

Count Regression Models with Application to Caries Experience for Children Attending Lady Northey Dental Clinic in Nairobi

Agnes Njambi Wanjau, Samuel Musili Mwalili

Department of Statistics and Actuarial Sciences, Jomo Kenyatta University of Agriculture and Technology, Nairobi, Kenya

Email address:

agiewanjau@gmail.com (A. N. Wanjau), samuel.mwalili@gmail.com (S. M. Mwalili)

To cite this article:

Agnes Njambi Wanjau, Samuel Musili Mwalili. Count Regression Models with Application to Caries Experience for Children Attending Lady Northey Dental Clinic in Nairobi. *American Journal of Theoretical and Applied Statistics*. Vol. 6, No. 4, 2017, pp. 176-181. doi: 10.11648/j.ajtas.20170604.12

Received: May 20, 2017; **Accepted:** May 31, 2017; **Published:** June 9, 2017

Abstract: Count regression models were developed to model data with integer outcome variables. These models can be employed to examine occurrence and frequency of occurrence. Four common types of count regression models are applied to caries data among children aged between three and six years attending Lady Northey Dental clinic between September, 2014 and November 2014. These models include Poisson, Negative Binomial (NB), Zero Inflated Poisson (ZIP) and Zero Inflated Negative Binomial (ZINB). The simplest count regression model, Poisson, was fitted first before considering other complex models. However, it did not perform better than its improved counterparts. The NB model proved to be the the simplest model that fits the data well according to Akaike Information Criterion (AIC), and was therefore employed to determine the important predictors of caries experience among the children. Model comparison was performed on the four models by use of AIC. Deviance values for various NB models were compared and the model with the least deviance value was considered to give a subset of best predictors of Early Childhood Caries (ECC). These predictors included age, gender, brushing frequency, feeding habit biscuits, feeding habit jam and highest education of the mother.

Keywords: Count Regression Models, Model Selection, AIC

1. Introduction

Dental caries, also known as tooth decay or cavity, is a bacterial infection that dissolves tooth enamel, the outer hard layer of a tooth. This dissolution is caused when acids from hydrolysis of food particles attack enamel surface and remove minerals from it (demineralisation). Alkali such as sodium bicarbonate present in saliva help in remineralisation, thus reversing the process. Predominance of demineralisation leads to dental caries. Enamel tissue breaks down progressively, leading to holes or cavities in the tooth (dental caries). The decay spreads eventually into the dentine and the pulp.

The first dental caries statistics were published at around 1900. They were however very low and therefore hard to interpret. A survey was prepared later, between 1950 and 1963 by the International Dental Association which revealed that several studies were done at that time. The dmft index was used to quantify the dental health status among children

with primary dentition.

Until sixties, the researchers drew samples either in areas where incidences were expected or in towns near dental schools [6]. Sampling methods were not usually specified.

Two articles on rates of caries in 14 sub-Saharan African nations, exempting South Africa, between 1945 and 1989 reported that dental caries is fairly stable and at a low level. However, these two articles did not consider the varying diagnostic methods used in the studies and the range of age groups was wide, i.e., from 6 to 20 years.

Tooth decay is a very common oral disease among adults as well as among children. It is still a global problem due to sucrose and sugars present in the diet [10]. The World Health Organization (WHO) reports that caries experience affects 60% to 90% of school children and a huge number of adults. It's also regarded as a highly prevalent oral disease in many countries of

Latin American and Asia. Previous studies portray an increase in dental caries prevalence as well as incidence in developing countries. This is mainly attributed to limited exposure to fluorides and increasing amounts of sugar in the diet. Recent reviews revealed that use of fluoride toothpastes and fluoridated water while rinsing the mouth reduces the prevalence of dental caries significantly.

Children are not spared the agony of dental caries during eruption of deciduous teeth. Early Childhood Caries (ECC) is an oral disease that starts when a child's teeth emerge. This decay spreads rapidly to other tooth surfaces that are not affected.

ECC is termed present if one or more teeth are decayed, filled or missing due to caries among children aged six years or below. This oral disease has as well been referred to as labial caries, nursing bottle syndrome, nursing caries, milk bottle caries, bottle rot and baby bottle tooth decay. The two main types of ECC are rampant caries and nursing caries, with the engagement of mandibular incisors in the former as opposed to the latter. Among others, factors contributing to ECC are prolonged use of sweetened foods and night-time prolonged bottle feeding of milk or juice [2]. The quality of life of young children has been observed to be affected by ECC as it (ECC) contributes to other health problems. According to [2], children with ECC are seen to weigh less than 80% of their expected weight. Prevalence rate of ECC in industrialized countries has been recorded to range between 1% and 12% while in developing countries it is reported to be 70%.

Nairobi County is one of the 47 counties in Kenya, consisting of 17 constituencies. After devolution, health care was devolved to the County Governments. Statistics reveal that dental caries prevalence rate in this county is 93.6% among 3-5 year old children and 89.4% among those above the range. [8] reported reduced prevalence rates of 63.5% among 3-5 year old children in Nairobi while [9] reported 59.5% among the same group in Kiambaa constituency, Kiambu County. Mean dmft scores of 1.27 and 2.95 among preschool children in Nairobi have been reported at different studies.

Lady Northey Dental Clinic is a public clinic in Nairobi County. It is concerned with children's oral health problems.

2. Methods

2.1. Data

The data were collected at Lady Northey Dental Clinic, a government run institution, among patients aged 3-6 years from September to November, 2014. This clinic is located in Westlands constituency, Nairobi County along State house Avenue. The study involved patients in full primary dentition, hence the age group. The sampling frame included only those children aged between 3 and 6 years, that is, children in full primary dentition and were accompanied by parents who accepted to take part in the study. A questionnaire was employed for data collection and a face-to-face interview with the parent/guardian. Interview was conducted by a trained research assistant. Only observations with all values for every variable were used.

2.2. Study Variables

The outcome variable is the number of teeth decayed, missing or filled due to dental caries (dmft). Predictor variables included demographic, dietary practices and oral hygiene practices, which are a mixture of continuous and categorical variables. The covariates considered for analysis included the following: Age(x1), Gender(x2), Highest education of the father(x3), Highest education of the mother(x4), Employment state of the father(x5), Feeding habit biscuits(x6), Feeding habit gum(x7), Feeding habit jam(x8), Feeding habit juice(x9), Feeding habit soda(x10), Feeding habit sweets(x11), Feeding habit tea with sugar(x12), Brushing frequency(x13), Use of flouridated toothpaste(x14).

2.3. Count Data Regression Models

2.3.1. Introduction

There exists a wide range of models that analyze count data. The most popular model for count data is the Poisson model, which is based on the property that the mean and variance of the dependent variable are equal. However, this is not always the case, as the variance sometimes exceeds the mean. This is referred to as overdispersion.

Overdispersion can be modelled using negative binomial (NB) regression model, but more models accounting for overdispersion exist. The negative binomial regression model assumes a gamma distribution for the Poisson mean with variation over the subjects [7]. Further, the response variable can be observed to show excess zero counts, contrary to what is expected, on the basis of Poisson or negative binomial distribution. According to [13], this is an implication that the count data are zero inflated. Zero-inflated models allow for overdispersion as well as modelling zero-inflated count data. The frequently used models for zero inflated count data are zero-inflated Poisson (ZIP) and zero-inflated negative binomial (ZINB).

2.3.2. Poisson Model

Let y_i be the number of dmft count.

$$y_i \sim \text{Poisson}(\mu_i) \quad (1)$$

$$p(Y = y) = \frac{e^{-\mu} \mu^y}{y!}, y = 0, 1, 2, \dots \quad (2)$$

$$E(Y) = \text{Var}(Y) = \mu \quad (3)$$

Where $E(Y)$ is the mean of Y and $\text{Var}(Y)$ is the variance of Y .

With reference to [1], Poisson regression model is a member of generalized linear models, an extension of the ordinary linear model that allows the mean of the output variable to depend on a linear predictor through a log link function.

2.3.3. Assumptions of Poisson Regression Model

- The response variable y has a Poisson distribution with mean μ and variance μ .

- b. The logarithm of the expected value of Y can be modelled by a linear combination of unknown parameter.

2.3.4. Negative Binomial Model (NB)

Negative binomial model has been employed by researchers to relax the restriction of the Poisson model, that the variance of the random variable equals the mean [4]. The NB distribution assumes that the means follow Gamma distribution.

$$\Pr(Y = y_i) = \frac{r(\alpha + y_i) \left(\frac{\alpha}{\alpha + \mu}\right)^\alpha \left(\frac{\mu}{\alpha + \mu}\right)^{y_i}}{y_i! r(\alpha)} \quad (4)$$

Where $y_i \geq 0$ and $\mu > 0$.

The shape parameter α measures the amount of overdispersion and the conditional mean is maintained.

$$E(Y) = \mu \quad (5)$$

$$\text{var}(Y) = \mu + \frac{\mu^2}{\alpha} \quad (6)$$

The mean response related to a vector of covariates through a log-linear model builds a NB regression model.

2.3.5. Zero Inflated Models

The concept of zero-inflated data led to derivation of zero-inflated distributions. Zero-inflated models have been of importance in describing data characterized by dominance of

zeros. In such a case, more zeros are observed than what Poisson, NB or Poisson-Gamma processes would predict. Based on the work of [5], zero-inflated models have a wide application in dental epidemiology. Two sources of data with excess zeros are believed to exist, thus the process has been referred to as dual state. One state is a normal count process while the other is a zero count state. This system assumes that the excess zeros are due to heterogeneity in the data and every observation has the same mean μ . Examples of regression models proposed in literature to model zero-inflated data are zero-inflated Poisson (ZIP) and zero-inflated negative binomial (ZINB).

Zero inflated distribution with a proportion p of zeros is given as follows;

Let $f(y_i|\emptyset)$ be a distribution function for count data, e.g., the Poisson and NB.

$$Y_i = \begin{cases} 0, & \text{with probability } p \\ f(Y_i = y_i|\emptyset), & \text{with probability } (1 - p) \end{cases} \quad (7)$$

Where \emptyset is a vector of unknown parameters.

Zero inflated distribution, $ZIF(y_i|\emptyset)$, is given by;

$$\Pr(Y_i = y_i|\emptyset) = \begin{cases} p + (1 - p)f(y_i = 0|\emptyset) \\ (1 - p)f(Y_i = y_i|\emptyset), y_i = 1, 2, \dots \end{cases} \quad (8)$$

The mean and variance of zero inflated distribution are given by:

$$E_{zif}(Y_i|p, \emptyset) = (1 - p)E_f(Y_i|\emptyset) \quad (9)$$

$$\begin{aligned} \text{Var}_{zif}(Y_i|p, \emptyset) &= (1 - p)\{E_f(Y_i^2|\emptyset)\} - \{(1 - p)E_f(Y_i|\emptyset)\}^2 = (1 - p)\{E_f(Y_i^2|\emptyset) - [E_f(Y_i|\emptyset)]^2 + p[E_f(Y_i|\emptyset)]^2\} \\ &= (1 - p)\{\text{Var}_f(Y_i|\emptyset) + p[E_f(Y_i|\emptyset)]^2\} \end{aligned} \quad (10)$$

When $p = 0$, $E_{zif}(Y_i|\emptyset) = E_f(Y_i|\emptyset)$ and $\text{Var}_{zif}(Y_i|\emptyset) = \text{Var}_f(Y_i|\emptyset)$

2.3.6. Zero Inflated Poisson

Zero-inflated Poisson distribution is given by:

$$f(y_i) = \begin{cases} p + (1 - p), & y_i = 0 \\ (1 - p) \frac{e^{-\mu} \mu^{y_i}}{y_i!}, & y_i > 0 \end{cases} \quad (11)$$

$$E(Y_i) = (1 - p)\mu \quad (12)$$

$$\text{Var}(Y_i) = \mu(1 - p)(1 + p\mu) \quad (13)$$

2.3.7. Zero Inflated Negative Binomial

$$f(y_i) = \begin{cases} p + (1 - p) \left(1 + \frac{\mu}{\alpha}\right)^{-\alpha}, & y_i = 0 \\ \frac{(1 - p) \{r(y_i + \alpha)(1 + \frac{\mu}{\alpha})^{-\alpha}\}}{\{r(y_i + 1)r(\alpha)(1 + \frac{\mu}{\alpha})^{y_i}\}}, & y_i > 0 \end{cases} \quad (14)$$

$$E(Y_i) = (1 - p)\mu \quad (15)$$

$$\text{Var}(Y_i) = \mu(1 - p) \left(1 + p\mu + \frac{\mu}{\alpha}\right) \quad (16)$$

3. Results and Discussions

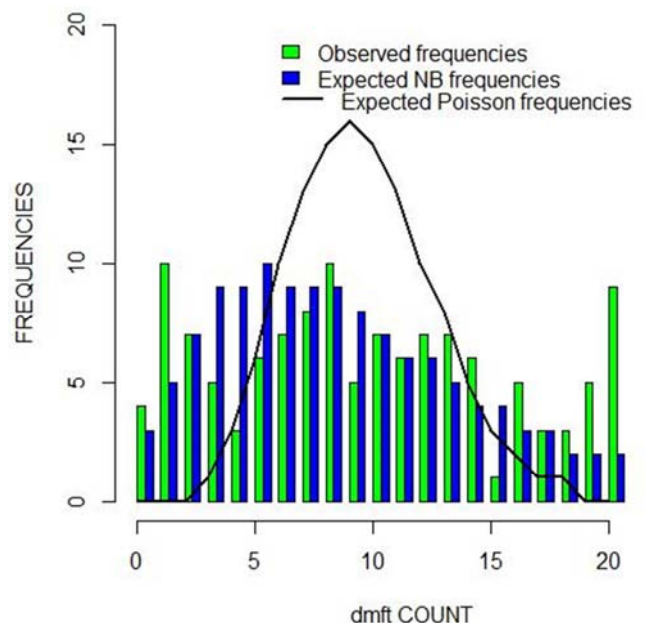


Figure 1. Expected frequencies from intercept-only NB and Poisson models compared with the observed frequencies.

Table 1. Observed distribution of dmft-count and expected frequencies obtained by fitting Poisson, NB, ZIP and ZINB using R.

DMFT	OBSERVED	EXPECTED FREQUENCIES			
		POISSON	NB	ZIP	ZINB
0	4	0	3	4	4
1	10	0	5	0	5
2	7	0	7	0	7
3	5	1	9	1	8
4	3	3	9	3	9
5	6	6	10	5	9
6	7	10	9	8	9
7	8	13	9	12	9
8	10	15	9	14	9
9	5	16	8	15	8
10	7	15	7	15	7
11	6	13	6	13	7
12	7	10	6	11	6
13	7	8	5	8	5
14	6	5	4	6	5
15	1	3	4	4	4
16	5	2	3	2	3
17	3	1	3	1	3
18	3	1	2	1	3
19	5	0	2	0	2
20	9	0	2	0	2
Total	124	124	124	124	124
Mu		9.467742	9.467742	9.782744	9.621402
Alpha			2.409		2.6985
P				0.03220496	0.01593773
Mean	9.467742	9.467742	9.467742	9.467691	9.468059
variance		9.467742	46.67743	12.45051	44.67794
AIC		1001	794.85	961.7486	795.9031

Table 2. Deviances for NB Models fitted to caries data and chosen at each stage.

Model	Deviance
x1+x2	91.71085
x1+x2+x13	91.702
x1+x2+x13+x6	91.69895
x1+x2+x13+x6+x8	91.51689
x1+x2+x13+x6+x8+x4	91.49009
x1+x2+x13+x6+x8+x4+x15	91.64665
x1+x2+x13+x6+x8+x4+x15+x5	91.78052
x1+x2+x13+x6+x8+x4+x15+x5+x10	91.79385
x1+x2+x13+x6+x8+x4+x15+x5+x10+x7	91.60804
x1+x2+x13+x6+x8+x4+x15+x5+x10+x7+x11	91.49541
x1+x2+x13+x6+x8+x4+x15+x5+x10+x7+x11+x9	92.29508
x1+x2+x13+x6+x8+x4+x15+x5+x10+x7+x11+x9+x3	93.09084
x1+x2+x13+x6+x8+x4+x15+x5+x10+x7+x11+x9+x3+x12	97.31071

Table 3. Regression coefficients with NB.

	Parameter	Estimate	Unadjusted Model			Adjusted Model			
			Std. Error	z value	Pr(> z)	Estimate	Std. Error	z value	Pr(> z)
Gender	Intercept	-	-	-	-	2.516131	0.850141	2.96	0.00308 **
	Age	-0.10408	0.07193	-1.447	0.148	-0.1189	0.078651	-1.512	0.13062
	Male	0	-	-	-	0	-	-	-
	Female	0.02629	0.15167	0.173	0.862	0.052574	0.158505	0.332	0.74013
Brushing Frequency	Several times a week	0	-	-	-	0	-	-	-
	Once a day	-0.0599	0.35096	-0.171	0.864	0.092067	0.373909	0.246	0.8055
	Two or more times a day	-0.05264	0.4189	-0.126	0.9	-0.00444	0.439192	-0.01	0.99194
	Feeding habit								

	Parameter	Estimate	Unadjusted Model			Adjusted Model			
			Std. Error	z value	Pr(> z)	Estimate	Std. Error	z value	Pr(> z)
biscuits	Never	0	-	-	-	0	-	-	-
	Several times a month	0.07199	0.27321	0.264	0.7922	0.212701	0.312878	0.68	0.49662
	Once a week	-0.15955	0.24133	-0.661	0.5085	-0.23191	0.248063	-0.935	0.34985
	Several times a week	-0.19416	0.24182	-0.803	0.422	-0.24379	0.250253	-0.974	0.32997
	Everyday	-0.05595	0.2272	-0.246	0.8055	-0.07879	0.234399	-0.336	0.73678
	Several times a day	-0.77632	0.46522	-1.669	0.0952	-0.703	0.473773	-1.484	0.13785
	Feeding habit jam								
Feeding habit jam	Never	0	-	-	-	0	-	-	-
	Several times a month	0.0736	2.90E-01	0.254	0.7999	0.236716	0.321967	0.735	0.46221
	Once a week	0.0000	0.317	0	1	-0.02824	0.318723	-0.089	0.92939
	Several times a week	-0.1.47	0.278	-0.53	0.5964	-0.07925	0.290462	-0.273	0.78496
	Everyday	0.0170	0.193	0.088	0.9298	0.059779	0.191284	0.313	0.75465
	Several times a day	-0.491	0.274	-1.791	0.0734	-0.40661	0.291676	-1.394	0.1633
	Highest education of mother								
Highest education of mother	Less than primary school	0	-	-	-	0	-	-	-
	Primary school completed	-0.04124	0.67738	-0.061	0.95145	0.374512	0.687405	0.545	0.58588
	Secondary school completed	-0.23293	0.67747	-0.344	0.730979	0.280243	0.683119	0.41	0.68163
	College/University	-0.27763	0.67807	-0.409	0.682214	0.162913	0.680414	0.239	0.81077

3.1. Discussions

The first step was to find a suitable distribution for the observed dependent variable, dmft. To illustrate this, Poisson, NB, ZIP and ZINB were fitted to the observed dependent variable data. Only the intercept was fitted for the four models and expected frequencies obtained, as shown in Table 1. The mean was observed to be equal across the four models. However, NB, ZIP and ZINB models demonstrate overdispersion, i.e. variance in each case exceeds the mean. Further, the values of alpha, the dispersion parameter, exceed zero in both cases of NB and ZINB. The probability of being an extra zero was almost nonexistent among children with primary dentition, this is evidenced by the low values of p. AIC was used to choose among the four models. Examination of fit across the four models portrayed the NB as the model that yielded the lowest value of AIC. NB is therefore a good fit for the dmft indices. From figure 1, which uses data from Table 1, it's clear that the NB model produced a good fit for most of the observed frequencies. On the other hand, Poisson, the simplest model, underfitted dmft values of 0, 1, 2, 19, and 20, while overfitting dmft values in the middle range. Majority of the children had dmft counts greater than zero, with only 3.226% being caries free.

3.2. Model Goodness of Fit

Fourteen potentially useful predictors were considered and deviance statistics employed for variable selection in order to build a good, predictive model. Y (dmft count) was regressed

on x1, x2 and x3 only, then on x1, x2 and x4 only, up to x1, x2 and x14 only. In each case, the model with the smallest value of deviance was picked. The process was repeated with initially chosen model and adding each of the remaining independent variables in turn.

Table 2 shows the results of fitting a variety of NB models to the caries data, while controlling for age and gender at each iteration. The model with deviance value of 91.49009 proved the best according to deviance criterion. It provides a good subset of covariates that predict the caries data well. These predictors include: Age(x1), Gender(x2), Brushing frequency(x13), Feeding Habit Biscuits(x6), Feeding Habit Jam(x8), Highest Education Mother(x4).

3.3. Negative Binomial Regression Coefficients

The NB model chosen under model goodness of fit was fitted first without adjustment for covariates, ie, the outcome variable, dmft, was regressed on each of the covariates one at a time. The second step was to regress dmft on all the covariates which gave the results for the adjusted model. (Age and Gender were controlled for in both cases). Analyzing both the adjusted and unadjusted models checked if there was any impact by adjustment for covariates. Categories with parameter estimate of zero were the reference categories and the estimates for the included levels were interpreted relative to them. Test of significance of regression coefficients for any dummy variable is a test of whether belonging to that category with reference to the reference category affects the outcome. Significant parameter estimates are presented with star(s) against the p values.

From Table 3, none of the regression coefficients of various categories is significant. For instance, being a female and not a male does not affect the dmft count. Similarly, eating biscuits several times a month, once a week, several times a week, daily or several times a day as compared to not eating at all has no significant effect on the dmft count. Most of the *p* values under adjusted model were observed to decrease except for brushing two or more times a day, eating biscuits several times a day, eating jam several times a day and several times a week.

Results are as reported in Table 3.

4. Conclusions

The number of caries free children is extremely low (3.226%), indicating large number of children under consideration with caries experience presenting at the clinic. NB regression model fitted the data well. Main risk factors of ECC according to the study were age, gender, brushing frequency, feeding habit biscuits, feeding habit jam and highest education of the mother. Being in any of the observed categories as compared to corresponding reference category has no significant effect on the outcome. Adjustment for the covariates was necessary as it reduced most of the *p* values.

Recommendations

There is need to increase the number of cavity free children. Dentists need to establish ways of treating early cavity as well as providing preventive measures to subjects. The county government together with the national government need to invest in early treatments as well as awareness among the citizens. This calls for concerted efforts including, but not limited to, oral hygiene campaign. It is also necessary to train pediatricians to enhance their skills on guidance about tasks promoting oral health. Additionally, caregivers should be advised on nutrition patterns that help decelerate caries among children with deciduous teeth. Association of dentist in Nairobi and in Kenya should consider global summit that bring researchers together from around the globe in order to lay strategic plans of dealing with changes in ECC prevalence. A national survey should be carried out to scientifically collect important cavity information for proper allocation of resources nationwide to improve living standards among citizens. More research resources can be committed to discovering new products, for treating and preventing ECC.

Acknowledgements

I express my gratitude to the Almighty for bringing me to this point in my studies. I thank my supervisor, Prof. S. M. Mwalili for his rich ideas and patient guidance. I also thank all the staff in the Department of Statistics and Actuarial Sciences, JKUAT for the platform offered to me in the course

of my studies. My deep gratitude to my mother Loise for believing in me and consistent encouragement, my brother Wambugu and my sister Wamaita. Appreciation to my fiancé, Wambaya for his partial financial support. I thank Elda and Linda for helping me with statistical analysis and my classmates. Finally, I extend my special thanks to my relatives and friends for their prayers and well wishes.

References

- [1] Agresti, A., & Kateri, M. (2011). *Categorical data analysis*. Springer.
- [2] Chepkwony, F. C. (2015). *Oral health status and treatment needs among 3-6 year old children attending lady northey dental clinic, nairobi city county*. Unpublished doctoral dissertation, University of Nairobi.
- [3] Cox, S., West, S. G., & Aiken, L. S. (2009). The analysis of count data: A gentle introduction to poisson regression and its alternatives. *Journal of personality assessment*, 91 (2), 121–136.
- [4] Greene, W. (2008). Functional forms for the negative binomial model for count data. *Economics Letters*, 99 (3), 585–590.
- [5] Lee, A. H., Wang, K., Scott, J. A., Yau, K. K., & McLachlan, G. J. (2006). Multilevel zero-inflated poisson regression modelling of correlated count data with excess zeros. *Statistical methods in medical research*, 15 (1), 47–61.
- [6] Marthaler, T. (2004). Changes in dental caries 1953–2003. *Caries research*, 38 (3), 173–181.
- [7] Mwalili, S. M., Lesaffre, E., & Declerck, D. (2008). The zero-inflated negative binomial regression model with correction for misclassification: an example in *caries research*. *Statistical Methods in Medical Research*, 17 (2), 123–139.
- [8] Ngatia, E., Imungi, J., Muita, J. G., et al. (2001). Dietary patterns and dental caries in nursery school children in nairobi, kenya. *East African medical journal*, (12), 673–677.
- [9] Njoroge, N., Kemoli, A., & Gatheche, L. (2010). Prevalence and pattern of early childhood caries among 3-5 year olds in kiambaa, kenya. *East African medical journal*, 87 (3), 134–137.
- [10] Osiro, K., Macigo, & Dienya. (2011). Knowledge, perception and practice of atraumatic restoration treatment among dentists in nairobi.
- [11] Pearl, J. (2015). Detecting latent heterogeneity. *Sociological Methods & Research*, 0049124115600597.
- [12] Sonfield, A., Hasstedt, K., Kavanaugh, M. L., & Anderson, R. (2013). The social and economic benefits of women's ability to determine whether and when to have children. New York: Guttmacher Institute.
- [13] Zuur, A. F., Ieno, E. N., Walker, N. J., Saveliev, A. A., & Smith, G. M. (2009). Zero-truncated and zero-inflated models for count data. In *Mixed effects models and extensions in ecology with r* (pp. 261–293). Springer.