# Hasselt University

## Master of Statistics

### Computer Intensive Methods

---

## **Assignment III**

---

*Students:*
Christine Jani (1747124)
Joe Gwatsvaira (1747457)
Bache Emmanuel Bache
(1747695)
Adama Kazienga (1747603)

*Lecturers:*
Prof.Ziv Shkedy

January 20, 2019

# Question 1

(a) **Use the Likelihood ratio test (LRT) to test the null hypothesis**

A likelihood ratio test (LR test) is a statistical test used for comparing the goodness of fit of two statistical models a null model against an alternative model. The test is based on the likelihood ratio, which expresses how many times more likely the data are under one model than the other. Interest is in testing the hypotheses $H_0 : \beta_1 = \beta_2 = 0$. This hypothesis states that the real per capita disposable income (dpi) does not depend on the $X_1$ (personal savings (sr)) and $X_2$ (percentage of the population under 15 (pop15)) . We build a likelihood ratio test for this hypothesis: Two nested models are constructed and the model under the null hypothesis is given by:

$$H_0 : dpi = \beta_0 + \beta_1 pop75 + \epsilon_i$$

and a larger model under the alternative hypothesis given by;

$$H_1 : dpi = \beta_0 + \beta_1 sr + \beta_1 pop15 + \beta_1 pop75 + \epsilon_i$$

The test statistic for the hypothesis is obtained by comparing the maximized log likelihoods for the two models i.e $\chi^2_{cal} = -2L_0-(-2L_1)$ where $L_0$ is the log likelihood of the null model and $L_1$ is the log likelihood of the model under the alternative hypothesis. Table 1 show the results of the LRT, $p = 0.34037$ we fail to reject the null hypothesis and conclude that there is insufficient evidence at 5% level of significance to reject the reduced model.

| Results | Df | log likelihood | $\chi^2_{cal}$ | pvalue |
|---|---|---|---|---|
| Saturated Model | 4 | 390.1452 | 2.15547 | 0.34037 |
| Reduced Model | 2 | 390.2229 | | |

Table 1: Results of LRT

(b) **Use non parametric bootstrap to test the null hypothesis in (a) with the LRT**

Under non parametric bootstrap no distributional assumptions are made on the error terms $\epsilon_i$. Under the null hypothesis we have that $\beta_1 = \beta_2 = 0$ thus there is no linear relationship between dpi and (sr, pop15). The model under $H_0$ is

$$dpi = \beta_0 + \beta_1 pop75 + \epsilon_i$$

Bootstrapping is done obeying the null hypothesis i.e we preserve the relationship between dpi and pop75 and break the relationship between dpi and (sr, pop15). A random bootstrap sample $Z^* = (dpi^*, pop75^*)$ of size 50 from $Z = (dpi, pop75)$ with replacement while fixing (sr ,pop15). B=1000 bootstrap resample are obtained and at each time fitting the reduced and larger(saturated )model. The Test statistic: $LRT = -2(L_0 - L_1)$

1

where $L_0$ is the log likelihood of the restricted model and $L_1$ is the full model loglikelihood, which is calculated for each bootstrap sample.

The bootstrap P-value is calculated as

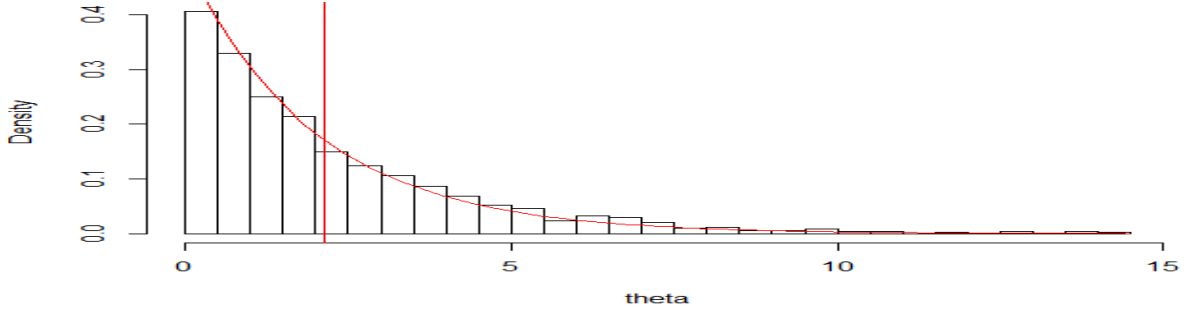$$P_{ASL} = \frac{1 + sum\ (LRT \geq LRT.observed)}{(B + 1)}$$



Figure 1: Histogram bootstrap replicates

As result,

1. As a conclusion, we failed to reject the null hypothesis ($P_{ASL} = 0.379$) hence we conclude that we have insufficient evidence at 5% significance level to say that at least 1 of the coefficients $\beta_1$ or $\beta_0$ is different from zero.

2. It is known from the classical theory that the test statistic LRT follow a $\chi^2$ distribution under the null hypothesis. This is shown by the right skewed histogram of the bootstap LRT statistics in Figure 1 overlayed with a curve of $\chi_2^2$

| Parameter Estimates | | | |
| --- | --- | --- | --- |
| Coefficients | Estimate | Std Error | Pvalue |
| Intercept | -278.55 | 179.44 | 0.127 |
| Pop75 | 604.15 | 68.35 | <0.0001 |

Table 2: Linear regression estimates

The estimated regression model is

$$\hat{dpi}_i = -278.55 + 604.15 pop75_i$$

The intercept does not make sense to be interpreted in this case. A percent increase in population over 75 years will increase the real per capita disposable income by 604.15.

(c) **The predictive model for a new country with** $sr = 7.5, pop15 = 32, pop75 = 2.51$

Predictive modelling is a process through which future values of an outcome are predicted based on the current data at hand. The first step is to see if there is a relationship between dpi and the explanatory variable, and also there is need to check if the model is a good fit for our data. This was done in part (a) and (b). We consider the model $dpi = \beta_0 + \beta_1 pop75 + \epsilon_i$ and hence make predictions. We estimate the real per capita disposable income for the new country by substituting $pop75 = 2.51$ in $\hat{dpi}_i = -278.55 + 604 pop75_i$ we get $\hat{dpi}_i = 1237.858$ with a confidence interval $[1059.715; 1416.001]$ and standard error $s.e = 88.600$. To check the distribution of the bootstrapped predicted values of $dpi_i$, B=1000 non parametric bootstraps samples with replacement(we re-sample rows from the original sample), were sampled. At each bootsrap we run the regression model $dpi = \beta_0 + \beta_1 pop75 + \epsilon_i$ and predict the $dpi$ for a country with $pop75 = 2.51$. The distribution of the bootstrapped predicted values is shown in Figure 2 it is symmetric around the predicted value from the original sample. The percentile confidence interval is given by $[1174.145; 1302.538]$ and a standard error $S.e = 96.814$ and the average of the predicted values is 1239.452.
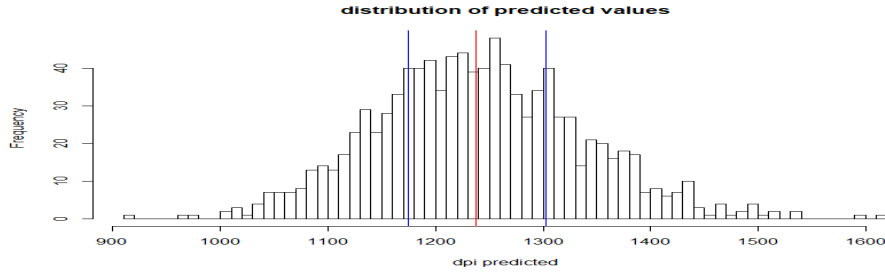


Figure 2: Distribution of predicted values

# Question 2

(a) **Formulate an appropriate model for the number of satellites. Fit the model and test the null hypothesis** $H_0 : \beta_2 = 0$ **against a two side alternative**

A Poisson model was fitted using a natural log link for the number of satellites $Y_i \sim Poisson(\mu_i)$. The model was formulated as follows:

$$log(\mu_i) = \beta_0 + \beta_1 \times Width + \beta_2 \times Dark$$

Moreover, an hypothesis test as follows: $H_0 : \beta_2 = 0$ vs $H_1 : \beta_2 \neq 0$. The LRT was used with saturated with

$$log(\mu_i) = \beta_0 + \beta_1 \times Width + \beta_2 \times Dark$$

3

as model whereas the model used in the reduced model is

$$log(\mu_i) = \beta_0 + \beta_1 \cdot Width$$

The Table 3 displays the hypothesis testing findings and suggests that there was enough evidence to reject the null hypothesis.

| Results | Df | log likelihood | $\chi^2_{cal}$ | pvalue |
|---|---|---|---|---|
| Saturated Model | 3 | 280.4788 | 6.92101 | 0.00851 |
| Reduced Model | 2 | 283.4788 | | |

Table 3: results of LRT

(b) **Use parametric and non parametric bootstrap to test the null hypothesis in (a)**

**Non- Parametric**

In non parametric bootstrap, a random sample with replacement was drawn under the null hypothesis, ie keeping the Dark variables fixed and sampled only the satellites and Width variables. Then, a saturated and reduced was fitted as in (a) and the statistic of interest was estimated.

**Parametric**

Assumption was made in the case of the parametric bootstrap that the distribution come from a Poisson distribution, and random sample was drawn accordingly using the appropriate parameters (size and the mean). As in the non parametric case, the resampling techniques was made under the null hypothesis and the statistic of interest estimated later on. We simulate data that is Poisson($\mu$) with $\mu$ given by the predictions in the small model i.e $log(\mu_i) = \beta_0 + \beta_1 \times Width$ that is $Y_i \sim Po(\mu_i)$ (Geyer,Shaw and Wagenious 2003). These bootstrapped Y values are then used to fit the full and reduced model.

A likelihood ratio test statistic is calculated at each bootstrap resample. i.e $LRT = -2(L_0 - L_1)$ where $L_0$ and $L_1$ are the loglikelihood values for the reduced and full models respectively. In both bootstrap methods, a total of B=1000 replicates was used. The LRT statistic values are then plotted and the Histograms are shown in Figure 3. The histograms show a right skewed distribution which confirms to the chi square distribution. The p values calculated for the classical, non parametric and parametric tests are shown in Table 4 we can see that the p values for the classical and the parametric bootstrap test are significant ( i.e the suggest the the null hypothesis is rejected and conclude that we have sufficient evidence at 5% level of significance that the saturated model fits the data well.) on the other hand the p value for the non parametric bootstrap test show an insignificant p value suggesting

4

that we fail to reject $H_0$. This is can be seen as a way of showing that a Poisson model does not work well when there is over dispersion in the data set or excess zeros. Though Poisson regression can produce consistent parameter estimates the standard errors may be incorrect. Thus it is recommended that one use models that take care of over dispersion or excess zeros (Agresti,2003).
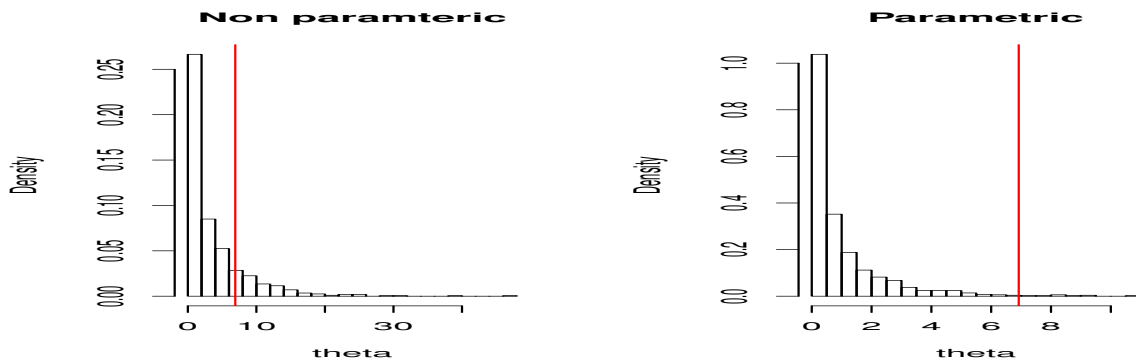


Figure 3: Histogram bootstrap replicates

| Test | Estimate(LRT) | P value |
|---|---|---|
| Classical | 6.921 | 0.0085 |
| Non Parametric | | 0.1608 |
| Parametric | | 0.009 |

Table 4: P values Likelihood ratio test

# References

[Agresti, 2003] Agresti, A. (2003). *Categorical data analysis*, volume 482. John Wiley & Sons.

[Charles J. Geyer, 2003] Charles J. Geyer, Ruth G. Shaw, S. W. (2003). A glm example. http://www.stat.umn.edu/geyer/5931/mle/seed2.pdf. [Online; accessed 15-January-2019].