

MODELING GRADUATE EMIGRATION IN NIGERIA USING LOG-LINEAR APPROACH

**¹O. S. BALOGUN, ²D. E. BRIGHT, ¹A. A. AKINREFON,
¹S. S. ABDULKADIR**

¹Modibbo Adama University of Technology, NIGERIA

²Michael Imoudu National Institute for Labour Studies, NIGERIA

Abstract. This research work was carried out as a result of the alarming rate at which graduates in Nigeria, which should form the bulk of our nation's workforce and ensure the growth and development of our great nation, leave our shores in search of greener pastures. The main focus of this research work is to develop models that can be used to study the trend of graduate emigration in Nigeria using log-linear modeling based on the results from of likelihood ratio (G^2), Akaike Information Criteria (AIC) , Saturated model has a perfect fit for modeling graduate emigration in Nigeria. This implies all the three factors involved (discipline, year and sex) has to be included in the model in order to have an appropriate result.

Keywords: contingency table, log-linear model, log-likelihood ratio, Akaike information criteria

Introduction

Over the years, the issue of graduate migration has become a source of concern to the government as an increasing number of its human resource leave in search of 'greener pastures' in foreign lands. Migration profiles were first proposed by the European commission in the 'Communication on Migration and Development' in 2005. According to this document, migration profiles should be a statistical report that provides information on a range of issues related to migration in the European Commission to provide information to community-assistance programs for third countries in the field of migration, as well as poverty-reduction strategies (Afolayan, 2009). Although Nigeria is traditionally an important destination for migrants in the West African region, there are more people emigrating from, than immigrating to Nigeria. The net migration rate (per 1000 people) has increasingly become negative in recent years, decreasing from -0.2 in 2000 to -3.3 in 2005, this trend is expected to continue. According to recent estimates, the net migration rate will decrease to -0.4 in 2010.¹⁾

The major focus of this research work is to develop models that can be used to study the trend/pattern of graduate emigration in Nigeria and to compare the various models and determine the best for achieving the above. The data used for this research is a secondary data consisting of the Discipline, Year of emigration and Gender of Nigerian graduate that leave the country. The Disciplines include Administration, Arts, Education, Medical, Science, Energy/Technology and Business between 2002/2003 and 2004/2005 for both males and females collected from the National Universities Commission.

This research work was prompted as a result of alarming rate at which graduates in Nigeria, which should form the bulk of our nation's workforce and ensure the growth and development of our great nation, leaves our shores in search of greener pastures. The data for this work span 3 years (from 2002/2003 to 2004/2005). From the final model, one can be able to predict the

expected number of graduates that leave and as such make adequate plans to stem this action and also make plans to harness their potentials for the development of both the nation and the graduates alike.

Definition of terms

Migration: this is the movement of large numbers of people from one place to another. It can be divided into: (i) *Immigration*: the process of coming to live permanently in a country that is not your own; (ii) *Emigration*: the process of leaving your own country to go and live permanently in another country.

Types of migration include those for: refugee/asylum-seekers, labor-migrants, students and tourists and visitors.

Contingency tables

The multidimensional table in which each dimension is specified by discrete variable or grouped continuous (range) variable gives a basic summary for multivariate discrete and grouped continuous data. If the cell of the table are number of observation in the corresponding values of the discrete variables then it is *contingency tables*. The discrete or grouped continuous variables that can be used to classify a table are known as *factors*. Examples include sex (male or female), religion (Christianity, Islam, traditional, etc.).

Types of Contingency Tables: (i) one dimensional ($1 \times J$) tables; (ii) two dimensional ($I \times J$) tables; (iii) square tables ($I \times I$); (iv) multidimensional tables

Log-linear models for contingency tables

The concept of log-linear analysis in contingency tables is analogous to the concept of analysis of variance (ANOVA) for the continuously distributed factor response variables while response observations are assumed to be

continuous with underlying Normal distributions in ANOVA, the log-linear analysis assume that the response observations are counts having Poisson distributions (Lawal, 2003).

Log-linear models are parts of the generalized linear models (GLMs) using the log link function with a Poisson response. The models specify how the expected counts depend on the levels of the categorical variables for that cell as well as association and interactions among variables (Agresti, 2002)

Common sampling distributions for two-way classification

The sampling scheme for different two-way contingency table varies from one table to another depending on the underlying assumptions that contribute to the table. Suppose we have observed counts n_{ij} in the $k=I \times J$ cells of a contingency table, such that each n_{ij} has expected value denoted by $m_{ij} = E(n_{ij})$.

Poisson sampling

We assume Poisson sampling plan for each cell n_{ij} when there is no restriction in the total sample size. The probability mass function for the n_{ij} is

$$p(n_{ij}) = \frac{m_{ij}^{n_{ij}} e^{-m_{ij}}}{n_{ij}!} \quad (1)$$

$$n_{ij} = 0, 1, 2, \dots$$

The mean and the variance of this mass function is m_{ij} $E(n_{ij}) = \text{var}(n_{ij}) = m_{ij}$. The Poisson model assumes that n_{ij} are independent. Thus the joint probability mass function for n_{ij} is the product of the probabilities in the K cells. The Poisson model is used for rare events which are independently distributed over disjoint classes.

Multinomial sampling

We assume multinomial sampling scheme when the total sample size $n_{..}$ is fixed. The restriction imposed on a series of independent Poisson distributions gives the multinomial distributions (Fisher, 1922) with probability density function:

$$F(n_{ij}) = \frac{n_{..}!}{\prod_{ij} n_{ij}!} \prod_{ij} \left(\frac{n_{ij}}{n_{..}} \right)^{n_{ij}} \quad (2)$$

This is often written with π_{ij} representing the probability of a count falling in cell (I, J). Substituting $\frac{n_{ij}}{n_{..}}$ for π_{ij} , we have

$$F(n_{ij}) = \frac{n_{..}!}{\prod_{ij} n_{ij}!} \prod_{ij} (\pi_{ij})^{n_{ij}}$$

The multinomial distribution also applies when n independent observations are taken from a probability distribution concentrated on a set of K categories.

Product multinomial sampling

Suppose we observe on a categorical variable Y at various levels of an explanatory variables X. if in the cell (X=I, Y=j) we have n_{ij} observations, then the resulting distribution in the two dimensional group totals in the product of independent multinomial:

$$F(n_{ij}) = \frac{n_{i.}!}{\prod_j n_{ij}!} \prod_j (\pi_{j/i})^{n_{ij}} \quad (3)$$

Special sampling plan

Apart from the above sampling schemes, we also have the hypergeometric sampling plan used when the marginal are fixed.

The Poisson and (Product) multinomial schemes are the commonly used in sampling distribution for categorical data.

Log-linear models for three categorical variables

Let n_{ijk} denote the number of sample units for which we observe $D = i, S = j, Y = k$, where $i = 1, 2, \dots, I, j = 1, 2, \dots, J$ and $k = 1, 2, \dots, K$ - levels for the respective categorical variables are. Also let $\pi_{ijk} = P(D = i, S = j, Y = k)$ and $\log(n\pi_{ijk}) = \log(m_{ijk})$. The most complex Log-linear model (saturated model) for this class of variables is

$$\log(m_{ijk}) = \mu + \lambda_i^D + \lambda_j^S + \lambda_k^Y + \lambda_{ij}^{DS} + \lambda_{jk}^{DY} + \lambda_{jk}^{SY} + \lambda_{ijk}^{DSY} \quad (4)$$

$i = 1, 2, \dots, I, j = 1, 2, \dots, J$ and $k = 1, 2, \dots, K$ where for identifiability, the λ_j terms are constrained to sum to zero over any subscript

$\sum_{i=1}^I \lambda_i^D = \sum_{j=1}^J \lambda_j^S = 0, \sum_{k=1}^K \lambda_k^Y = 0$ and $\sum_{i=1}^I \lambda_{ik}^{DY} = \sum_{j=1}^J \lambda_{jk}^{SY} = 0$ and so on. Simpler models set higher order interactions to zeros.

Statistical independence

The statistical independence between row and column variables means that the joint probabilities $\{\pi_{ij}\}$ of the observed falling into a cell are equal to the product of the marginal probabilities.

Mathematically, it is expressed as

$$\pi_{ij} = \pi_i \cdot \pi_j \quad \forall i = 1, 2, \dots, I \text{ and } j = 1, 2, \dots, J$$

So in terms of expected frequencies (cell counts)

$$M_{ij} = n\pi_{ij} \quad \forall i, j \text{ (Some text use } \mu_{ij} \text{ for } m_{ij} \text{)}$$

The probability $\{\pi_{ij}\}$ and the expected frequencies are $\{(\mu_{ij}) = m_{ij} = n\pi_{ij}\}$. Log-linear model formulas use $\{m_{ij}\}$ as their response variable rather than the cell probabilities $\{\pi_{ij}\}$, therefore the random component in Poisson (Agresti, 2002; Adejumo, 2005).

For the definition of statistical independence, the model for the expected number of counts in multiplicative.

$$m_{ij} = n\pi_{ij} = n\pi_{i.}\pi_{.j} \quad \forall i = 1, 2, \dots, I \text{ and } j = 1, 2, \dots, J$$

and taking logarithms gives us

$$\log(m_{ij}) = \log(n) + \log(\pi_{i.}) + \log(\pi_{.j})$$

which can be further expressed as

$$\log(m_{ij}) = \mu + \lambda_i^D + \lambda_j^S \tag{5}$$

It is known as the “log-linear model of independence” for 2-ways contingency tables where μ represents an “overall” effect or a constant. This term ensures $\sum_i \sum_j m_{ij} = n$; λ_i^D represents the “main” or marginal effect of the row variable D. It represents the effect of effect of classification in row i. λ_i^D 's ensure that $\sum_j m_{ij} = m_{i.} = n_{i.}$; λ_j^S represents the main of marginal effect of the column variable S. It represents the effect of classification in column j. this term ensures that

$$\sum_i m_{ij} = m_{.j} = n_{.j}$$

Also, $\sum_{i=1}^I \lambda_i^D = 0 \quad \sum_{j=1}^J \lambda_j^S = 0$ with $df = (I - J)(J - 1)$

Saturated log-linear models

Statistically dependent variables satisfy a more complex log-linear model

$$\log(m_{ij}) = \mu + \lambda_i^D + \lambda_j^S + \lambda_{ij}^{DS} \quad (6)$$

where $\mu, \lambda_i^D, \lambda_j^S$, are the overall and marginal effect terms as defined terms as defined earlier λ_{ij}^{DS} represents the association or interaction between D and S. It reflects deviations from independence (Agresti, 2002). It further ensures that $m_{ij} = n_{ij}$.

$$\text{Also, } \sum_{i=1}^I \lambda_i^D = 0, \sum_{j=1}^J \lambda_j^S = 0 \text{ and } \sum_{i=1}^I \lambda_{ij}^{DS} = \sum_{j=1}^J \lambda_{ij}^{DS} = 0$$

The right hand side of Eq. (6) resembles the formula for cells means in a two-way ANOVA allowing interaction. The λ_{ij}^{DS} represent interaction between D and S; whereby the effect of one variable on m_{ij} depends on the level of the other. From the above constraints, the model that include the interaction terms, fits the data perfectly, such that fitted values are exactly equal to observed values and has many unique parameters as there are number of cells in the table. That model is called “saturated model”. This has independence as a special feature when the interaction effect is zero, that is

$$\lambda_{ij}^{DS} = 0 \quad \forall i \text{ and } j$$

Direct relationship exist between the model parameters and log odds ratio, which is how we are defining and measuring interactions

$$\begin{aligned} \log(\theta_{ii^*jj^*}) &= \log\left(\frac{m_{ij}m_{i^*j^*}}{m_{i^*j}m_{ij^*}}\right) \\ &= \log(m_{ij}) + \log(m_{i^*j^*}) - \log(m_{i^*j}) - \log(m_{ij^*}) \\ &= (\mu + \lambda_i^D + \lambda_j^S + \lambda_{ij}^{DS}) + (\mu + \lambda_{i^*}^D + \lambda_{j^*}^S + \lambda_{i^*j^*}^{DS}) - (\mu + \lambda_{i^*}^D + \lambda_j^S + \lambda_{i^*j}^{DS}) - (\mu + \lambda_i^D + \lambda_{j^*}^S + \lambda_{ij^*}^{DS}) \\ &= \lambda_{ij}^{DS} + \lambda_{i^*j^*}^{DS} - \lambda_{i^*j}^{DS} - \lambda_{ij^*}^{DS} \end{aligned}$$

The odds-ratio measures the strength of the association and depends only on the interaction terms λ_{ij}^{DS} . There is need to completely characterize (I-1), (J-1) and (I-1) (J-1) main effect and association effect in a (I×J) table, which is the number of unique parameters that are there. This is necessary to overcome the issue of over parameterization (that is, having more parameters than the underlying probabilities). This is done by imposing some constraints or identifiability conditions on the parameters. Below are three forms of constraints usually imposed to solve this problem (Lawal, 2003): (i) Sum-to-zero constraints on the parameters. Here the parameters are constrained to sum to zero either row-wise or column-wise for main and interaction effects; (ii) Only the parameters of the last category of each variable and corresponding interaction terms are set to be zero; (iii) Only the parameter of the first category of each variable and corresponding interaction terms is set to zero.

Log-linear model fitting

After selecting a log-linear model, the observed data are used to estimate model parameters, cell probabilities, and expected frequencies. Although alternative methods of estimation are sometimes useful, the maximum likelihood (ML) method offers several advantages. First of all, the MLEs for hierarchical log-linear models are relatively easy to compute, since the estimates satisfy certain intuitive marginal constraints. In addition, the ML method can be used when data are sparse, that is, when there are several observed cell counts of zero (Note that marginal totals, however, cannot be equal to zero) (Stokes et al, 2000).

Minimal sufficient statistics

For three-ways tables, the Poisson probability that cell counts $(Y_{ijk} = n_{ijk})$ is

$$\prod_i \prod_j \prod_k \frac{\ell^{-\mu_{ijk}} \mu_{ijk}^{n_{ijk}}}{n_{ijk}}$$

where the product refers to all cells of the table. The kernel of the log-likelihood is

$$L(\mu) = \sum_i \sum_j \sum_k n_{ijk} \log \mu_{ijk} - \sum_i \sum_j \sum_k \mu_{ijk} \quad (7)$$

The fitted values for a model are solutions to the likelihood equations using the general representation ($\log \mu = X\beta$) for a log-linear model. For a vector of counts n with $\mu = E(n)$, the model is $\log \mu = X\beta$, for which $\log(\mu_i) = \sum_j x_{ij} \beta_j$ for all i .

Extending Eq. (7) for Poisson sampling, the log-likelihood is

$$\begin{aligned} L(\mu_i) &= \sum_i n_i \log \mu_i - \sum_i \mu_i \\ &= \sum_i n_i \left(\sum_j x_{ij} \beta_j \right) - \sum_i \exp \left(\sum_j x_{ij} \beta_j \right) \end{aligned}$$

The sufficient statistic for β_j is its coefficient, $\sum_i n_i x_{ij}$. Since

$$\begin{aligned} \frac{\partial}{\partial \beta_j} \left[\exp \left(\sum_j x_{ij} \beta_j \right) \right] &= x_{ij} \exp \left(\sum_j x_{ij} \beta_j \right) = x_{ij} \mu_i \\ \frac{\partial L(\mu)}{\partial \beta_j} &= \sum_i n_i x_{ij} - \sum_i \mu_i x_{ij}, j = 1, 2, \dots, p \end{aligned}$$

The likelihood equations equate these derivatives to zero.

The parameters can be estimated via direct or iterative methods which include Newton-Raphson method, iterative proportional fitting, Fisher's scoring iterative, etc.

Testing goodness of fit

After the ML estimates have been obtained, the test statistic is then approximately χ^2 with the number of degrees of freedom reduced by the number of estimated parameters.

Under

$$H_0 : \pi_{ij} = \pi_i \pi_j \quad \forall i, j$$

we have the following test.

Pearson chi-square statistic

Having obtained fitted cell counts, we can assess model goodness of fit by comparing them to the observed cell counts. We use Chi-squared statistic to test the hypothesis that population expected frequencies satisfy a given model.

The likelihood ratio Chi-squared equals

$$\chi^2 = \sum \frac{(n_{ij} - m_{ij})^2}{m_{ij}} \approx \chi_{(r-1)(c-1), 1-\alpha}^2 \quad (9)$$

The likelihood ratio test

Under

$$H_0 = \pi_{ij} = \pi_i \pi_j$$

vs

$$H_1 = \pi_{ij} \neq \pi_i \pi_j$$

$$\Lambda = \frac{\prod_{i=1}^I \prod_{j=1}^J (n_{i \cdot} n_{\cdot j})^{n_{ij}}}{n^n \prod_{i=1}^I \prod_{j=1}^J (n_{ij})^{n_{ij}}}$$

It follows that Wilks's G^2 is given by

$$G^2 = -2 \ln \Lambda \quad (9)$$

$$G^2 = 2 \sum_{i=1}^I \sum_{j=1}^J n_{ij} \ln \left(\frac{n_{ij}}{m_{ij}} \right)$$

$$\text{with } m_{ij} = \frac{n_{i.} n_{.j}}{n_{..}} \text{ (estimate under } H_0 \text{)}$$

If H_0 holds, Λ will be large (i.e., near 1) and G^2 will be small. This means that H_0 is to be rejected for large G^2 .

Criterion for model selection

Akaike Information Criterion (AIC): It judges a model by how close its fitted values tend to be to the true values, in terms of a certain expected value. Even though a simple model is farther from the true model than a more complex model, it may be preferred because it tends to provide better estimates of certain characteristics of the true model, such as cell probabilities. Thus, the optimal model is the one that tends to have fit close to reality.

AIC = -2(Maximum Likelihood – number of parameter in the model).

This penalizes model for having many parameter.

Methodology

Emigration from Nigeria has been on the increase; most highly skilled persons leave the country creating skill shortages in sensitive sectors such as technology and health. The education sector lacks the capacity to replace the skills that leave the country. In this research work, we will use log-linear model approach to analyze the pattern of graduate movement out of the country. A program was written in the R programming language to execute this.

The general log-linear model for the three-way is

$$\log(m_{ijk}) = \lambda + \lambda_i^D + \lambda_j^S + \lambda_k^Y + \lambda_{ij}^{DS} + \lambda_{ik}^{DY} + \lambda_{jk}^{SY} + \lambda_{ijk}^{DSY} \quad (10)$$

where D denotes discipline, S denotes sex and Y denotes year.

Model in focus

Based on the data used for this research work, information on discipline (D), year (Y) and sex (S) on the trend of graduate emigration in Nigeria is available.

Log-linear model formation is

$$\text{Model 5 (Saturated Model): } \log(m_{ijk}) = \mu + \lambda_i^D + \lambda_j^S + \lambda_k^Y + \lambda_{ij}^{DS} + \lambda_{ik}^{DY} + \lambda_{jk}^{SY} + \lambda_{ijk}^{DSY}$$

$$\text{Model 4 (DY.DS.SY)} \quad : \log(m_{ijk}) = \mu + \lambda_i^D + \lambda_j^S + \lambda_k^Y + \lambda_{ij}^{DS} + \lambda_{ik}^{DY} + \lambda_{jk}^{SY}$$

$$\text{Model 3 (DY.SD)} \quad : \log(m_{ijk}) = \mu + \lambda_i^D + \lambda_j^S + \lambda_k^Y + \lambda_{ij}^{SD} + \lambda_{ik}^{DY}$$

$$\text{Model 2 (SD.Y)} \quad : \log(m_{ijk}) = \mu + \lambda_i^D + \lambda_j^S + \lambda_k^Y + \lambda_{ij}^{SD}$$

$$\text{Model 1 (D.S.Y)} \quad : \log(m_{ijk}) = \mu + \lambda_i^D + \lambda_j^S + \lambda_k^Y$$

Data presentation and analysis

Data

Below is a summary table showing the cross classification of the graduate discipline, sex and year of emigration. These were used in the computation of the parameter estimates for the respective log-linear models.

Discipline	2002/2003		2003/2004		2004/2005	
	Male	Female	Male	Female	Male	Female
Administration	12670	6843	7777	5116	5315	3215
Art	7532	7107	5771	4798	3147	2363
Education	5313	4713	3958	4405	3031	2634
Medical	3811	2039	3687	1686	1628	711
Sciences	8654	5666	7257	3965	5782	1894

Energy/Technology	6199	1028	4989	819	1824	188
Business	10693	6662	8853	5269	4170	3113
Total	58872	34058	42292	26058	24897	14118
Grand Total	88930		68350		39015	

Analysis

The result of the analysis of graduate emigration in Nigeria in respect of their discipline, sex and year of emigration is shown below:

$$Model\ 1: \log(m_{ijk}) = \mu + \lambda_i^D + \lambda_j^S + \lambda_k^Y$$

Statistic	X^2	df	$p(> X^2)$
Likelihood ratio	9103.5	32	< 0.00001
Pearson	8611.7	32	< 0.00001

AIC = 9543.8

$$Model\ 2: \log(m_{ijk}) = \mu + \lambda_i^D + \lambda_j^S + \lambda_k^Y + \lambda_{ij}^{SD}$$

Statistic	χ^2	d.f.	$p(> X^2)$
Likelihood ratio	2226.1	26	< 0.00001
Pearson	2213.3	26	< 0.00001

AIC = 2680.4

$$Model\ 3: \log(m_{ijk}) = \mu + \lambda_i^D + \lambda_j^S + \lambda_k^Y + \lambda_{ij}^{SD} + \lambda_{ik}^{DY}$$

Statistic	χ^2	d.f.	$p(> X^2)$
Likelihood ratio	836.1	14	< 0.00001
Pearson	817.7	14	< 0.00001

AIC = 1315.2

$$\text{Model 4: } \log(m_{ijk}) = \mu + \lambda_i^D + \lambda_j^S + \lambda_k^Y + \lambda_{ij}^{SD} + \lambda_{ik}^{DY} + \lambda_{jk}^{SY}$$

Statistic	χ^2	d.f.	$p(> X^2)$
Likelihood ratio	723.5	12	< 0.00001
Pearson	715.7	12	< 0.00001

AIC = 1205.8

$$\text{Model 5: } \log(m_{ijk}) = \mu + \lambda_i^D + \lambda_j^S + \lambda_k^Y + \lambda_{ij}^{SD} + \lambda_{ik}^{DY} + \lambda_{jk}^{SY} + \lambda_{ijk}^{DSY}$$

Statistic	χ^2	d.f.	$p(> X^2)$
Likelihood ratio	0	0	1
Pearson	0	0	1

AIC = 506.3

Summary table

	G^2	d.f.	AIC	P-value
Model 1	9103.5	32	9543.8	< 0.00001
Model 2	2226.1	26	2680.4	< 0.00001
Model 3	836.9	14	1312.2	< 0.00001
Model 4	723.5	12	1205.8	< 0.00001
Model 5	0	0	0	1

Summary and conclusion

This research work examined the fitted models to the emigration of graduate in Nigeria in respect to their discipline, year of emigration and sex. Comparing the likelihood ratio statistic, model 1 has the highest G^2 value followed by model 2, 3, 4 and the least is model 5 which is saturated model. Also, using Akaike Information Criteria (AIC) to confirm the best model, saturated model fit best because it has the minimum AIC value. Model 5 (saturat-

ed model) is also chosen because it has the lowest G^2 value as well as the least AIC value compare to the other models.

Based on the results on the likelihood ratio G^2 and Akaike Information Criteria (AIC), model 5 which is a saturated model is more appropriate for modeling graduate emigration in Nigeria. Hence, to study the graduate emigration in Nigeria, all the three factors involved (discipline, sex and year) have to be included in the model to have a reasonable and an appropriate result.

NOTES

1. <http://www.un.org/esa/population/publications/migration/WorldMigrationReport2009.pdf>

REFERENCES

- Adejumo, A.O. (2005). *Modeling generalized linear (log-linear) models for raters' agreement measure*. Frankfurt: Peter Lang.
- Afolayan, A. (2009). *Migration in Nigeria: a country profile*. Geneva: International Organization for Migration.
- Agresti, A. (2002), *An introduction to categorical data analysis*. New York: John Wiley & Sons.
- Fisher R.A. (1922). The goodness of fit of regression formulae and the distribution of regression coefficient. *J. Roy. Stat. Soc.*, 85, 597-612. .
- Lawal, B. (2003). *Categorical data analysis with SAS and SPSS application*. London: Lawrence Erlbaum Associates.
- Stokes, M.E., Davis, C.S. & Koch, G.G. (2000). *Categorical data analysis using the SAS system*. Cary: SAS.

✉ Mr. O.S. Balogun (corresponding author)
Department of Statistics and Operations Research
Modibbo Adama University of Technology
PMB 2076, Yola, Adamawa State, Nigeria
E-Mail: stapsalms@yahoo.com

