n - 2\widehat{p}\sum_{n=1}^{\infty} y_n + \sum_{n=1}^{\infty} \widehat{p}(1-\widehat{p}) - 2n\widehat{p} + n\widehat{p}^2\right]$ $= \frac{1}{n-1}n(\widehat{p}(1-\widehat{p})) = \frac{n}{n-1}\widehat{p}(1-\widehat{p})$
 - (d) Is the sample variance s^2 unbiased estimator for σ^2 ? Find $E(s^2)$ Under SRS without replacement, we have seen an dominativate in class. that when estimating μ , the sample variance s^2 is becased Estimator for Γ^2 and $E[s^2] = \frac{1}{N-1}s^2$, which is a
- (e) Compute $V(\widehat{p})$. Given the result in (d), find an unbiased estimator for $V(\widehat{p})$ in term of \widehat{p} .

From the estimation of μ , we have demonstrated that using SR5 $V(\bar{y}) = \frac{N-n}{N-1} \frac{S^2}{n}$. For $\mu=p$, $V(\bar{p}) = \frac{N-n}{N-1} \frac{S^2}{n}$. Rep,

- Based on (d), an unbiased estimator for 6^2 is $\frac{N-1}{N}$ 6^2 thus, unbiased Estimator for $V(\hat{p})$ is $\frac{N-n}{N-1} \cdot \frac{N-1}{N} \cdot \frac{5^2}{N} = \frac{N-n}{N} \cdot \frac{5^2}{N} \cdot \frac{N-n}{N} \cdot \frac{$
- Using $S = \frac{u}{u-1} \hat{p}(u-\hat{p})$, => & become $\sum_{N=1}^{N-1} \frac{1}{N} \frac{u}{u-1} \hat{p}(u-\hat{p})$ Estimat $d_{V}(\hat{p}) = V(\hat{p}) = (u-\frac{n}{N}) \frac{\hat{p}(u-\hat{p})}{u-1} V(\hat{p}) = \frac{N-u}{N} \frac{1}{n} \frac{u}{u-1} \hat{p}(u-\hat{p})$

- Problem 2. To estimate the proportion of voters in favor of a controversial proposition, a simple random sample of XXXXX eligible voters was contacted and questioned. Of these, 552 reported that they favored the proposition. The study also reports a margin of error \pm 3%, 19 out of 20. The number of eligible voters in the population is approximately 1,800,000.
 - (a) What "a margin of error \pm 3%, 19 out of 20" means? \pm 2 means 95% Margin of Error is 0.03, e.g. For at least 95% of all random samples, the Bound on error (the amount by which the sample proportion \hat{p} is expected to differ from the true popular proportion \hat{p} is expected to differ from the true popular proportion \hat{p} is 0.03.

(b) Complet the study by computing the missing sample size.

The property of the study by computing the missing sample size.

The property of the property of

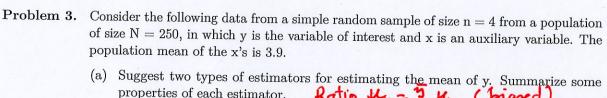
 $\phi = \frac{552}{1111} = 0.497$

(d) Give a 95% confidence interval for population proportion.

a 95% C.I for p is $(\hat{p}-ME, \hat{p}+ME)$ (0.497-0.03, 0.497+0.03)

(0.467, 0.527)

Interprétation: we are 95% confidence that the true prop b lies between 0.467 and 20.527



properties of each estimator. Ratio H = 1 Kx (biased) Simple: y= 1 Zy Regregaion firen = ig+b(Kz-zi)

(b) The data obtained were recorded in variables y and x, analyzed in R and produced results presented below (see next page). Based on the R output, answer the following questions.

(i) Estimate mean of y using the simple estimator, and estimate the variance of the estimator. $\frac{590}{4} = 147.5$

V(4)=(1-1) 54 - 416.501

(ii) Estimate mean of y using the ratio estimator, and estimate the variance of the estimator.

for = 177 = = - 177

(iii) Estimate mean of y using the regression estimator, and estimate the variance of the estimator.

line 24 $\hat{\mu}_{reg} = \bar{\eta} + b(\hat{\mu}_{z} - \bar{z}) = 171.72$ line 26 $V(\hat{\mu}_{reg}) = (1 - \frac{n}{2}) \frac{MS\bar{t}}{M} = \sqrt{154.3}$ = (1-1) MSE = 154.3

(iv) Estimate mean of y using the difference estimator, and estimate the variance of the estimator.

 $= \sqrt{4} + (42 - 2) = \frac{590}{4} + (3.9 - \frac{15}{4}) = 590.65$ $V_{i}(\hat{\mu}_{0}) = (1 - \frac{n}{N}) \frac{1}{n} = \frac{\sum (gl_{i} - \overline{d})}{N-1}$ (not provided in the Routput)

Based on the data, which estimator appears preferable in this situation?

We would prefere Rotor estimator which provides small variance.

R output:

```
> # Population
       > N<-250
      > mu_x<-3.9
       > # SRS of size n=4
       > n<-4
       > y<-c(150, 100, 200, 140)
       > x<-c(4,2,4,3)
       > ysum<-sum(y); ysum
       [1] 590
       > xsum<-sum(x); xsum
       [1] 13
       > s2_y<-var(y); s2_y
       [1] 1691.667
       > r<-(ysum/xsum); r
       [1] 45.38462
       > (ysum/xsum)*mu_x
       [1] 177
       > s2_r<-var(y-r*x); s2_r
       [1] 478.501
(4 > (1 - n/N)*(s2_y/n)
      [1] 416.15
15 > (1 - n/N)*(s2_r/n)
       [1] 117.7112
       > (1 - n/N)*(1/mu_x)^2*(s2_r/n)
       [1] 7.739069
       > (1 - n/N)*(1/mu_x)*(s2_r/n)
       [1] 30.18237
       > cor(x,y)
       [1] 0.8676399
       > # regression of y on x
       > fitreg<-lm(y~x)
       > coef(fitreg)
       (Intercept)
          26.36364
                     37.27273
       > yhat<-fitted(fitreg)</pre>
       > ehat<-y-yhat
24 > mean(y) + coef(fitreg)[2]*(mu_x-mean(x))
       171.7273
      > MSE<-sum( ehat^2 )/(n-2)
26 > (1 - n/N)*(MSE/n)
       [1] 154.3091
      > (1 - n/N)*(MSE/(n-1))
       [1] 205.7455
      > (1 - n/N)*(MSE/(n-2))
       [1] 308.6182
```