## *Research*

*Corresponding author
**Paul Juneau, MS***
Statistical Services Group
Truven Health Analytics
21937 Greenbrook Drive
Boyds, MD 20841
USA
E-mail: paul.juneau@truvenhealth.com

# Analyzing Pregnancy Costs with Finite Mixture Models: An Opportunity to More Adequately Accommodate the Presence of Patient Data Heterogeneity

**Paul Juneau***

*Statistical Services Group, Truven Health Analytics, 21937 Greenbrook Drive, Boyds, MD 20841, USA*

## ABSTRACT

The choice of a model in the analysis of patient health care costs and utilization is critical for a clear understanding of the behavior and estimation of quantities like incremental costs or cost-effectiveness. In studying heath care claims related to pregnancy, it would not be surprising that a small portion of the women have costs associated with their care and treatment that might be extreme or outlying. Many strategies exist for accommodating outliers; however, is one approach superior to the others because it may be implemented over a broader set of conditions without making unreasonable assumptions about the prevailing data characteristics? In this study, the author will show an example of a data set based on the medical claims for over 300K pregnant women, aged 15-49, where the traditional, or widely used Generalized Linear Model (GLM) approach to modeling costs may be less than optimal due to the presence of patients with very large, or very small expenditure values. These values, in some sense "contaminate" the typically employed GLM and cause it to violate its underlying requisite statistical assumptions.

Finite Mixture Models (FMMs) have been employed in other areas of clinical research to model health care utilization. The author will introduce FMMs as an alternative to the commonly used GLM model and show that in his example data set, the fit of the FMMs is superior for the modeling of maternity expenditures in the presence of extreme or outlying cost values.

**KEYWORDS:** Maternal health care expenditures; Statistical model; Generalized linear model; Gamma distribution; Log link; Outlier; Residual; Finite mixture model; Akaike Information Criteria.

**ABBREVIATION:** GLM: Generalized Linear Model; FMMs: Finite Mixture Models; AIC: Akaike Information Criterion.

## INTRODUCTION

It has been said that "beauty is in the eye of the beholder", but, so too is an investigator's understanding of phenomena a function of the lens that he or she uses to look at data. The estimates, inferences and/or conclusions that one draws from data are highly related to the way that the data are analyzed, modeled and presented. If an analyst makes a particular assumption about the prevailing characteristics of a data set and these conditions are absent, it is to no one's surprise that the subsequent downstream estimates, inferences and conclusions are at best, imprecise; at worst, erroneous.

Which analysis or model one uses to study health care costs has been the subject of debate amongst expert analysts in health outcomes research and resulted in numerous recom-

mendations.[1-7] Regardless of the philosophical position that one decides to take with respect to the analysis and/or modeling of health care costs, a critical part of the endeavor is to study the adequacy of the underlying assumptions that serve as the basis for this activity.[8,9] Without such an examination, the results derived from the data can be open to criticism and skepticism.

It has been the experience of the author, after over 26 years as a data analyst in many diverse areas of biomedical research, that the data characteristics necessary for an analysis or model to perform correctly are not frequently studied, but often assumed to be true and in some sense "robust" against departures from said features in the sample or samples. Moreover, it is also the author's contention that the verification of the requisite assumptions for an analysis or modeling exercise are often not shared with the reviewers and readers of published medical research.

A common model used to study total health care costs over many disease indications is a generalized linear model (GLM), assuming that the log of the mean costs describes the set of predictors or covariates in a linear fashion and that a gamma probability model adequately describes the distribution pattern of the observed data: its central tendency, spread, shape, etc. The author has employed such a model under many circumstances after checking that its assumptions were appropriate. However, what happens if the accepted GLM does not adequately represent the behavior of the total health care costs under study? Do other approaches exist to accommodate these departures for the requisite underlying statistical theory?

One circumstance where the standard GLM model has the potential to perform less than optimally is in the presence of extreme cost values, whether it is skewing to the high end (right), to the low (left), or in both the upper and lower range of the data set (longer "tails"); i.e., in a range than would not be reasonably predicted by the model. Such a circumstance can occur in disease indications where complications and/or particular co-morbidity patterns have the potential to increase treatment cost and expand the cost range dramatically.[10-14] The medical and pharmacy claims for women who are experiencing their first pregnancy are one such example where complications related to pregnancy can produce claims with high costs that expand or skew the cost distribution to a degree not anticipated by the conventional model.

In the face of extreme or outlying values, an analyst has a number of practices that he or she may engage in to reduce their influence or leverage on the chosen model. One approach is to analyze the data with the extreme values in and out of a model as a form of "sensitivity" analysis to see how the results vary by using these values and then removing them from the analysis. If the outlying values greatly influence the modeling results, it is common for analysts to examine the data for "assignable cause" (i.e., Are these extreme costs related to patients with unusual clinical characteristics that set them apart from most of the other

patients?), subject the extreme values to some form of outlier test[15] or compare them with a known reference range.[16] Any one of these approaches can result in a loss of sample size because the investigators may decide to discard these patients from consideration after judging their costs relative to a frame of reference related to clinical experience ("assignable cause"), a cutoff point in a statistical test (outlier test) or a reference range gleaned from the related medical literature. This decreased sample size can then have downstream effects on the operating characteristics of statistical significance tests. Also, the outlying cases may have an important place in the context of the investigation, scientifically.[17,18]

Another approach to managing the outlying or extreme values involves the use of robust statistical methods to "downweight" their influence.[19] These methods are useful if extreme values exist in both the lower and upper ranges of the cost distribution. However, techniques like the application of a Winsorized or trimmed mean to remove outliers lose their statistical optimality[20] (e.g., unbiasedness) when a distribution is asymmetrical, which is often the case for cost data.[21-23]

An approach that may be employed, whether data are in the upper, lower, or both extreme ends of the cost range, and without having to resort to the removal of patient data from a sample is the application of a *finite mixture model*. A finite mixture model (FMM) consists of two or more underlying assumed distributions, such that each contributes to the understanding of overall data pattern a specific proportion of the time, with the sum of the proportions totaling to 1 or 100%.[24,25] For the sake of simplicity, suppose that the cost distribution can be described by two different probability structures. Then the contribution of one is $p$% (e.g., say, 25%) and the other is (100-p)% (in this example, that would be [100-25]% =75%). If $f$ represents the overall probability density function for the data, it could be described by:

$$f = 0.25f_1 + 0.75f_2, \qquad (1)$$

where $f_1$ and $f_2$ are referred to as "component" probability density functions.

To further help the reader develop an intuition for the concept, suppose that we have an FMM consisting of underlying component distributions that are bell-shaped, Gaussian or normal. Figure 1 shows an example of an FMM, where 25% is from a normal with mean = 2.5 and standard deviation = 5 and the remaining three-quarters is from a normal with mean = 0 and standard deviation = 1.

Figure 1 shows an example of a *homogeneous* mixture. It is homogeneous because the two component distributions are bell-shaped, Gaussian, or normal. It is also possible to have a *heterogeneous* mixture of two different underlying distributions. Figure 2 shows an example of a mixture of 50% normal and 50% gamma.
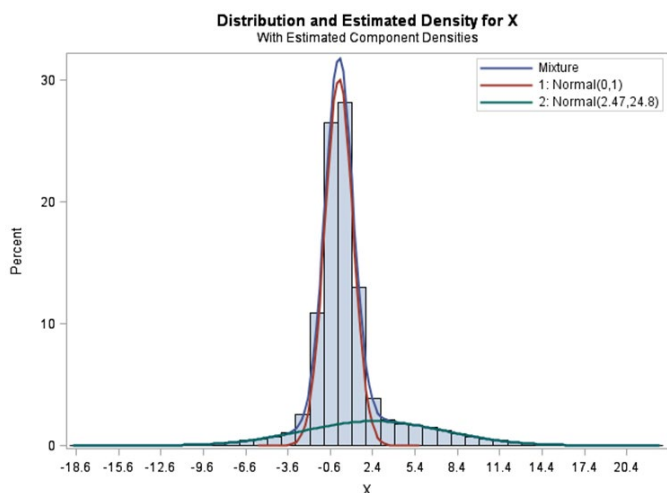
**Figure 1:** An example of a mixture of two normally distributed measurements.
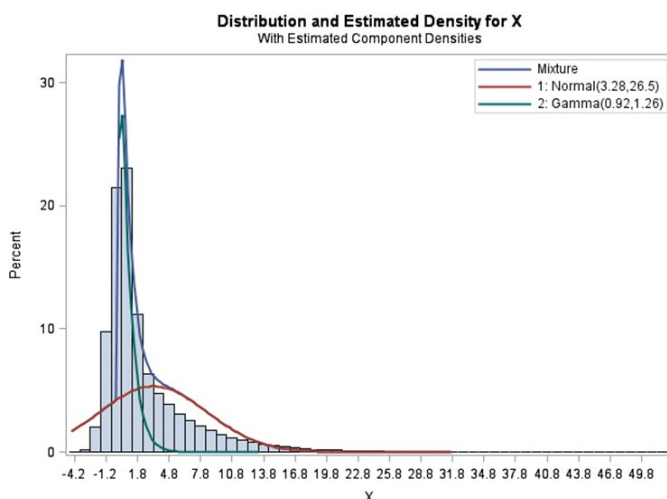


**Figure 2:** An example of a mixture of two measurements of different component distributions.

The FMM approach also allows an analyst to perform regression modeling; e.g., one could look at a mixture of cost values and see how well a set of co-morbidity or demographic characteristics predict the mean costs.

One the surface, FMM may seem like an interesting statistical or mathematical curiosity, but this approach to study cost data has been applied with success in various areas of clinical research.[26-29] Given that probability of a clinical complication during pregnancy is not zero, if one were to look at a large enough data set, he or she might end up with patients whose costs are on the order of millions of dollars. Certainly, such costs would exert influence on a regression model where most of the costs are within a lower anticipated range.

The advantage of the FMM approach is that the analyst does not really ever have to worry about whether a value is an outlier or not. If sufficient outliers exist, they may be modeled as part of a smaller, or minority component distribution in a cost model. The burden of examining the values for assignable cause

is eliminated. Extreme data can occur at the high end, low end, or both in the cost range and studied with an FMM without the limitations imposed by robust statistical methods such as trimming.

The work in this paper is, to the knowledge of the author, the first application of an FMM to cost data in the setting of pregnancy. The author will share an example of a real data set that has properties that will demonstrate the benefits of the FMM approach when a portion of the costs are extremely low and high. These values effectively ruin the fit when standard cost models are employed.

### DATA SOURCE

The data set for analysis consisted of 322,107 pregnant women aged 15-49 years using de-identified medical and pharmacy claims from the Truven Health MarketScan® Commercial Claims and Encounters database incurred between January 1, 2007 and December 31, 2011. The total health care costs were calculated from the date of the first pregnancy-related claim through to 3 months post-delivery, adjusted to 2011 dollars.

### METHODS

The constructed data set was examined for its fidelity to a set of assumptions typically used in health care cost models that the costs could be adequately modeled with a gamma distribution.[30,31] First, the distribution of the costs was fit assuming a gamma distribution and compared with a kernel density estimate of the data distribution, an approach using more general assumptions and not imposing nor assuming a particular form for the cost data. A plot containing a histogram with the two superimposed distributions was used to make an initial assessment of the adequacy of the gamma assumption. Second, these data were fit with a small set of co-morbidity predictors and the model residuals were subject to an examination for the aptness of the gamma assumption a second time and for consideration of the appropriateness of the standard belief that the log of the mean costs could be related to the predictors in a linear fashion (i.e., that assuming a log link was plausible). After these two assessments, the data were refit using both heterogeneous and homogenous finite mixture models to compare their fit with the gamma assumption. The models with the lowest Akaike Information Criterion (AIC) were considered to be a better fit or description of the observed costs.

All data analysis, models and graphics were conducted using various SAS v 9.4 (TS1M2) procedures.

### RESULTS

Figure 3 shows a histogram with two superimposed curves. One is a kernel density estimate of the distribution of the costs with minimal assumptions about its shape, scale, etc.

The second curve was fitted assuming that the costs follow the standard assumption that they may be described by a gamma probability model. Figure 3a shows the data in its entirety. As the presence of patients with costs in the millions of dollars are part of the data set, the plot was run a second time with the data truncated at $100K. Figure 3b was only created for the sake of illustration purposes in this instance.

Figure 3b illustrates the main shortcoming of the choice of a gamma distribution for use as a cost model in this setting. The blue curve in Figure 3b is how an assumed gamma distribution would fit if it were used to describe the maternal cost data. Note that if an analyst were to use this model for these data the assumed distribution would over-predict how often the costs were on the lower and higher ends of the distribution and under-predict how often the costs would be around the mode, or most frequent value.
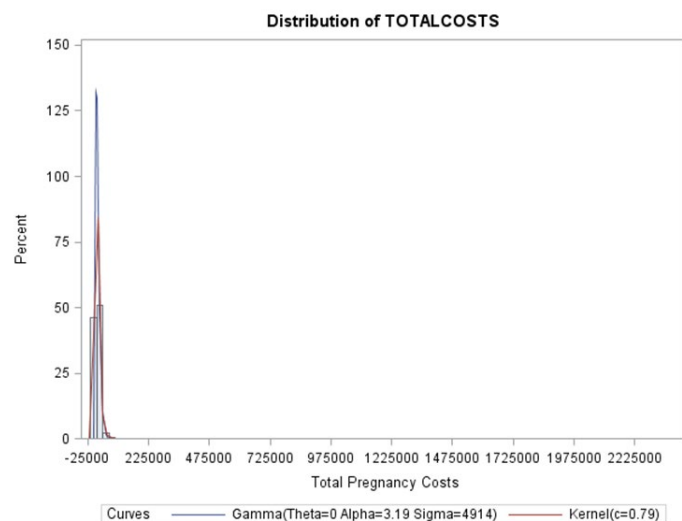
The initial plausibility of the appropriateness of the standard cost model is somewhat in doubt. The other feature of the standard model assumes that the predictors may be adequately described by a linear relationship with the log of the mean costs. In a GLM this assumed relationship is called the "link" between the mean costs and predictors. In Figure 4 the model was fit with a small set of 9 binary co-morbidity variables (presence v absence) and the assumption was checked by a special type of plot called a *cumulative residual plot*. Details regarding the statistics and construction of the plot may be found in other sources.[32,33] The main point of this plot is that the behavior of the actual data (dark blue solid curve – indicated by solid arrows) should fall within a set of bounds found (dashed, lighter blue region) by re-sampling the cost data several times. For the most part, the actual data do well and fall within the bounds, with the exception of the area in the purple-shaded boxes.
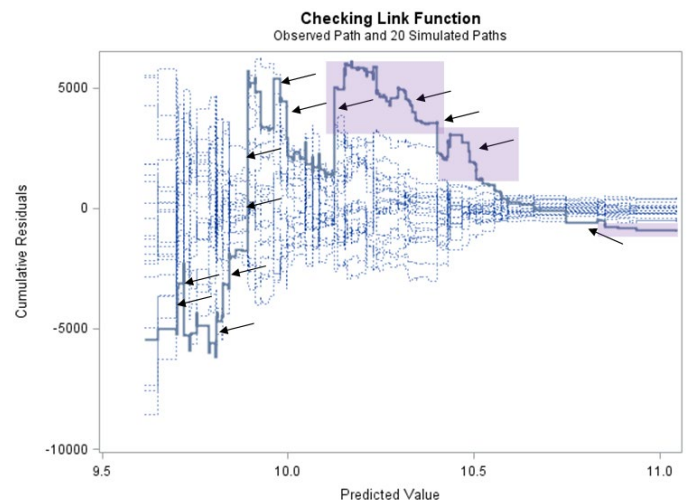


**Figure 3a:** A first look at the maternity total cost data.



**Figure 4:** Cumulative residuals for a model assuming a log link.

Now, consider a simple approach to correct the shortcomings of the gamma and log link model. A standard regression technique is to try and look at the log of each patient's individual costs. A plot of the distribution of the log-transformed costs is shown in Figure 5.

In Figure 5 the log costs have been plotted and a superimposed kernel and normal curve were drawn over the data. The log transformed data that roughly follow a normal or bell curve are called *log normally distributed* data. This approach suffers from the similar problems to the gamma model; however, the pattern of over-estimation is reverse (for smaller values, the white area between the blue and red curves is larger on the left-hand side than was the case for the assumed gamma model in Figure 3b).

The short comings of these two models suggest that some heterogeneity exists in the cost data; *viz.*, a single probability model or distribution will not adequately describe the behavior of the entire data set. Let's now attempt to fit an FMM to see if the description of such a model is superior to the ones
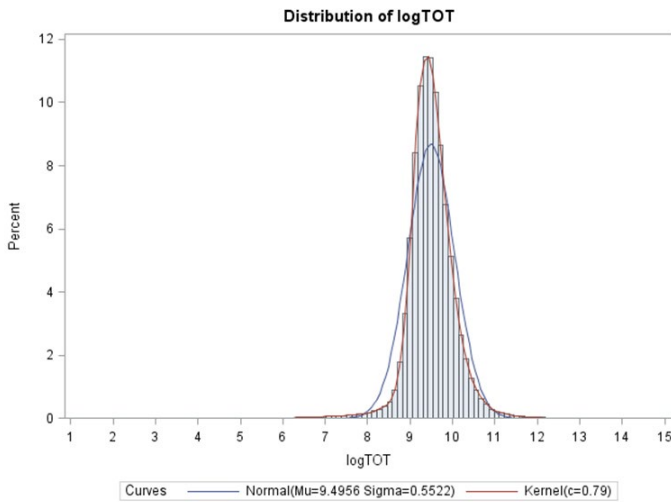


**Figure 3b**: A second look at the maternity total cost data – Truncated at $100K.

previously employed.



**Figure 5:** The distribution of the log-transformed costs.



**Figure 7a:** A heterogeneous mixture of a normal and log-normal probability model.

Figure 6 consists of various FMM fit to the maternal cost data on the log scale. First attempts were with homogeneous FMMs. Recall that a homogeneous FMM consists of two or more probability models or distributions, but only differing in the values of their parameters, like the example in Figure 1 where both are bell-shaped curves, only differing in their means and standard deviations.

In Figure 6 the fitted model appears to be superior to that of the earlier single gamma or normal distributions (the latter fit on the log-transformed data). Using a SAS® procedure called PROC FMM, the results suggest that the first component (slightly larger mean, variance about 10x as larger) describes about 90% of the mixture, while another component describes the remaining 10%.
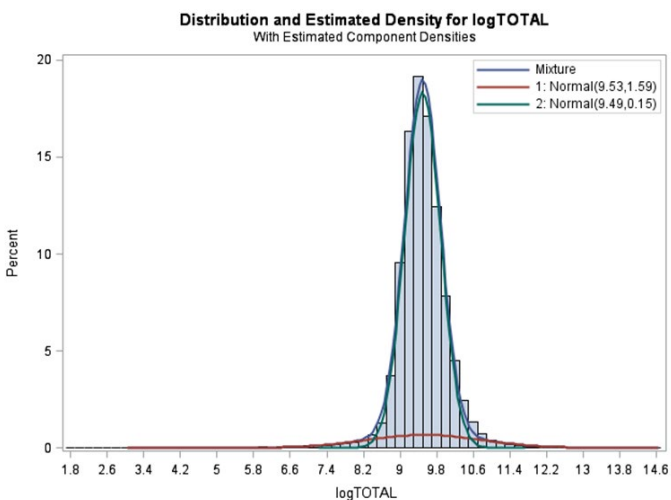


**Figure 7b:** A heterogeneous mixture of a normal and gamma probability model.



**Figure 6:** Homogeneous FMM fit to the maternal cost data using two normal distributions.



**Figure 7c:** A heterogeneous mixture of a normal and gamma probability model.

Figure 7 shows attempts to use various heterogeneous FMMs to describe the data. Only two component models were fit for this analysis.
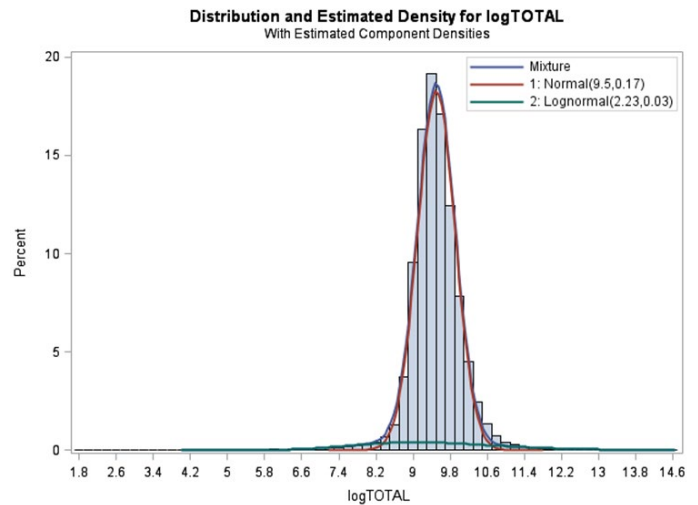
After examining Figures 6 and 7 it seems that all of these models are an improvement over the standard models, *visually*. Indeed, comparing one plot with another seems like an-

other "beauty contest" to find the best model. Fortunately, an analyst may use a statistic to differentiate between the models for a more objective choice.

The Akaike Information Criteria[34] (AIC) is a statistic that is often used to compare various models and their fit to data. Crudely, the AIC is a measure of the balance between underfitting a data set versus over-fitting by using too large of a number of parameters or a more complex model.[35,36] The AIC may be thought of as the sum of two quantities: the lack of fit for the model and a "penalty" for potential unreliability introduced by a more complicated model.[37] In a set of competing models, the model with the lowest AIC is considered to be the optimal one; in the context of FMM, AIC may be used to guide the selection of an optimal model.[38] Table 1 shows the AIC values for some models considered for the fit of the maternal cost data.

By virtue of the AIC values, the homogeneous FMM based on the log costs assuming two underlying normal probability models would seem like the best model to describe the maternal costs.

## DISCUSSION

Modeling a phenomenon involves the process of reducing it down to a set of features by detail abstraction and assuming that underlying conditions persist. Thus, for cost data, the analyst assumes a probability model that, in itself, describes the behavior of the cost data distribution and presumes that it has a specific shape, scale, or other characteristics that may be readily identifiable, mathematically (i.e., may be specified by a formula). When an analyst picks a model without checking its adequacy, he or she is imposing features on the data or "viewing it through a lens" that might distort reality because of assumptions that cannot be supported upon closer examination of the data. All estimates and inferences derived under these faulty or unchecked assumptions can be imprecise, or in error, respectively. For medical indications where a great deal of underlying heterogeneity exists, FMMs have been shown to more adequately describe the data characteristics and reflect the reality of the cost data than standard models. The author contends that their application in the study of costs related to pregnancy may be another area where the FMM approach more adequately describes the data, especially if specific medical complications exist and occur infrequently, but often enough to undermine the appropriate use of a more traditional, single distribution or probability model. His study is a first case analysis of the potential for FMMs in modeling of costs in gynecology and obstetrics.

In an age when increasing health care costs are falling under greater scrutiny by payers, a sensible starting place for greater understanding is to check the way that data are used to make estimates of incremental costs or cost-effectiveness. Are the results derived from a model based on a defensible or rational approach that is faithful to the features of the data? Are the prevailing conditions observed in the clinical environment addressed in the model? Equipped with the appropriate tools, the analyst is allowed to share a more accurate vision of the behavior of costs associated with pregnancy, which may lead to a more precise estimates of incremental costs or cost-effectiveness, and ultimately, serve the best interests for the treatment and care of pregnant women.

## CONFIDENTIALITY/CONSENT STATEMENT

The data used for this study did not involve the interaction or interview with any subjects and the data does not include any individually identifiable data (e.g. does not include names, addresses, social security or medical record numbers or other obvious identifiers) and as such is not research involving human subject as defined at 45 CFR 46.102(f)(2). Furthermore, this study used existing fully de-identified and the investigator(s) cannot be identified, directly or through identifiers linked to subjects

| Model | Model Type | AIC |
|---|---|---|
| Gamma on Untransformed Costs | Traditional Single Probability Model | 6690934 |
| Normal on Log-Transformed Costs | Traditional Single Probability Model | 531525 |
| 10% Normal(mean=9.53, variance=1.59) + 90% Normal(mean=0.15, variance =0.15) on Log-Transformed Costs | Homogeneous FMM | 459427 |
| 92% Normal(mean=9.50, variance=0.16) + 8% Gamma(intercept=2.25, scale= 44.04) on Log-Transformed Costs | Heterogeneous FMM | 461699 |
| 98% Normal(mean=9.51, variance=0.21) + 2% Exponential(intercept=2.18) on Log-Transformed Costs | Heterogeneous FMM | 484566 |
| 26% Normal(mean=9.9, variance=0.35) + 74% Weibull(intercept=2.26, scale=0.04) on Log-Transformed Costs | Heterogeneous FMM | 481570 |

**Table 1:** AIC values comparing single, homogeneous and heterogeneous FMMs.

and as such is exempt from 45 CFR 46.101(b)(4) from all 45 CFR part 46 requirements. Consequently, IRB approval is not required.

## REFERENCES

1. Diehr P, Yanez D, Ash A, Hornbrook M, and Lin DY. Methods for analyzing health care utilization and costs. *Annual Rev Pub Health*. 1999; 20: 125-144. doi: 10.1146/annurev.publhealth.20.1.125

2. Blough DK, Ramsey SD. Using generalized linear models to assess medical care costs. *Health Services and Outcomes Res Method*. 2000; 1(2):185-202. doi: 10.1023/A:1012597123667

3. Manning WG, Mullahy J. Estimating log models: to transform or not to transform? *J Health Econ*. 2001; 20: 461-494. doi: 10.1016/S0167-6296(01)00086-8

4. Buntin MB, Zaslavsky AM. Too much ado about two-part models and transformation? comparing methods of modeling medicare expenditures. *J Health Econ*. 2004; 23: 525-542. doi: 10.1016/j.jhealeco.2003.10.005

5. Montez-Rath M, Christiansen CL, Ettner SL, Loveland S, Rosen AK. Performance of statistical models to predict mental health and substance abuse costs. *BMC Medical Res Method*. 2006; 6: 53 doi: 10.1186/1471-2288-6-53

6. Mullahy M. Econometric modeling of health care costs and expenditures: a survey of analytical issues and related policy considerations. *Med Care*. 2009; 47(7): S104-S108. doi: 10.1097/MLR.0b013e31819c9593

7. Mihaylova M, Briggs A, O'Hagan A, Thompson SG. Review of statistical methods for analysing healthcare costs and resources. *Health Econ*. 2010; 20(8): 897-916. doi: 10.1002/hec.1653

8. McCullagh P, Nelder, JA. Generalized linear models, 2/e. New York, NY USA: Chapman & Hall; 1991.

9. Myers RH, Montgomery DC, Vining GC. Generalized linear models in engineering and the sciences. New York, NY USA: John Wiley & Sons, Inc.; 2002.

10. Williams MD, Braun LA, Cooper LM, et al. Hospitalized cancer patients with severe sepsis: analysis of incidence, mortality, and associated costs of care. *Critical Care*. 2004; 8(5): R291-R298 doi: 10.1186/cc2893

11. Carls GS, Lee DW, Ozimnkowski RJ, Wang S, Gibson TB, Stewart E. What are the total costs of surgical treatment for uterine fibroids? *J Women's Health*. 2008; 17(7): 1119-1132. doi: 10.1089/jwh.2008.0456

12. Brem H, Maggi J, Nierman D, et al. High cost of stage IV pressure ulcers. *Am J Surg*. 2012; 200(4): 473-477. doi: 10.1016/j.amjsurg.2009.12.021

13. Cardozo ER, Clark AD, Banks NK, Henne MB, Stegmann BJ, Segars JH. The estimated annual costs of uterine leiomyomata in the United States. *Am J Obset Gynecol*. 2012; 206(3): 211.e1-211.e9. doi: 10.1016/j.ajog.2011.12.002

14. Yeaw J, Halinan S, Hines D, et al. Direct medical costs for complications among children and adults with diabetes in the US commercial payer setting. *Appl Health Econ Health Policy*. 2014; 12(2): 219-230. doi: 10.1007/s40258-014-0086-9

15. Barnet V, Lewis T. Outliers in statistical data. New York, NY USA: John Wiley & Sons; 1979.

16. Harris EK, Boyd JC. Statistical bases of reference values in laboratory medicine. New York, NY USA: Marcel Dekker, Inc.; 1995.

17. Kandel R. Our changing climate. New York, NY USA: McGraw Hill; 1991.

18. Vitello P. Joseph Farman, 82 is dead; discovered ozone hole. New York Times. May 20, 2013: B9.

19. Maronna RA, Martin RD, Yohai VJ. Robust statistics: theory and methods. New York, NY USA: John Wiley & Sons; 2006.

20. Rivest L-P. Statistical properties of winsorized means for skewed distributions. *Biometrika*. 1994; 81(2): 373-383.

21. Dunn G, Mirandola M, Amaddeo F, Tansella M. Describing, explaining or predicting mental health care costs: a guide to regression models: methodological review. *British J Psych*. 2003; 183: 398-404. doi: 10.1192/bjp.183.5.398

22. Griswold M, Parmigiani G, Potsky R, Lipscomb J. Analyzing health care costs: a comparison of statistical methods motivated by medicare colorectal cancer charges. *Biostatistics*. 2004; 1(1): 1-23.

23. Başer O. Modeling transformed health care costs with unknown heteroskedasticity. *App Econ Res Bull*. 2007; 01: 1-6.

24. Kessler D, McDowell A. Introducing the FMM procedure for finite mixture models. https://support.sas.com/resources/papers/proceedings12/328-2012.pdf ; Accessed July 9, 2015.

25. McLachlan G, Peel D. Finite mixture models. New York, NY USA: John Wiley & Sons; 2000.

26. Deb P, Trivedi PK. Demand for medical care for the elderly: a finite mixture approach. *J Appl Econometrics*. 1997; 12: 313-336.

27. Deb P, Holmes AN. Estimates of use and cost of behavioral health care: a comparison of standard and finite mixture models. In: Jones A and O'Donnell O. ed. Econometric analysis of health data. Chichester, West Sussex, UK: John Wiley and Sons; 2002: 87-99.

28. Lourenço OD, Ferreira PL. Utilization of public health centres in Portugal: effect of time costs and other determinants. finite mixture models applied to truncated samples. *Health Econ*. 2005; 14: 939-953. doi: 10.1002/hec.1046

29. Rein DB. A matter of classes: stratifying health care populations to produce better estimates of inpatient costs. *Health Serv Res*. 2005; 40(4): 1217-1233. doi: 10.1111/j.1475-6773.2005.00393.x

30. Pierce DA, Schafer DW. Residuals in generalized linear models. *J Amer Stat Assoc*. 1986; 81(396): 977-986. doi: 10.2307/2289071

31. Dunteman GH. Introduction to generalized linear models. London, UK: Sage Publications; 2006.

32. Lin DY, Wei LJ, Ying Z. Model-checking techniques based on cumulative residuals. *Biometrics.* 2002; 58: 1-12. doi: 10.1111/j.0006-341X.2002.00001.x

33. The GENMOD procedure. SAS/STAT user's guide, 2/e. http://support.sas.com/documentation/cdl/en/statug/63033/HTML/default/viewer.htm#statug_genmod_sect060.htm; Accessed July 9, 2015.

34. Akaike H. Likelihood of a model and information criteria. *J Econometrics.* 1981; 16: 3-14.

35. Bozdogan H. Akaike's information criteria and recent developments in information complexity. *J Math Psychol.* 2000; 44: 62-91. doi:10.1006/jmps.1999.1277

36. Anderson DR, Burnham KP, White GC. Comparison of akaike information criteria and consistent akaike information criteria for model selection and statistical inference from capture-recapture studies. *J App Stat*. 1998; 25(2): 263-282. doi: 10.1080/02664769823250

37. Mutua FM. The use of akaike information criterion in the identification of an optimum flood frequency model. *Hydrological Sciences J.* 1994; 39(3): 235-244.

38. The FMM Procedure. SAS/STAT User's Guide. 13.2. http://support.sas.com/documentation/cdl/en/statug/67523/HTML/default/viewer.htm#statug_fmm_overview.htm . Accessed July 9, 2015.