# STA303 Assignment 1

Rebecca Li & Ebrahim Mohamud

## Introduction

Family size is a powerful reflection of a country's social and economic landscape, making it a crucial area of study. Changes in family size over time influence birth rates, fertility rates, and population growth—key drivers of economic stability and social welfare. Analyzing the factors that shape family size provides us with valuable insight into demographic shifts which leads to effective policy implementation.

While fertility and personal preferences are well-established factors behind the number of children that a couple has, the influences behind family size span beyond personal or biological reasons, often tracing to socio-economic and time factors. Freedman (1963) argues that the high birth rates in underdeveloped countries can largely be attributed to social factors such as the age of marriage and cultural norms around abortions. This is especially pronounced in developing countries, where children are utilized as a source of labor and less likely to survive to adulthood compared to their developed country counterparts (Freedman, 1963). These findings are relevant to the context of the present report, which analyzes data from the underdeveloped country of Portugal in order to investigate the primary factors behind family size, allowing for more accurate projections and predictions of future population trends.

Freedman (1963) also highlights the importance of age of marriage on the number of children had, citing that underdeveloped areas witness high birth rates to compensate for high mortality rates prompting people to marry and having children at a younger age. Furthermore, related studies of undeveloped countries demonstrate that literacy can also play a significant role. In a study of 178 males in Chakwal city, Punjab, Mahmood et al. (2016) found that a minority of uneducated men were in favour of smaller families (19.2%) compared to a majority of educated men (77.0%). Mahmood et al. also established that one leading motivation behind larger families in uneducated males was the desire for a son. Together, these two factors played a primary role in motivating our research question.

Finally, while a country's stage of development is crucial for understanding variations in family size, internal factors such as differences between urban and rural regions also play a significant role. Gutkind (1962) conducted an observational study of married African women in rural and urban areas, concluding that urban populations tend to have smaller families than rural populations, attributing this phenomenon to high cost of urban living. The significance of these results encouraged us to consider an additional variable of *region* in our model. The present report expands on the results from the aforementioned studies, using them to guide our analysis of fertility data from a 1979 study of Portuguese families in order to answer the motivating research question: "How do literacy, region, and age of a marriage affect family size?". Specifically, our analysis employs a generalized linear regression model to explain variations in family size

based on the aforementioned variables. We will justify our model choice, interpret the results, and discuss potential limitations, ultimately providing a comprehensive answer to our research question.

## Methods

To answer our research question, we first required a specific functional form for the generalized linear model we are employing. Since the number of children is a count variable distributed over a unit of time, the only appropriate model for this data is a Poisson model. We conducted a preliminary exploratory data analysis (EDA) by visualizing the data through a histogram and observed whether the data satisfied the properties of a Poisson model. Additionally, we compared the mean and variance of each age group within the *age married* variable to assess whether they met the Poisson assumption that the mean must equal the variance. Overdispersion is present when the mean and variance don't equate, therefore, we decided to create two models: A Poisson regression model featuring *age married*, *literacy*, and *region* as predictor variables and a Negative Binomial model which features the same predictors while accounting for overdispersion.

We included *region* in addition to the *age married* and *literacy* variables specific to the research question because of the surrounding literature by Gutkind (1962) which emphasized that rural families have significantly more children than their urban counterparts. Additionally we included an offset in both models to account for variation in exposure times. The offset ensure that each observation yields an estimated rate with a standard time unit of one year, thereby maintaining a consistent interpretation of the intensity of births across our data. Adding the offset is essential to our research question since meaningfully estimating the birth intensity across each family fully specifies our regression model.

We conclude our analysis by comparing the two models based on their estimated coefficients, p-values, confidence intervals, and standard errors. This evaluation allows us to determine which model better captures the relationship between the predictors and the number of children while accounting for overdispersion. The results of the better model according to the above metrics will be interpreted, providing meaningful insights that directly address our research question.

## Results

Our dataset consists of 5,148 observations of married couples from 1979 Portugal, capturing key demographic details such as age at marriage, literacy status, number of children, months since marriage, and region. Since we are interested in how age at marriage, literacy status, and region contribute to variations in family size, we plotted a histogram visualizing our data with respect to our response (the number of children). The histogram (as in Figure 1) exhibits a large right skew with a mean (center) of 2.26 and a standard deviation (spread) of 1.86. To establish further support for our chosen Poisson GLM, we fit a Poisson fitted line with parameters based on the global mean, as well as the means from the various levels within the *age married* predictor to determine if they matched the data. This initial visual assessment illustrates that a Poisson model is a reasonable fit both overall and across the different categories of the age at marriage predictor variable. To improve balance across categories, we redefined the age brackets to ensure each level contained a roughly equal number of observations. This adjustment helped prevent sparsely populated categories, reducing the risk of unstable estimates and enhancing the reliability of our analysis.
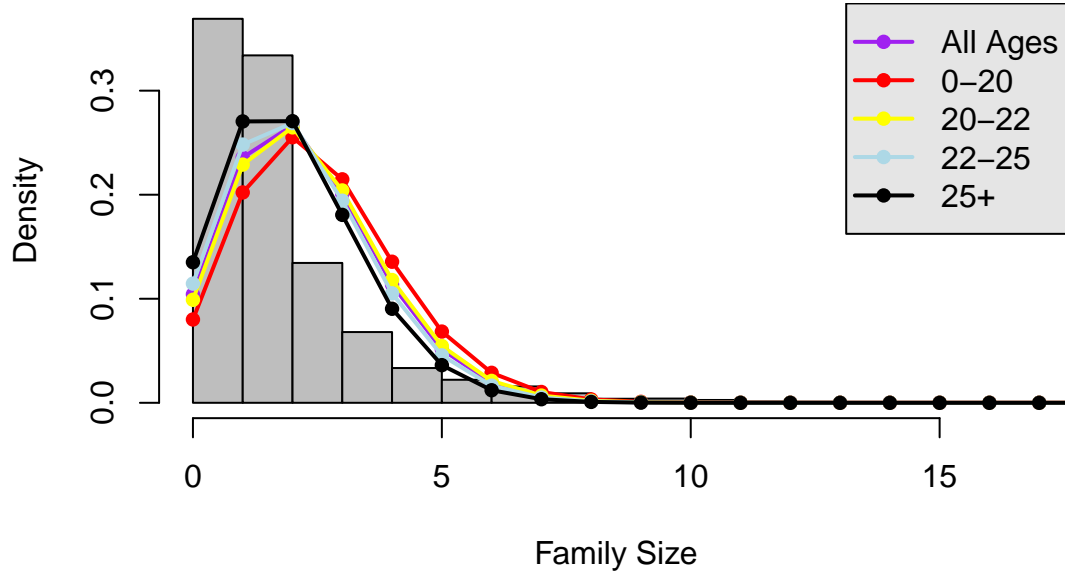
## Histogram of the Number of Children



Figure 1: . . .

Before committing to a Poisson model, we needed to ensure all the Poisson assumptions were met within our data. Therefore, we proceeded by checking whether the mean and variance of each age group within the *age married* variable equated or whether they didn't. As Table 1 illustrates, there is a positive gap between the variances and means suggesting possible overdispersion in the data. Therefore, our analysis compares a Poisson regression model against a Negative Binomial model according to their estimated coefficients, p-values, confidence intervals, and standard errors to determine which model more appropriately answers our research question.

Table 1: . . .

| Age Married | Mean Num Children | Variance Num Children | Variance-Mean Gap |
|---|---|---|---|
| 0to20 | 2.525 | 4.667 | 2.142 |
| 20to22 | 2.313 | 3.522 | 1.209 |
| 22to25 | 2.167 | 2.847 | 0.680 |
| 25+ | 2.002 | 2.550 | 0.548 |

After fitting the Poisson and Negative Binomial models, we summarized and displayed the coefficient estimates and confidence intervals on the natural scale, along with the standard errors, test statistics, and p-values, in Figures 2 and 3, respectively. Comparing the two tables, we notice that both models have similar coefficient estimates, with differences occurring in the hundredths decimal place. Both models had similarly low standard errors (with slightly larger standard errors for the Negative Binomial model), and both models indicate that the variables *Age Married 20-22*, *Age Married 25+*, and *Region Porto* are insignificant in explaining family size (at a significance level of 0.05).
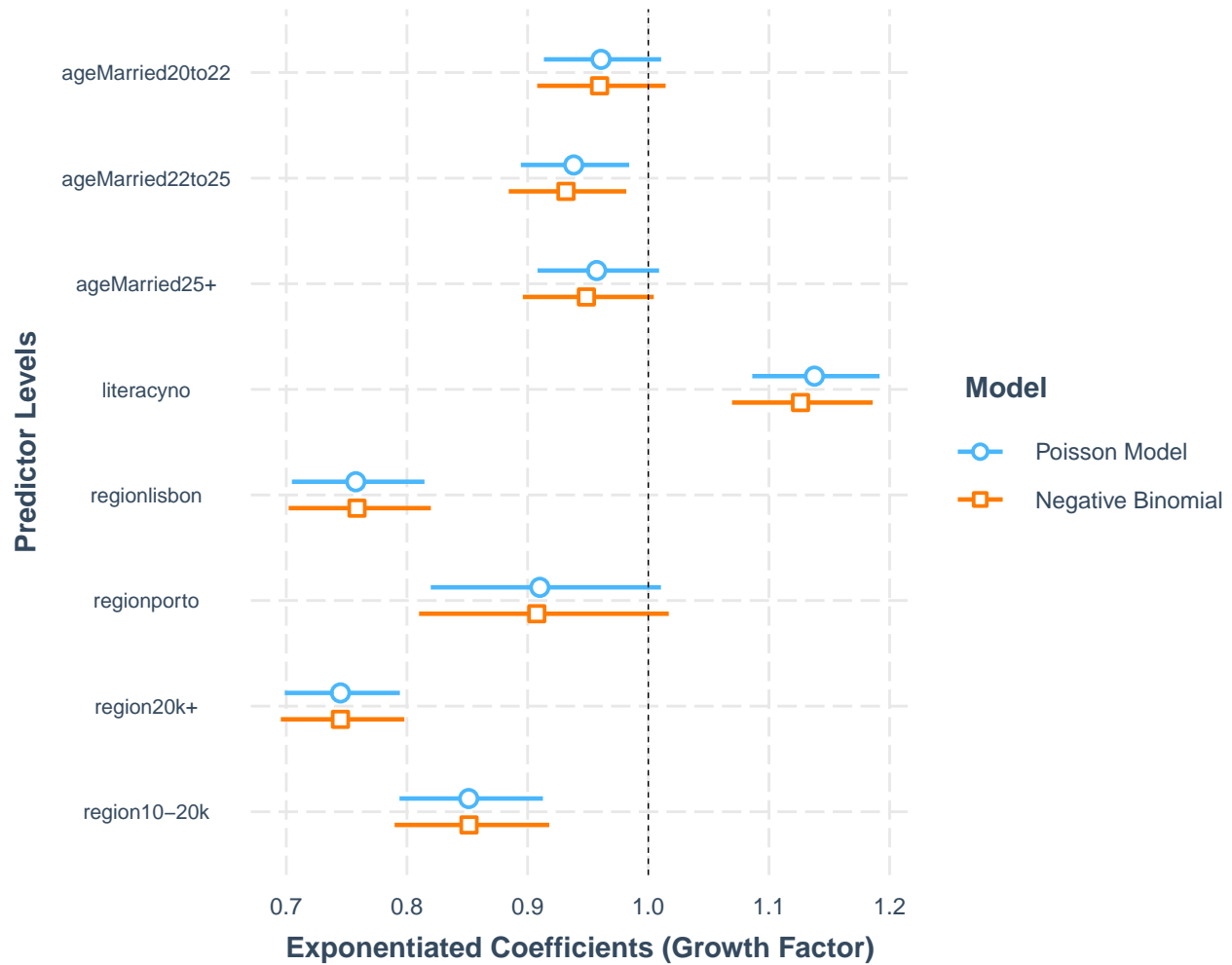
Table 2: . . .

|  | 2.5 % | 97.5 % | Estimate | Standard Error | Test Statistic | P_Value |
|---|---|---|---|---|---|---|
| Intercept | 0.183 | 0.196 | 0.189 | 0.003 | -90.239 | 0.000 |
| No Literacy | 1.086 | 1.191 | 1.138 | 0.027 | 5.455 | 0.000 |
| Age Married: 20-22 | 0.913 | 1.011 | 0.961 | 0.025 | -1.550 | 0.121 |
| Age Married: 22-25 | 0.894 | 0.984 | 0.938 | 0.023 | -2.613 | 0.009 |
| Age Married: 25+ | 0.908 | 1.009 | 0.957 | 0.026 | -1.628 | 0.103 |
| Region: Lisbon | 0.704 | 0.814 | 0.758 | 0.028 | -7.512 | 0.000 |
| Region: Porto | 0.818 | 1.009 | 0.910 | 0.049 | -1.765 | 0.078 |
| Region: Population 20k+ | 0.698 | 0.794 | 0.745 | 0.024 | -9.021 | 0.000 |
| Region: Population 10-20K | 0.793 | 0.912 | 0.851 | 0.030 | -4.527 | 0.000 |

Table 3: . . .

|  | 2.5 % | 97.5 % | Estimate | Standard Error | Test Statistic | P-Value |
|---|---|---|---|---|---|---|
| Intercept | 0.186 | 0.202 | 0.194 | 0.004 | -80.196 | 0.000 |
| No Literacy | 1.069 | 1.186 | 1.126 | 0.030 | 4.503 | 0.000 |
| Age Married: 20-22 | 0.908 | 1.014 | 0.960 | 0.027 | -1.457 | 0.145 |
| Age Married: 22-25 | 0.884 | 0.982 | 0.932 | 0.025 | -2.649 | 0.008 |
| Age Married: 25+ | 0.896 | 1.005 | 0.949 | 0.028 | -1.805 | 0.071 |
| Region: Lisbon | 0.702 | 0.820 | 0.759 | 0.030 | -6.977 | 0.000 |
| Region: Porto | 0.810 | 1.017 | 0.908 | 0.053 | -1.671 | 0.095 |
| Region: Population 20k+ | 0.695 | 0.798 | 0.745 | 0.026 | -8.406 | 0.000 |
| Region: Population 10-20K | 0.790 | 0.918 | 0.851 | 0.033 | -4.191 | 0.000 |

Figure 2 compares the exponentiated coefficients (rate ratios) with confidence intervals for Poisson and Negative Binomial models. Both models yield similar coefficient estimates and trends. However, the Negative Binomial model generally displays slightly wider confidence intervals, reflecting its ability to account for overdispersion. For the *Age Married 25+* category, the Negative Binomial model's confidence interval covers less of the line at 1.0 compared to the Poisson model, indicating improved precision. Also the model estimated an overdispersion parameter of 24.4% (CI 21.4%, 27.8%) indicating that birth rates vary by about 24.4% between families. Therefore, we believe a Negative Binomial model is better at explaining variation in family size and it serves as our model choice.

## Conclusion

Overall, the results from our data analysis align with ideas from pre-existing literature and demonstrates a positive progression in selecting an appropriate model to fit the data. Both our original Poisson model and the Negative Binomial model demonstrated generally significant p-values for variables such as literacy and region which is supported by the results of Gutkind's (1962) and Mahmood et al.'s (2016). However, contrary to our previous literature review, some variable levels did not prove to be as significant as initially thought, such as *Age Married 20-22*, *Age Married 25+*, and *Region Porto*.

In particular, the variable *literacy* under the Negative Binomial model had an estimated coefficient of 1.126 on the natural scale (see Figure 3). Since this represents a growth factor; the birth rate of uneducated couples compared to their educated counterparts increases by a factor of 1.126, holding all else constant. This trend matches the findings from Mahmood et al.'s (2016) which found that 77% of educated males favoured smaller families compared to only 19.2% of uneducated males ($p < 0.001$). Additionally, the variable at level *Age Married 22-25* has an estimated coefficient 0.932 (see Figure 3) which is interpreted as a decrease in birth rates by a factor of 0.932 for couples married at an age between 22 to 25 compared to couples married at the age between 0 to 20, holding all else constant. This trend matches the findings from Freedman (1963)

which also found an inverse relationship between age of marriage and the birth rates.

Our final model suggests that both age at marriage, literacy status, and region play significant roles in shaping family size. Specifically, the results indicate that marrying later as well as higher literacy status is generally associated with having fewer children. Our Negative Binomial model seems to be an appropriate fit for the 1979 Portugal fertility data, demonstrating a significant relationship between literacy, region, and family size, with room for further investigation of the age of a marriage, in order to deeper understand the population behaviors and trends.

## Citation:

Freedman, R. (1963). Norms for family size in underdeveloped areas. Proceedings of the Royal Society of London. Series B. Biological Sciences, 159(974), 220–245. https://doi.org/10.1098/rspb.1963.0074

Gutkind, Peter C. W. (1962). African Urban Family Life: Comment on and Analysis of Some Rural-Urban Differences. Cahiers d'Études Africaines, 3(10), 149–217. http://www.jstor.org/stable/4390830

Mahmood, H., Khan, Z., & Masood, S. (2016). Effects of male literacy on family size: A cross sectional study conducted in Chakwal city. J Pak Med Assoc, 66(4), 399-403.