

Minimum Bias, GLMs and Credibility in the Context of Predictive Modeling

Christopher Gross and Jonanthan Evans

Abstract:

When predictive performance testing, rather than testing model assumptions, is used for validation, the needs for detailed model specification are greatly reduced. Minimum bias models trade some degree of statistical independence in data points in exchange for statistically much more tame distributions underlying individual data points. A combination of multiplicative minimum bias and credibility methods for predictively modeling losses (pure premiums, claim counts, and/or average severity, etc.) based on explanatory risk characteristics is defined. Advantages of this model include grounding in longstanding and conceptually lucid methods with minimal assumptions. An empirical case study is presented with comparisons between multiplicative minimum bias and a typical generalized linear model (GLM). Comparison is also made with methods of incorporating credibility into GLM.

Keywords: predictive modeling, minimum bias, credibility, ratemaking, generalized linear models

1. INTRODUCTION

As predictive models that relate losses (pure premiums, claim counts, and/or average severity, etc.) to explanatory risk characteristics become ever more commonplace, some of the practical problems that frequently emerge include:

- Models often use complex techniques that are effectively “black boxes” without a lucid conceptual basis.
- Models may require very detailed parametric or distributional assumptions. Invalid assumptions may result in biased parameters.
- A highly Frequentist approach, usually involving Maximum Likelihood Estimation (MLE), can lead to overfitting sparsely populated data bins.

Some longstanding methods can be combined to overcome these problems:

- Minimum Bias Iterative fitting of parameters is simple, longstanding in practice, and non-parametric in specification.
- Credibility methods are similarly simple and longstanding. Credibility directly solves the sparse bin problem.

Most importantly, properly done predictive testing, in contrast with testing model assumptions, makes highly detailed model specification generally unnecessary.

1.1 Research Context

The minimum bias criteria and iterative solution methodology were introduced by Bailey and Simon in [2] and [3]. Brown in [5] substituted the minimum bias criteria with MLE of Generalized Linear Models (GLM), an approach further explored by Mildenhall in [10]. Venter in [13] further discusses credibility issues related to minimum bias methods. The basic contemporary reference on credibility methods is Klugman, S., et al. [9]. Nelder and Verrall in [11] and Klinker in [8] discuss incorporating random effects into GLM to implement credibility adjustments. Brosius and Feldblum provide a modern practical guide to Minimum Bias Methods in [4]. A similar practical guide to GLM is provided by Anderson, et al. in [1]. A demonstration of predictive model fitting and testing can be found in Evans and Dean [6], particularly the predictive testing methods that will be used in this paper. “Gibbs Sampling” is a term we will use for Markov Chain Monte Carlo (MCMC) methods, as these are implemented using Gibbs Sampling software, such as BUGS, WinBUGS, or JAGS. Scollnik in [12] introduces MCMC. Particularly relevant to this paper is the recent book on predictive modeling for actuaries Frees, E., et al. [7]. This book contains very detailed information on GLM, particularly incorporating credibility through Gibbs Sampling. This paper represents in a certain sense an opposite perspective from [7] and [12], by emphasizing very simple models combined with rigorous predictive testing as described in [6].

1.2 Outline

The remaining sections of this paper are:

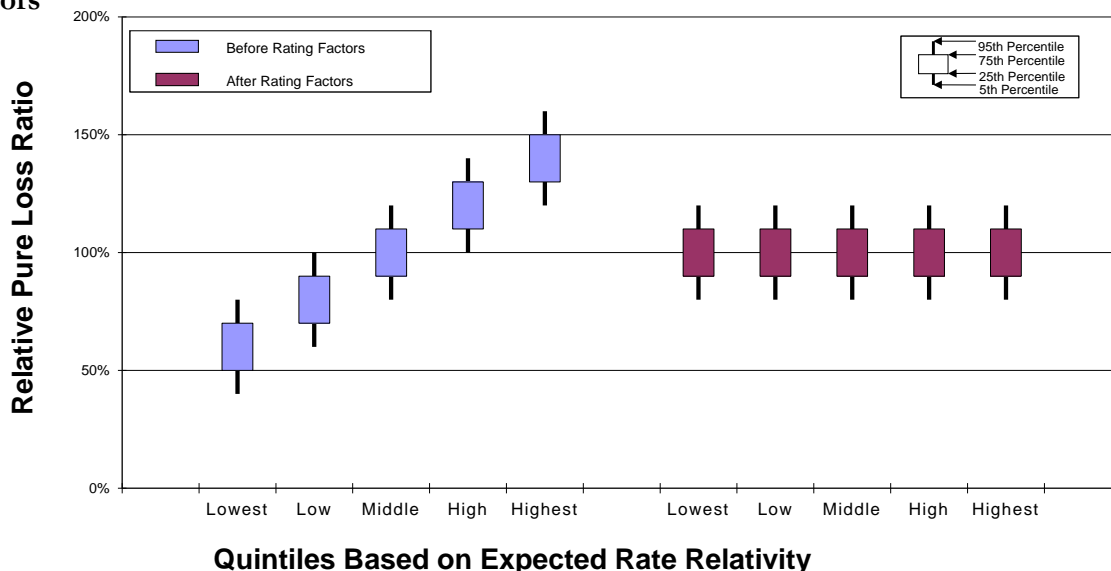
2. Predictive Performance as the Modeling Objective
 3. Multiplicative Minimum Bias Iteration
 4. Incorporating Credibility
 5. Anchoring And Iteration Blending For Practical Iterative Convergence
 6. Testing of Individual Explanatory Variables
 7. Empirical Case Study
 8. Summary Discussion
- Appendix A. Details of Empirical Case Study
- Appendix B. Gibbs Sampling Model Code

2. PREDICTIVE PERFORMANCE AS THE MODELING OBJECTIVE

Traditionally, statistical models tend to use the same data for both fitting and validation. Validation tends to involve testing of model assumptions. For example, a linear regression of the form $Y = mX + b + \xi$, where $\xi \sim \text{Normal}(0, \sigma^2)$, might be fit, using least squares, to a set of data points (x_i, y_i) , $i = 1, \dots, n$. Validation tests would check to verify that the residuals ξ_i are normally distributed with constant variance and are independent of x_i , y_i , and each other. Hypothesis tests would then be performed to confirm that the probability is sufficiently remote that the actual data set would result if $m = 0$ or $b = 0$ (null hypotheses). This framework relies on detailed assumptions, without which validation testing would not be possible.

Modern predictive models split available data into multiple sets for separate fitting and validation. In the previous example, the parameters m and b might be fit to the points (x_i, y_i) , $i = 1, \dots, k$, using any method, and then tested on the points (x_i, y_i) , $i = k+1, \dots, n$. The test would only be concerned with how well $\hat{y}_i = \hat{m}x_i + \hat{b}$ predicts y_i for the test set. A bootstrap quintile test might be used, where the validation points are sorted by the value \hat{y}_i into 5 equal-sized groups. The average value of y_i should ascend with the quintile groups and for each group the average value of y_i should be close to the average value of \hat{y}_i . Figure 1 is a hypothetical example of a quintile test, with bootstrap confidence intervals added, as described by Evans and Dean in [6], validation of rating factors. Note, the assumption $\xi \sim \text{Normal}(0, \sigma^2)$ and other implicit assumptions of linear regression are unnecessary here.

Figure 1. Hypothetical Example of Bootstrap Quintile Test Predictive Validation of Rating Factors



In practice, predictive modelers often split data into three or more sets (i.e., training, testing, and validation), but only the distinction between two separate data sets for fitting and validation will be covered in this paper.

In the predictive framework, detailed model assumptions are not necessary. A model, even if its assumptions seem unjustified or erroneous, is valid as long as it performs well at predicting outcomes for data that were not used to fit its parameters. This comes with the caveat that care must be taken that both the fitting and validation data should be representative of – effectively random samples of – the loss process. For example, predictive testing might be misleading if both the fitting and validation data occurred in a single year influenced by a somewhat rare catastrophe, such as a hurricane.

3. MULTIPLICATIVE MINIMUM BIAS ITERATION

Suppose the basic data available consists of actual losses $L_{i_1, \dots, i_n} \geq 0$ and exposures $P_{i_1, \dots, i_n} \geq 0$, ($P_{i_1, \dots, i_n} = 0 \Rightarrow L_{i_1, \dots, i_n} = 0$) where $i_j = 1, \dots, n_j$ indexes the individual classes within the classification dimension j and i_1, \dots, i_n denotes the cell corresponding to the intersection of a single class in each classification dimension. Also the total exposure in any class is positive, $\sum_{i_j=k} P_{i_1, \dots, i_n} > 0$, otherwise it would make sense to exclude the class entirely from estimating rating parameters. A multiplicative minimum bias model assumes that $L_{i_1, \dots, i_n} = B_{i_1, \dots, i_n} + P_{i_1, \dots, i_n} \prod_{j=1, \dots, n_j} X_{j, i_j}$. The parameters X_{j, i_j} are fit with the goal of minimizing some bias function, or functions, of the residual errors B_{i_1, \dots, i_n} .

The minimum bias goal is that the sum of the residual errors for each class $\sum_{i_j=k} B_{i_1, \dots, i_n}$ should be 0. A corresponding iterative sequence of parameter estimates can be formed whose convergence corresponds to convergence to the goal:

$$\begin{aligned}
 X_{j,k,1} &= 1 \\
 X_{j,k,t+1} &= \frac{\sum_{i_j=k} L_{i_1, \dots, i_n}}{\sum_{i_j=k} P_{i_1, \dots, i_n} \prod_{l \neq j} X_{l, i_l, t}} \quad (3.1)
 \end{aligned}$$

The effective sample is now $\sum_{j=1, \dots, n} n_j$ data points with values $\sum_{i_j=k} L_{i_1, \dots, i_n}$, which reduces to $\sum_{j=1, \dots, n} n_j - (n-1)$ linearly independent numbers. There is a corresponding $(n-1)$ dimensional degeneracy in the parameters. If the parameters X_{k, i_k} are multiplied by a constant $c > 0$ and the parameters X_{l, i_l} are divided by c , where $0 \leq k < l \leq n$, then $\prod_{j=1, \dots, n_j} X_{j, i_j}$ will be unchanged.

The Central Limit Theorem implies that the distribution of $\sum_{i_j=k} L_{i_1, \dots, i_n}$ can be expected to more closely resemble a Normal distribution, with a generally lower coefficient of variation than the individual cell values L_{i_1, \dots, i_n} . However, whereas the cellular values L_{i_1, \dots, i_n} can reasonably be assumed to be statistically independent of each other, the aggregated values $\sum_{i_j=k} L_{i_1, \dots, i_n}$ include many statistical dependencies since there is an overlap of cells between classes in different dimensions. So, a tradeoff is made for a minimum bias iteration model. Statistical independence of sample data points, a desirable property, is partially sacrificed in exchange for the benefit of a more Normal distribution, generally having a lower coefficient of variation than the distributions underlying each sample data point. This taming of the distribution of data points means that it becomes less necessary to specify the distribution of the individual cellular loss values, or as may be the case the distributions of individual loss observations within the cells, as would be necessary for a GLM.

Example 1

Suppose there are three classification dimensions, each with 10 classes, resulting in 1,000 individual cells. We can expect about 100 times as much data underlying each class as for each cell, and correspondingly an average coefficient of variation by class that is only about 10% as much as by cell. Two classes in different dimensions overlap in 10 cells and thus actual losses between them will have a correlation coefficient of about 10%.

Multiplicative minimum bias effectively aims toward the same parameters estimates as a GLM with a logarithmic link function and Poisson likelihood function. The logarithmic link converts the sum of linear explanatory factors into a multiplicative product of their exponentials. The Poisson likelihood leads to equations for MLE that correspond to a fixed limit point of the minimum bias iteration, as pointed out by Brown in [5].

However, the Poisson distributional assumption is usually unrealistic and not a part of the minimum bias model. Data are generally not restricted to integer values. The Poisson coefficient of variation (CV) is not scale independent (it is 10 times greater when applied to dollar amounts versus when applied to the same amounts measured as pennies) and implodes for large nominal means (mean of 1,000,000 implies a CV of 0.1%). So, the Poisson assumption is important only in the optimization equations it implies for MLE.

4. INCORPORATING CREDIBILITY

Credibility adjustments $0 \leq Z_{j,i_j} \leq 1$ can be easily and directly incorporated into the iteration equations:

$$\begin{aligned}
 X_{j,k,1} &= 1 \\
 X_{j,k,t+1} &= (1 - Z_{j,k}) + Z_{j,k} \frac{\sum_{i_j=k} L_{i_1, \dots, i_n}}{\sum_{i_j=k} P_{i_1, \dots, i_n} \prod_{l \neq j} X_{l, i_l, t}} \tag{4.1}
 \end{aligned}$$

Note, other than the constraint of the interval $[0, 1]$, nothing has been specified about the determination of $Z_{j,i}$. There are many possibilities for Z_{j,i_j} , including functions of the sum of exposure $P_{j,k} = \sum_{i_j=k} P_{i_1, \dots, i_n}$. The ultimate test will be the predictive performance of the final model regardless of whether $Z_{j,i}$ itself satisfies any traditional goals of credibility theory, such as limiting fluctuation or greatest accuracy.

For GLM, the basic and common protection against fitting parameters to data that is not credible is to throw away explanatory variables whose parameters are not statistically distinct from 0, those variables with high p-values.

To add a true credibility, or “shrinkage”, adjustment is complicated. The two main approaches are:

1. General Linear Mixed Models. At least some rating factors are assumed to be random rather than fixed effects, but an MLE-like fitting method is still used. Numerical solution is rather difficult and, in practice, functions in R or procedures in SAS are used, very much as black boxes. See [7], [8] and [11] for background.
2. Bayesian Networks and Gibbs Sampling. Rating factors in each class dimension follow a prior distribution. The parameters of the prior distributions follow distributions that are very diffuse. Numerical solution is performed using a Gibbs Sampling program, such as JAGS or WinBUGS. The model itself is elaborately specified and lucid to an audience sophisticated enough read the specification. See [7] and [12] for background.

In Section 7, we will demonstrate an example of the second approach.

5. ANCHORING AND ITERATION BLENDING FOR PRACTICAL ITERATIVE CONVERGENCE

In practice the convergence of the iterative algorithms can be a problem even after the application of credibility. For one thing there is still the problem of $(n-1)$ dimensional degeneracy previously mentioned. Also, highly correlated dimensions can also contribute to non-convergence or slow convergence in practice. Other than the automatic degeneracy we will not attempt to deal with the more general convergence issue in a precise mathematical way, which appears to be an open problem for multiplicative minimum bias. From a practical point of view *anchoring* and *iteration blending* can effectively provide timely convergence.

Anchoring directly eliminates the degeneracy. One approach is to fix one of the parameters in each of $(n-1)$ dimensions to the value of 1.0; or to fix such a parameter in each of n dimensions and add a single overall base rate parameter. Another approach is to use a single overall base rate and rescale the parameters in each dimension to a weighted average of 1.0 at the end of each iteration.

Example 2

If $P = \begin{pmatrix} 1 & 1 \\ 1 & 1 \end{pmatrix}$ and $L = \begin{pmatrix} 1 & 2 \\ 3 & 4 \end{pmatrix}$ then parameter iterations will oscillate back and forth between the values $X = \begin{pmatrix} 1.5 & 3.5 \\ 2.0 & 3.0 \end{pmatrix}$ and $X = \begin{pmatrix} 0.6 & 1.4 \\ 0.8 & 1.2 \end{pmatrix}$. However, if we “anchor” one parameter at 1.00

the iterations will converge to $X = \begin{pmatrix} 1.000 & 2.333 \\ 1.200 & 1.800 \end{pmatrix}$.

Iteration blending can be implemented to accelerate convergence by modifying the iterative equations to be:

$$X_{j,k,t+1} = \alpha \left[(1 - Z_{j,k}) + Z_{j,k} \frac{\sum_{i_j=k} L_{i_1, \dots, i_n}}{\sum_{i_j=k} P_{i_1, \dots, i_n} \prod_{l \neq j} X_{l, i_l, t}} \right] + (1 - \alpha) X_{j,k,t+1} \quad (5.1)$$

where $0 < \alpha < 1$ is a selected constant blending parameter.

As an extreme illustration of correlation, let one classification dimension be replicated or made once redundant. Setting $\alpha = 0.5$ will allow the model to converge. Each one of the replicated dimensions will end up sharing equally in the observed predictive relationship, combining together to provide the appropriate prediction. In the case of full credibility, they will exactly reproduce the result from not replicating the dimension. With less than full credibility, the result will not be exactly the same from not replicating the dimension, but will be similar.

6. TESTING OF INDIVIDUAL EXPLANATORY VARIABLES

Sometimes predictive modeling techniques are used specifically to determine whether or not individual explanatory variables, or equivalently classification dimensions, are statistically significant. As mentioned earlier, when using GLM techniques, it is common to consider the p-values of the estimated parameters. These p-values are calculated under the distributional and other assumptions, such as independence of the GLM model being used.

Whether distributional assumptions are made (as with GLM) or not (as with minimum bias), tests of predictive performance can be performed and compared with and without a given classification dimension. In cases where the improvement is insignificant the dimension should be removed for the sake of parsimony.

7. EMPIRICAL CASE STUDY

The empirical data used in this case study consists of 371,123 records of medical malpractice payments obtained from the National Practitioner Data Bank. Three explanatory variables will be used

for modeling payment amounts: *Original Year*, *Allegation Group* and *License Field*. The records will be randomly split into two sets for model fitting and validation, respectively. Further details are included in Appendix A.

7.1 GLM Model Specifications

For our GLM model we will consider:

1. The logarithmic link function, which causes the fit factors to act multiplicatively.
2. Several likelihood functions: Gaussian, Poisson, Gamma, and Inverse Gaussian. These correspond to assumptions that variance σ^2 is related to mean μ as $\sigma^2 = \text{constant}$, $\sigma^2 \propto \mu$, $\sigma^2 \propto \mu^2$, and $\sigma^2 \propto \mu^3$, respectively.
3. Initially we will ignore credibility considerations, aside from reviewing p-values, and later we will use Gibbs Sampling to incorporate credibility.

7.2 Comparison of GLM and Minimum Bias Model Results

Figures 2 and 3 and Table 1 show the bootstrap quantile testing results of fitting and performance testing models. Optimal noise-to-signal estimates along the lines described in [6] suggested using 20 quantiles. Also, see [6] for details on the definitions of the test statistics. The “old statistic” test measure is the ratio of the variance of the relative average payments after rating factors are applied to the same variance before rating factors are applied, lower being better. The “new statistic” test measure is essentially the square root of the difference between these two variances, higher being better.

Although Figures 2 and 3 only correspond to the Minimum Bias fits, Table 1 demonstrates that the Log-Poisson GLM was identical to the Minimum Bias approach, and the best fitting model. In fact, we checked the individual predicted values and verified that they were numerically identical. Log-Gaussian and Log-Gamma were almost as good. The MLE for our run of Log-Inverse Gaussian failed to converge, almost certainly driven by its unrealistic variance assumption.

Figures 4 and 5 correspond to “Traditional” univariate rate relativities for the three explanatory variables. Rating factors are calculated separately and independently in each classification dimension. The Traditional method clearly performs much worse than Minimum Bias and the convergent GLMs, but is still a great improvement over no adjustment.

Figure 2. Bootstrap 20 Quantiles Test Validation of Minimum Bias Rating Factors

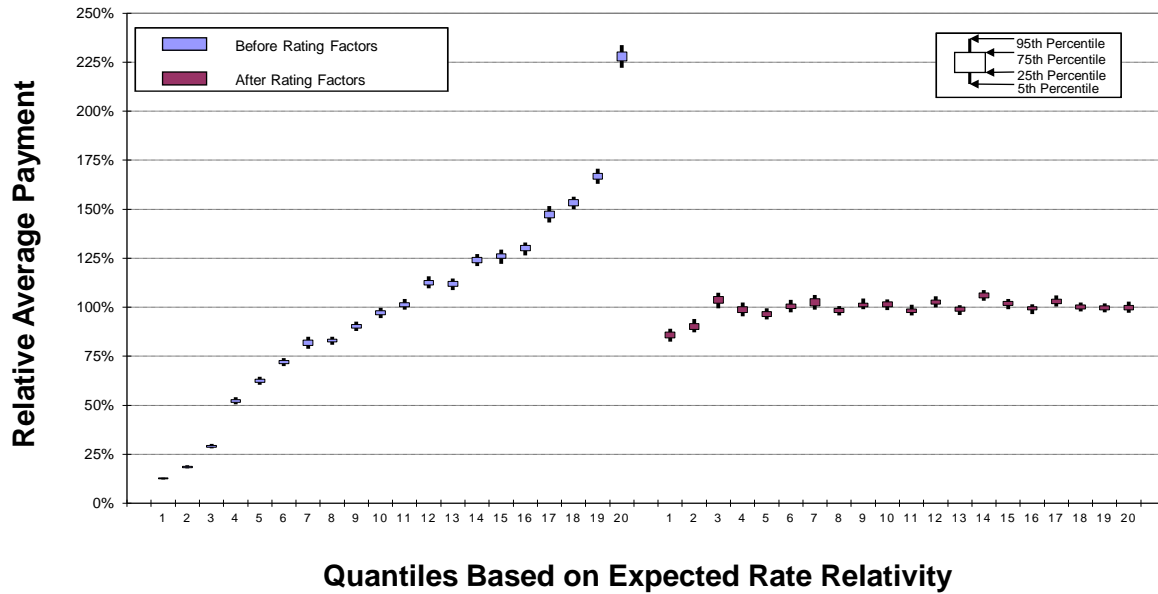


Figure 3. Allegation Nature - Bootstrap Test Validation of Minimum Bias Rating Factors

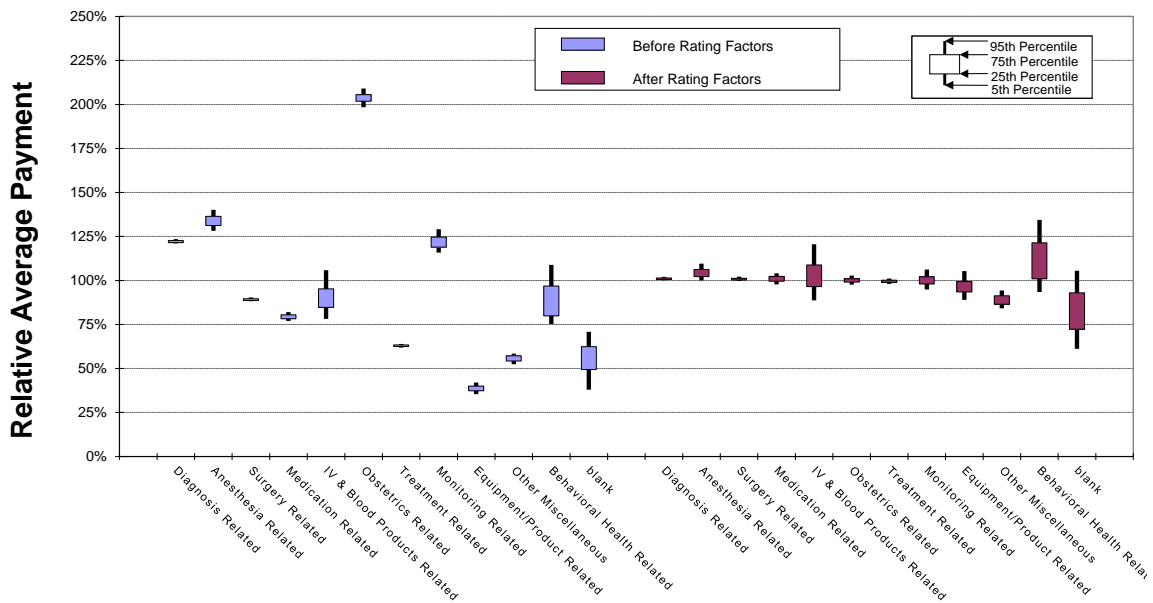


Figure 4. Bootstrap 20 Quantiles Test Validation of Traditional Rating Factors

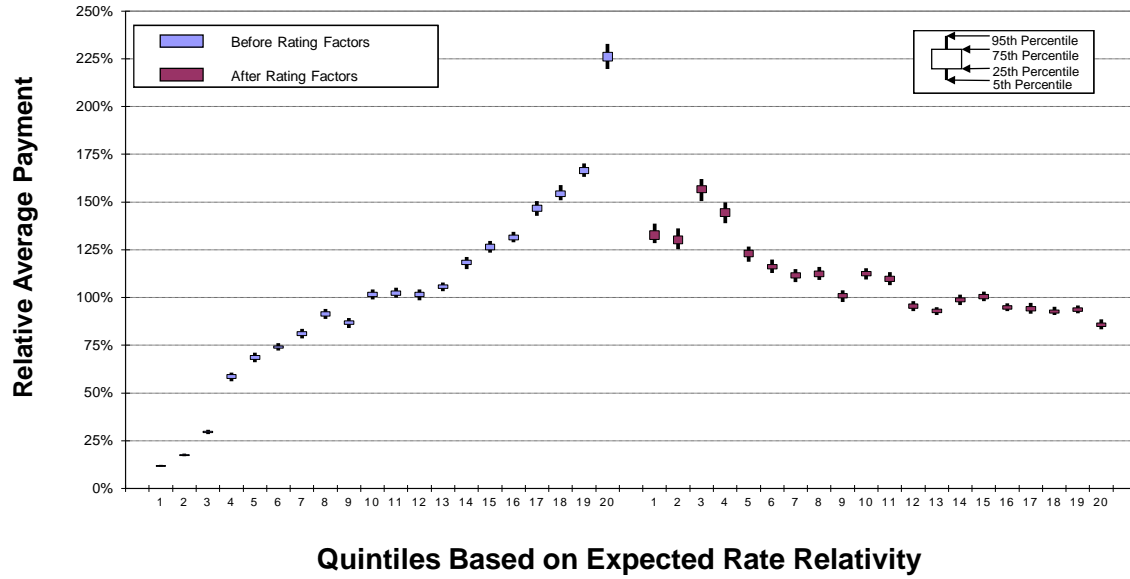


Figure 5. Allegation Nature - Bootstrap Test Validation of Traditional Rating Factors

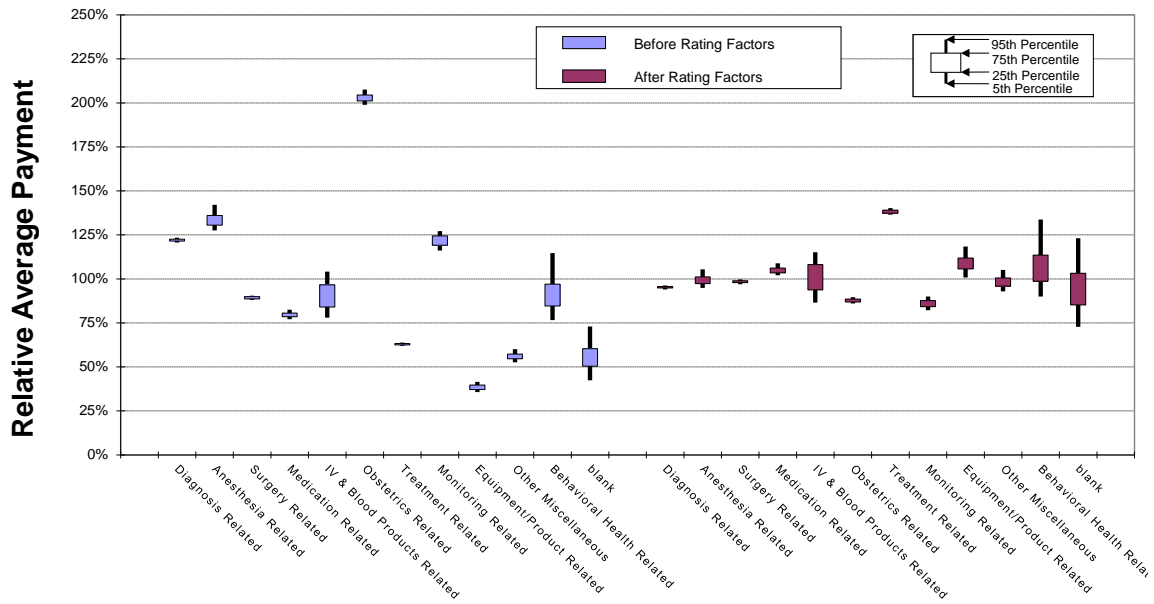


Table 1. Predictive Performance Statistics for Various Models

	20 Quantiles		Allegation Nature	
	Old Statistic	New Statistic	Old Statistic	New Statistic
Mult. Minimum Bias	0.007	0.512	0.023	0.425
GLMs				
Log-Gaussian	0.010	0.511	0.041	0.422
Log-Poisson	0.007	0.512	0.023	0.425
Log-Gamma	0.009	0.511	0.033	0.422
Log-InverseGaussian	Failed to Converge		Failed to Converge	
Traditional	0.135	0.470	0.089	0.408

At this point we have a clear picture of the relative predictive performance of the different models. However, we have not specifically tested the validity of any of the model assumptions, such as likelihoods, independence assumptions, etc. The optimal performance of Minimum Bias/Log-Poisson is likely due to the general validity of its implicit connection to the Central Limit Theorem as discussed earlier.

The GLM assumption that all risks are identically distributed is potentially problematic taken together with the log-link function.

Figures 6 through 8 illustrate the lack of distributional consistency for this dataset. We have broken the observations in the training data into 20 quantiles weighted by modeled values, sorted by actual/modeled result. Using the same breakpoints, determined from the entire training dataset, we then calculated the amount of summed modeled values for each allegation group. If the errors were identically distributed for each allegation group there should be only a random fluctuation around the 5% of total expected for each bin.

Figure 6 shows all allegation natures and naturally each bin demonstrates no differences in the weighted proportion. Figure 7 shows that the anesthesia related allegation group has a much higher percentage in the lowest bin than what would have been expected from the overall population,

Figure 6. All Allegation Nature 20 Value Weighted Quantile Bins

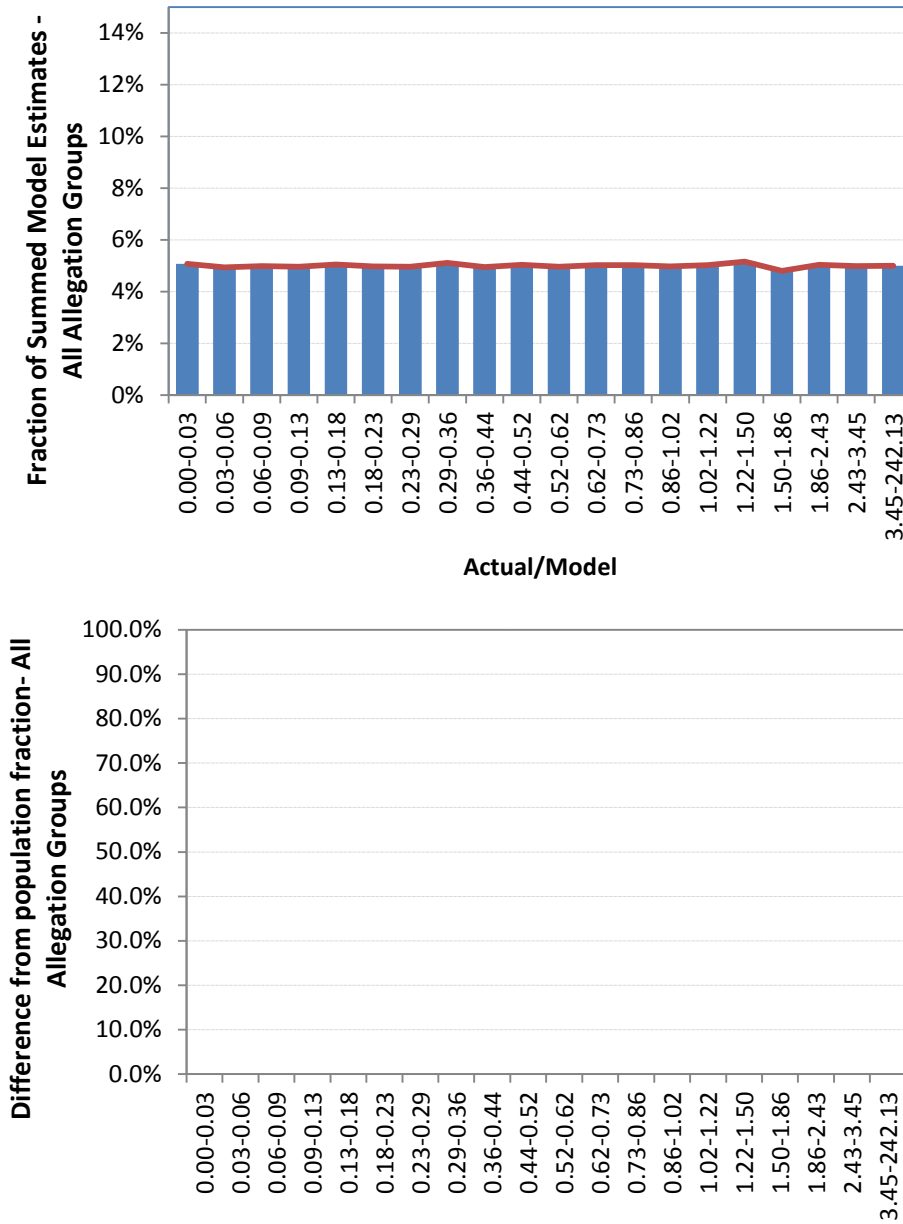


Figure 7. Anesthesia Allegation 20 Value Weighted Quantile Bins

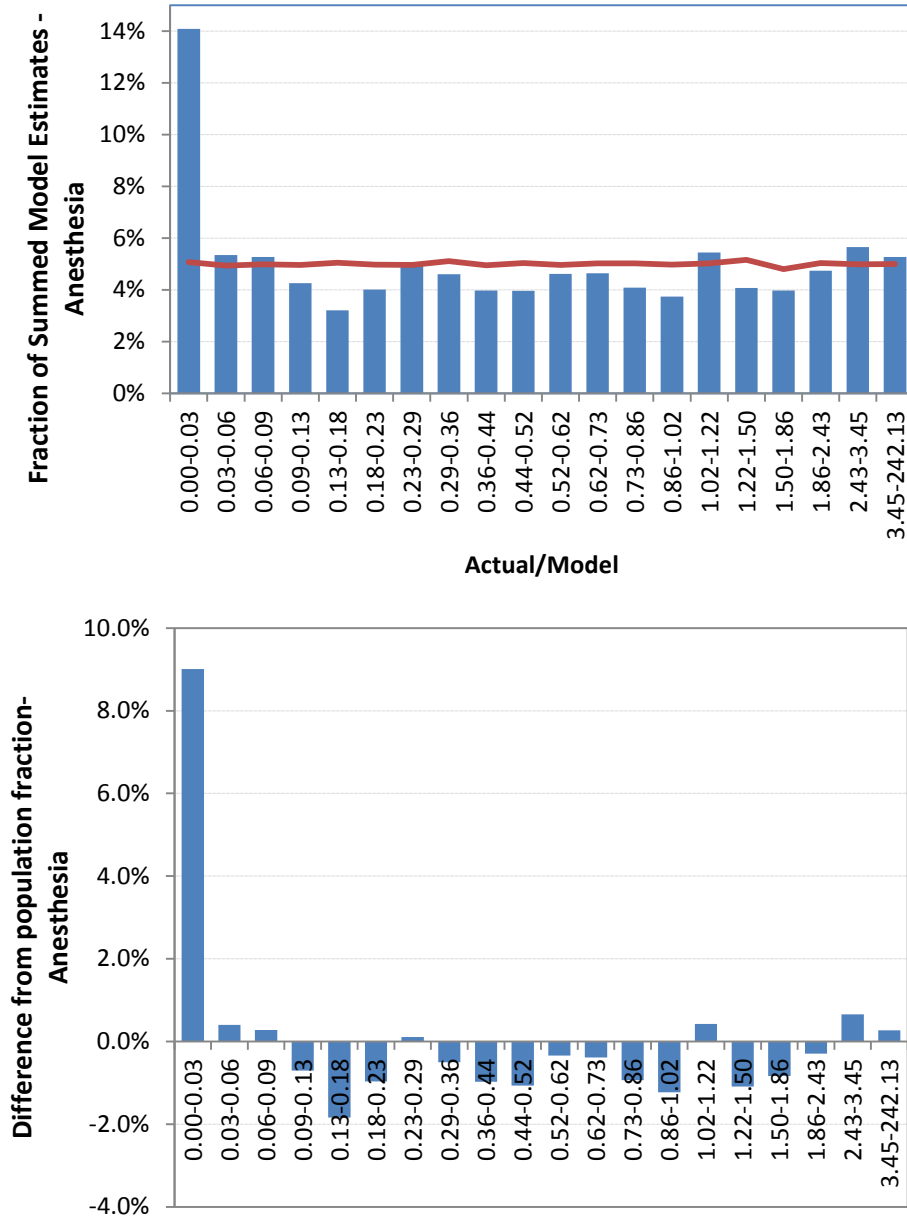
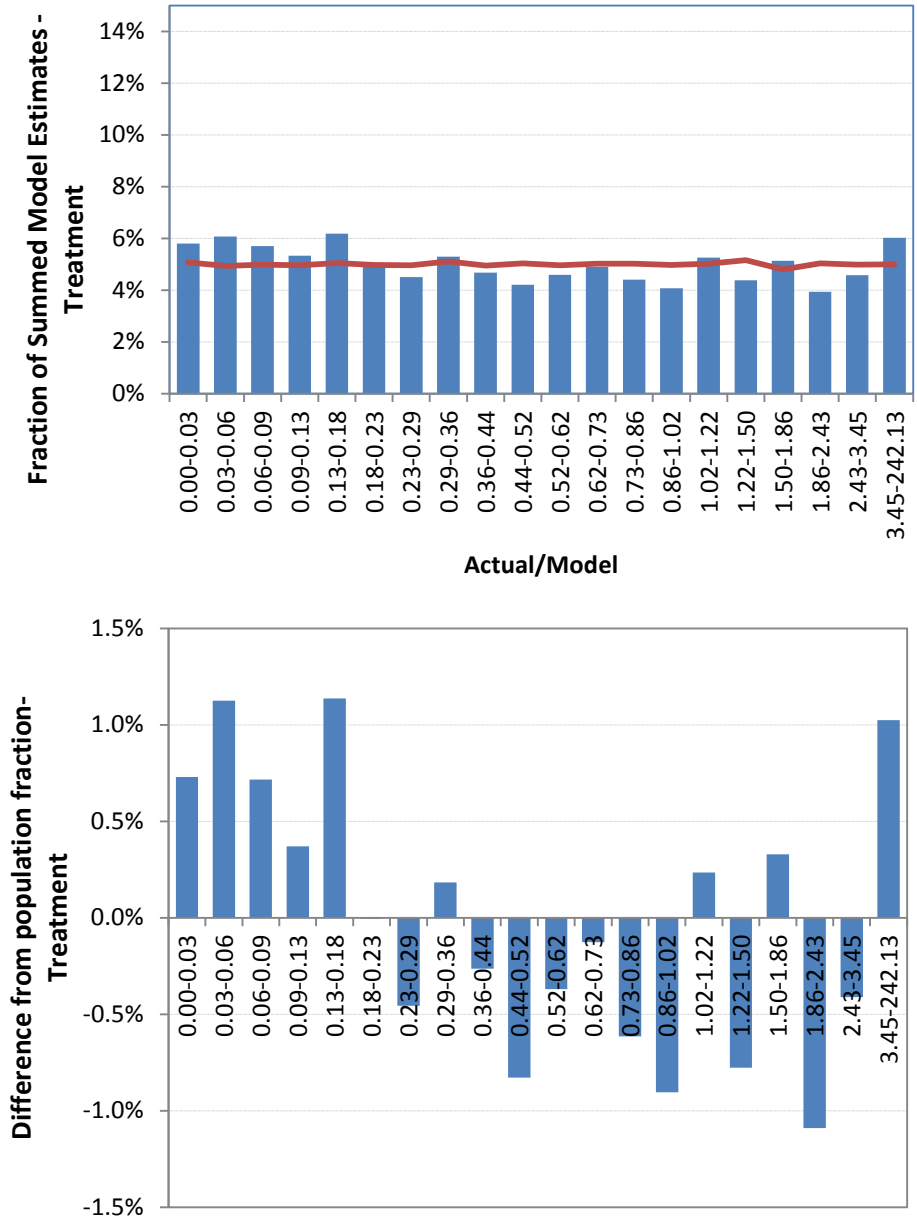


Figure 8. Treatment Allegation 20 Value Weighted Quantile Bins



of the error distribution. Figure 8 shows that, while not as dramatic, the treatment related allegation group shows greater variation than the overall error distribution, with more of the highest and lowest values.

This is far from uncommon with highly-skewed insurance data. The problem is compounded by the multiple dimensions of data. Error distributions could be, and likely are, differently distributed across many of the dimensions, if not every dimension being analyzed. Without adjustment, the basic assumption in a GLM is that the errors are identically distributed. The use of the log-link function, in conjunction with maximum likelihood estimation, puts a great deal of faith in the distributional assumption, inferring conclusions about results in the tail, based on the more voluminous observations at the lower parts of the distribution. But it is the tail itself that is of primary interest in most insurance questions, with the majority of the aggregate losses being caused by the minority of claims. Despite the unreasonable implied assumption of a log-Poisson GLM, because it happens to have effectively the same parameter estimation formulas as the multiplicative minimum bias approach, which has the associated Central Limit Theorem advantages previously described, it is less vulnerable to these distributional differences.

Table 2 shows a comparison of the model biases by allegation group on the validation data using multiplicative minimum bias with full credibility vs. GLM with a log-Gaussian assumption, by comparing actual aggregated results by allegation group to aggregated modeled results over a number of bootstrapped test sets. Despite the log-Gaussian assumption better characterizing the distribution of the data than does the log-Poisson assumption, it ultimately produces estimates that are more vulnerable to distributional differences. The only allegation group with a worse log-Gaussian mean bias was that for Equipment/Product related payments, and in that group, both sets of bootstrapped ranges contained zero, suggesting that the bias measure was inconclusive.

Table 2. Bootstrapped (Actual – Modeled)/Modeled By Allegation Nature

	Multiplicative Minimum Bias			Log-Gaussian		
	Mean	5th %	95th %	Mean	5th %	95th %
Diagnosis	1.0%	0.1%	2.0%	1.3%	0.4%	2.3%
Anesthesia	4.3%	0.0%	9.5%	7.1%	2.5%	11.9%
Surgery	0.8%	-0.3%	2.1%	1.1%	-0.2%	2.5%
Medication	0.9%	-2.2%	4.0%	2.2%	-0.6%	5.4%
IV & Blood Products	3.0%	-11.3%	20.5%	3.6%	-6.8%	15.9%
Obstetrics	0.1%	-2.4%	2.8%	-0.4%	-2.3%	1.8%
Treatment	-0.5%	-2.0%	1.1%	-2.5%	-4.0%	-1.0%
Monitoring	0.2%	-5.1%	6.2%	0.9%	-4.3%	5.7%
Equipment/Product	-3.4%	-11.0%	5.4%	0.0%	-9.3%	8.7%
Other	-11.1%	-15.8%	-5.7%	-14.3%	-19.6%	-8.9%
Behavioral Health	11.9%	-6.5%	34.4%	13.2%	-10.4%	40.9%
Blank	-17.0%	-38.8%	5.5%	-20.7%	-37.7%	-0.6%

7.3 Incorporating Credibility into Minimum Bias

Although the overall predictive performance without any credibility adjustments was very good, there are reasons to explore credibility. In some sparsely populated classes for License Field, rating variables might be so unreliable as to lead to adverse selection problems in real world applications.

In the previous example, the p-values for the rating factors in the Log-Poisson were all infinitesimally low (the largest p-value $\sim 10^{-204}$). This is likely due to the problematic general phenomenon that p-values tend to always implode with very large volumes of data, such as the volume in the example. In stark contrast, most of the p-values for the Log-Gaussian and Log-Gamma models were high, from 1% to approaching 100%. Whether these p-value results indicate any of the likelihood selections are valid, or not, they demonstrate the generally awkward nature of trying to use p-values and class consolidation to handle the lack of credibility in sparsely populated classes.

Rather than attempt a p-value based class consolidation, we will explore the impact of a very simple credibility adjustment for Minimum Bias. We select the very simple form $Z_{j,i_j} = \frac{P_{j,i_j}}{P_{j,i_j} + K}$ where P_{j,i_j} is the number of records where the i_j class for classification dimension j and $K \geq 0$ is a judgmental selection. Table 3 shows that this simple credibility adjustment only tends to erode overall predictive value for this large dataset with only truly predictive variables included.

Table 3. Predictive Performance Statistics for Credibility Adjusted Multiplicative Minimum Bias

	20 Quantiles		Allegation Nature	
	Old Statistic	New Statistic	Old Statistic	New Statistic
Mult. Minimum Bias				
K = 0	0.007	0.512	0.023	0.425
K = 1	0.009	0.511	0.032	0.425
K = 10	0.010	0.511	0.030	0.423
K = 25	0.009	0.510	0.029	0.425
K = 50	0.010	0.511	0.022	0.424
K = 100	0.011	0.511	0.028	0.425
K = 200	0.013	0.509	0.031	0.423
K = 700	0.023	0.505	0.082	0.414

Minimum Bias, GLMs, and Credibility in the Context of Predictive Modeling

To construct a smaller example where credibility is more relevant, we will use a random set of only 5,000 records for fitting and another random set of 5,000 records for testing, shown in Tables 4 and 5 and Figures 9 through 12. We will also do a full test using all the remaining 366,123 records not used for fitting, shown in Tables 6 and 7 and Figures 13 and 14.

Table 4. Smaller Sample Predictive Performance Statistics for Various Models

	6 Quantiles		Allegation Nature	
	Old Statistic	New Statistic	Old Statistic	New Statistic
Mult. Minimum Bias	0.021	0.463	2.216	-0.683
GLMs				
Log-Gaussian	0.041	0.448	3.252	-0.785
Log-Poisson	0.021	0.463	2.216	-0.683
Log-Gamma	0.052	0.445	2.245	-0.704
Log-InverseGaussian	Failed to Converge		Failed to Converge	
Traditional	0.524	0.302	2.419	-0.751

Table 5. Smaller Sample Predictive Performance Statistics for Credibility Adjusted Multiplicative Minimum Bias

	6 Quantiles		Allegation Nature	
	Old Statistic	New Statistic	Old Statistic	New Statistic
Mult. Minimum Bias				
K = 0	0.021	0.463	2.216	-0.683
K = 1	0.016	0.457	1.138	-0.419
K = 10	0.012	0.461	0.454	0.246
K = 25	0.022	0.458	0.394	0.316
K = 50	0.043	0.450	0.376	0.338
K = 100	0.068	0.449	0.373	0.345
K = 200	0.093	0.432	0.384	0.345
K = 700	0.255	0.387	0.479	0.319

As Tables 4 through 7 and Figures 9 through 14 show, the incorporation of credibility was particularly important when distinguishing differences between the allegation groups. Actuaries are regularly asked to provide estimates of the impact of rating variables despite having less than fully

Figure 9. Smaller Sample Bootstrap 6 Quantiles Test Validation of Minimum Bias Rating Factors

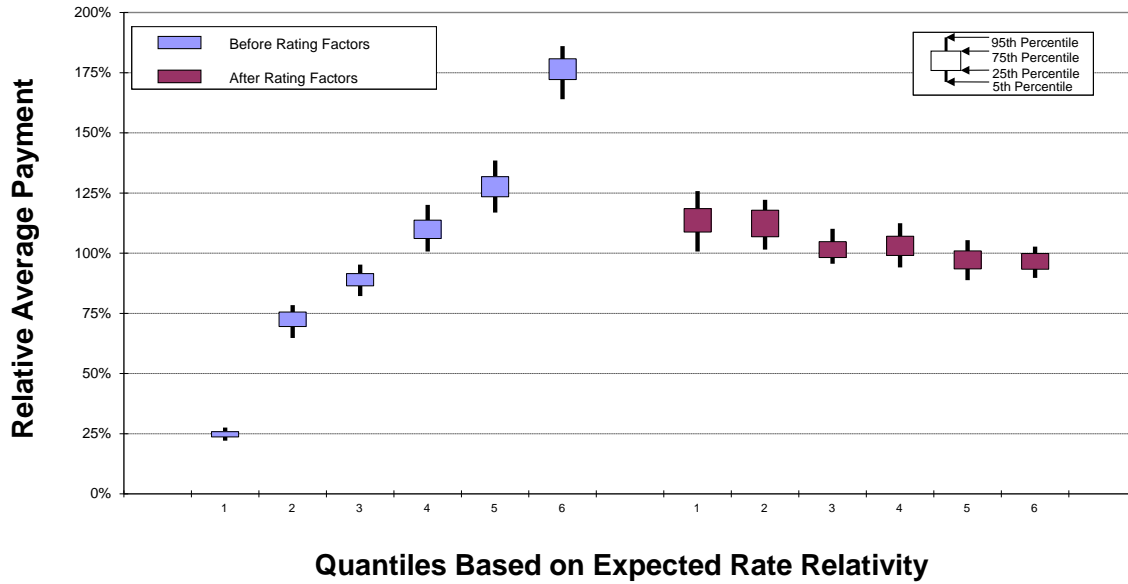


Figure 10. Smaller Sample Allegation Nature - Bootstrap Test Validation of Minimum Bias Rating Factors

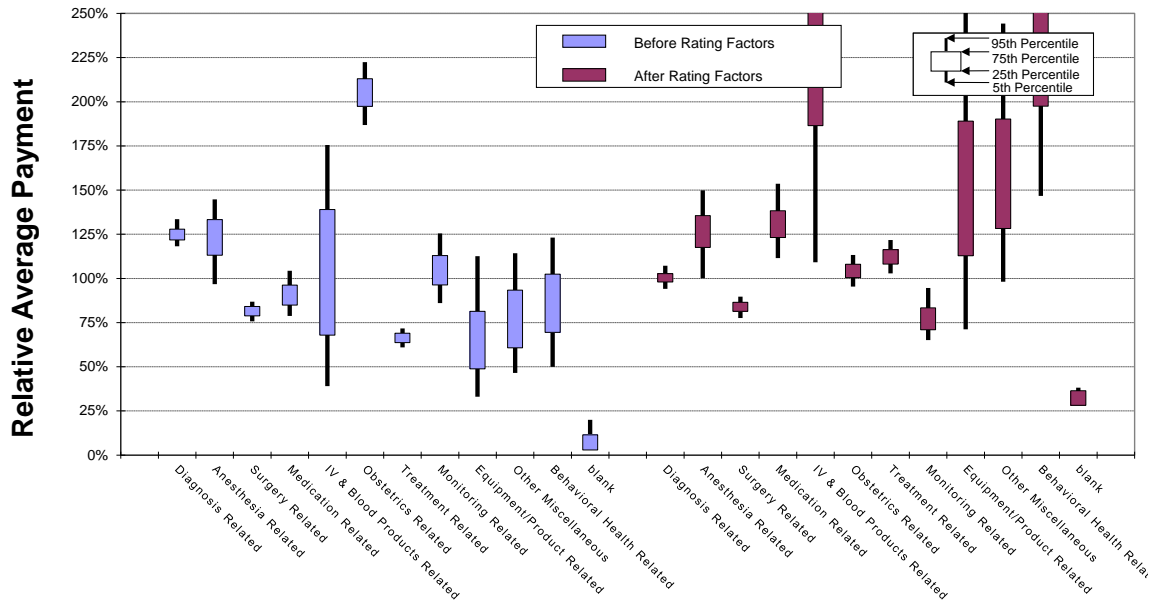


Figure 11. Smaller Sample Bootstrap 6 Quantiles Test Validation of Minimum Bias (Credibility K = 10) Rating Factors

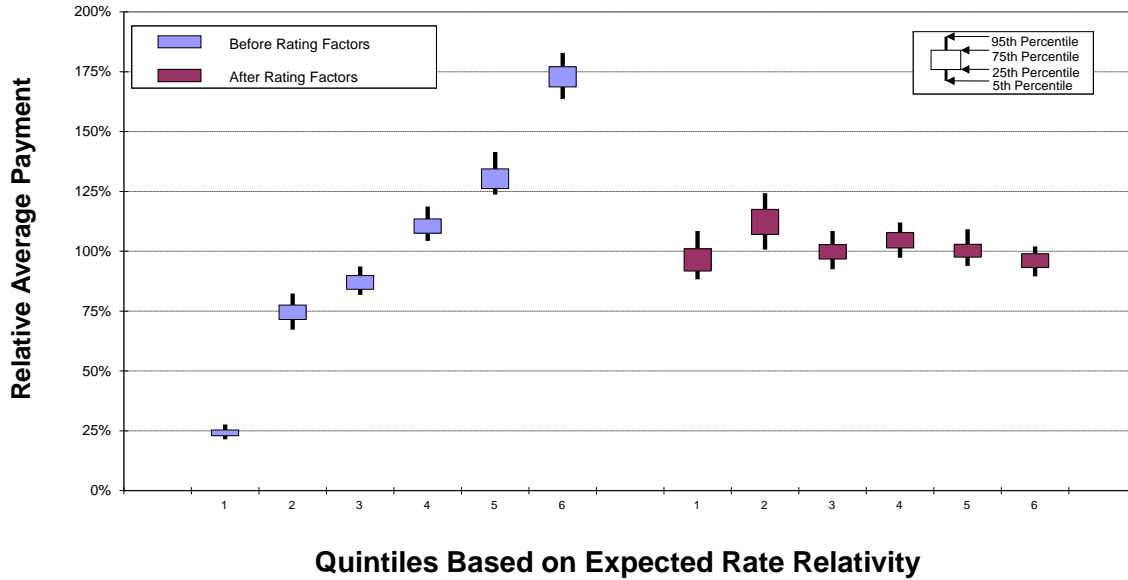


Figure 12. Smaller Sample Allegation Nature - Bootstrap Test Validation of Minimum Bias (Credibility K = 10) Rating Factors

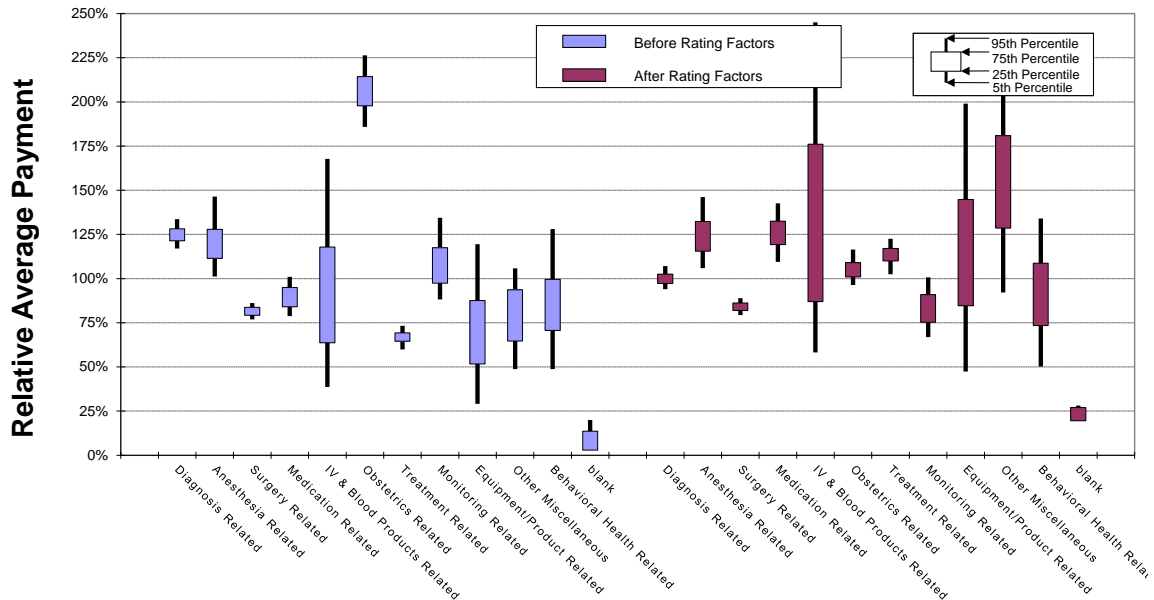


Table 6. Full Test of Smaller Sample Predictive Performance Statistics for Various Models

	20 Quantiles		Allegation Nature	
	Old Statistic	New Statistic	Old Statistic	New Statistic
Mult. Minimum Bias	0.031	0.488	1.906	-0.403
GLMs				
Log-Gaussian	0.038	0.482	2.673	-0.556
Log-Poisson	0.031	0.488	1.906	-0.403
Log-Gamma	0.072	0.474	3.256	-0.653
Log-InverseGaussian	Failed to Converge		Failed to Converge	
Traditional	0.489	0.350	2.158	-0.471

Table 7. Full Test of Smaller Sample Predictive Performance Statistics for Credibility Adjusted Multiplicative Minimum Bias

	20 Quantiles		Allegation Nature	
	Old Statistic	New Statistic	Old Statistic	New Statistic
Mult. Minimum Bias				
K = 0	0.031	0.488	1.906	-0.403
K = 1	0.020	0.492	0.835	0.139
K = 10	0.012	0.494	0.169	0.380
K = 25	0.013	0.493	0.187	0.379
K = 50	0.026	0.489	0.215	0.372
K = 100	0.063	0.479	0.246	0.364
K = 200	0.117	0.460	0.289	0.355
K = 700	0.300	0.399	0.427	0.317

Figure 13. Full Test of Smaller Sample Bootstrap 6 Quantiles Test Validation of Minimum Bias (Credibility K = 10) Rating Factors

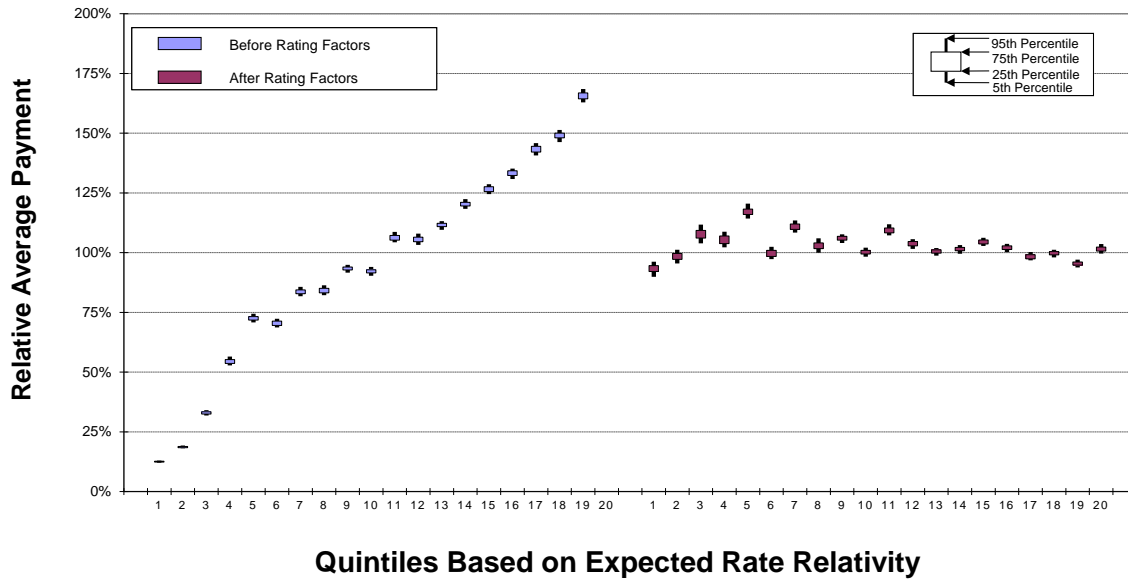
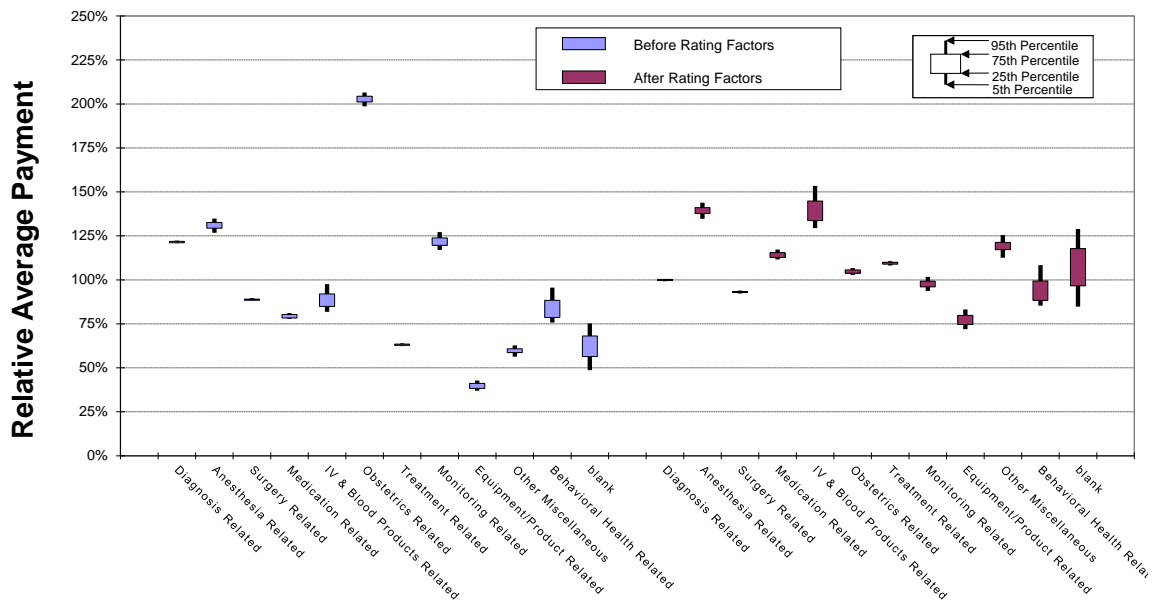


Figure 14. Full Test of Smaller Sample Allegation Nature - Bootstrap Test Validation of Minimum Bias (Credibility K = 10) Rating Factors



credible data. While the overall result may appear to be relatively unaffected by increasing the credibility standard, the ability to more robustly differentiate between them is illustrated.

8.4 Incorporating Credibility Into GLM

We can incorporate credibility, or “shrinkage” of parameter estimates, into a GLM model by defining a hierarchical Bayesian Network of random variables:

$$U_{1,j} = 0 \quad j = 1,2,3$$

$$U_{1,4} = \text{Uniform}(0,20)$$

$$U_{i,1} \sim \text{Normal}(-\sigma_1^2 / 2, \sigma_1^2) \quad i = 2, \dots, 83$$

$$U_{i,2} \sim \text{Normal}(-\sigma_1^2 / 2, \sigma_1^2) \quad i = 2, \dots, 12$$

$$U_{i,3} \sim \text{Normal}(-\sigma_1^2 / 2, \sigma_1^2) \quad i = 2, \dots, 9$$

$$\sigma_1^2 \sim \text{Lognormal}(0,10)$$

$$\sigma_2^2 \sim \text{Lognormal}(0,10)$$

$$\delta_k \sim \text{Normal}(-\sigma_2^2 / 2, \sigma_2^2) \quad k = 1, \dots, n$$

$$Y_k \sim \text{Poisson}(\text{Exp}(\delta_k + U_{1,4} + U_{i_{1,k},1} + U_{i_{2,k},2} + U_{i_{3,k},3})) \quad k = 1, \dots, n$$

Y_k are the individual actual claim amounts to be fit. $U_{i,j}$ are parameters in log space, with $U_{1,4}$ being a constant and the other $j=1,2$, or 3 corresponding to *License Field*, *Allegation Group*, and *Original Year*, respectively. $i_{j,k}$ is an index of which class the Y_k observation falls into in each classification dimension. δ_k is a random over-dispersion for each observation which itself has variance σ_2^2 . σ_1^2 is the parameter variance for each class parameter. Since $U_{1,4}$, σ_1^2 , and σ_2^2 follow highly diffuse Casualty Actuarial Society *E-Forum*, Winter 2017

distributions they will effectively be “fitted” parameters when Gibbs Sampling is performed. σ_1^2 , and σ_2^2 conceptually correspond to parameter and process variances in credibility.

We will also defined a simpler form of this model eliminating the over-dispersion arising from σ_1^2 and σ_2^2 . Running this simpler model numerically produced the same parameters as the MLE Log-Poisson/Minimum Bias with no credibility adjustment, confirming that our Gibbs Sampling model is constructed and coded on the right track up to to the point of adding credibility adjustments.

When the model including the δ_k and σ_2^2 was run numerically there was a shrinkage effect observed in the set of parameters. Table 8 shows that the range of the $U_{i,1}$ contracted significantly with over-dispersion. There was a slight broadening of the ranges for $U_{i,2}$ and $U_{i,3}$, which is not unreasonable as none of the corresponding classes in these dimensions are sparsely populated.

Table 8. Shrinkage Effect in Range of Gibbs Sampled Parameter Fits

	$U_{i,1}$		$U_{i,2}$		$U_{i,3}$	
	Min	Max	Min	Max	Min	Max
Large Split						
w/o overdispersion	-4.103	0.775	-0.920	0.473	0.000	0.691
w overdispersion	-2.173	0.550	-0.975	0.742	-0.040	0.494
Smaller Sample						
w/o overdispersion	-6.570	2.234	-1.405	0.432	0.000	0.742
w overdispersion	-2.033	0.963	-1.992	0.318	-0.069	0.691

Unfortunately, although there was a credibility-like shrinkage affect, the predictive performance actually deteriorated. Figures 15 and 16 show the deteriorating situation when the Gibbs Sampling with over-dispersion is included in the large split of the data. Table 9 shows the deterioration in test statistics for both the large split and smaller sample.

There are potential criticisms of the Bayesian network model as we have defined it. For example, the anchoring of the parameters for the first classes $U_{1,j} = 0 \quad j = 1,2,3$; offsetting the prior distributions on parameters so as to have mean 1 after exponentiation $U_{i,1} \sim \text{Normal}(-\sigma_1^2 / 2, \sigma_1^2) \quad i = 2, \dots, 83$; the same parameter variance σ_1^2 was used for all three classification dimensions; etc. However, the authors experimented with a myriad of alterations to the model definition, even going so far as to convert the likelihood function into a Negative Binomial distribution to capture the impact of over-dispersion of the Poisson more directly. In all cases tried performance deteriorated further or did not improve. The earlier presented multiplicative minimum

bias model with incorporated credibility would be vulnerable to similar or more extensive potential criticisms. Yet implementing it went quickly and easily produced desirable results.

This failed modeling experience in no way proves that a well performing Gibbs Sampled Bayesian model cannot be defined in this context. Obviously, well performing examples for much simpler situations, such as one classification dimension and an identity link function, are well known and easy to construct. Nor is the point that the theory behind these models does not provide deep insights into understanding modeling and statistical estimation. However, in this case, orders of magnitude more input of resources both in time and sophistication in effort than was used for minimum bias produced inferior predictive performance. Though neither author of this paper is a specialist in Gibbs Sampling methods, one author (Evans) has used them occasionally for over 10 years and informally consulted several more experienced specialists (in Acknowledgements). As of this writing, we have not been able to diagnose why the model as defined performs so much more poorly than a regular MLE GLM with no shrinkage effect. Whether the model is in some way poorly designed or, much less likely, one of the many technical choices made in running the Gibbs Sampling software should be tuned differently, does not alter the key conclusion. Namely, that the tremendous additional resource and intellectual burdens of such detailed and sophisticated models may offer no advantage, or may even be disadvantageous, in many practical situations of predictive modeling.

Figure 15. Full Test of Smaller Sample Bootstrap 20 Quantiles Test Validation of Gibbs Sampled Rating Factors with Shrinkage

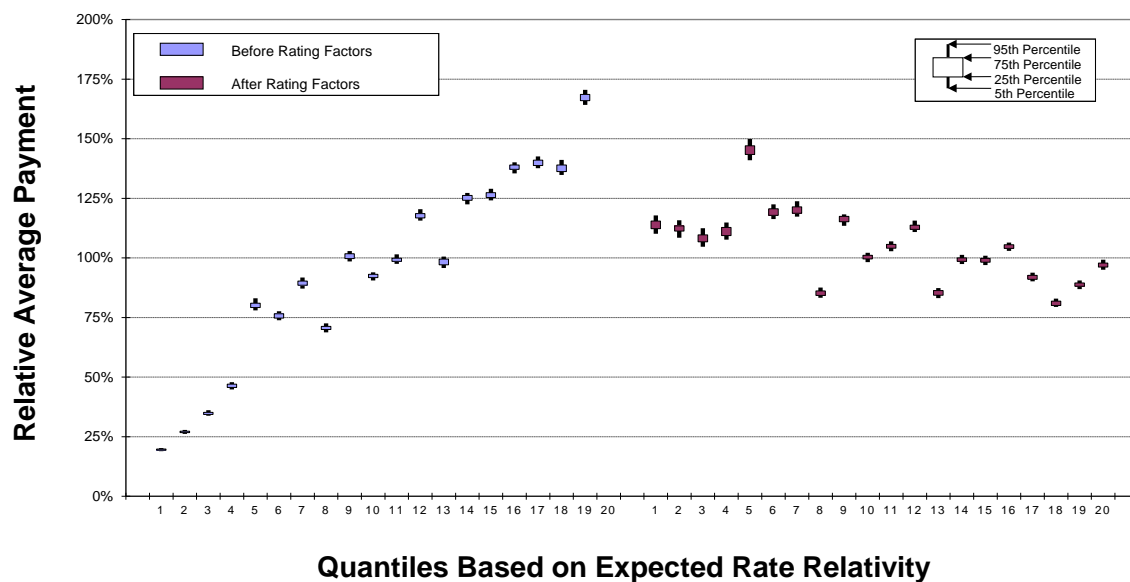


Figure 16. Full Test of Smaller Sample - Allegation Nature - Bootstrap Test Validation of Gibbs Sampled Rating Factors with Shrinkage

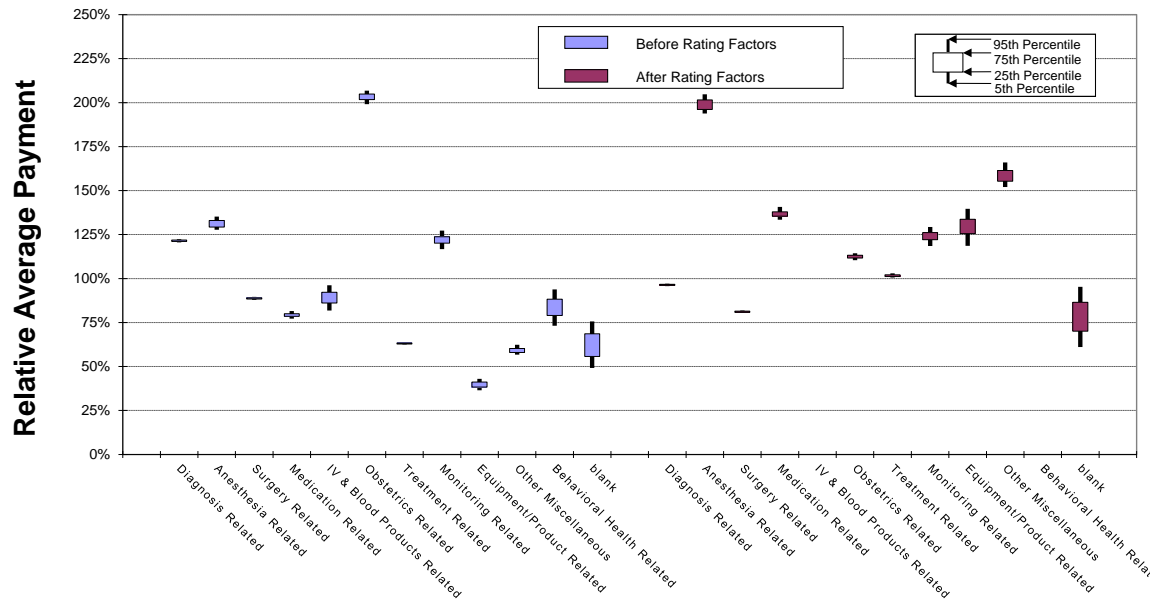


Table 9. Test Statistics for Sampled Rating Factors

Gibbs

	Quantiles		Allegation Nature	
	Old Statistic	New Statistic	Old Statistic	New Statistic
Large Split (20 Quantiles)				
w/o overdispersion	0.007	0.512	0.023	0.425
w overdispersion	0.102	0.463	0.219	0.376
Smaller Sample (6 Quantiles)				
w/o overdispersion	0.021	0.463	2.216	-0.683
w overdispersion	0.101	0.403	3.616	-0.943
Full Test Smaller Sample (20 Quantiles)				
w/o overdispersion	0.031	0.488	1.906	-0.403
w overdispersion	0.098	0.448	4.723	-0.818

8. SUMMARY DISCUSSION

The predictive modeling framework greatly reduces the burdens of model specification, because models are validated based on their predictive performance rather than hypothesis testing of model assumptions. Minimum bias models transform basic data in such a way as to partially sacrifice sample independence in exchange for much tamer distributions of individual data points that are much less

need of detailed distributional specification. The combination of multiplicative minimum bias iteration with a generic incorporation of credibility as presented in this paper demonstrates that a very simple model, without complete distributional specification, in practice can provide comparable or better predictive value than a far more complex model, such as a typical GLM.

GLM models are fit to individual data points and require specification of the distributions underlying each data point. Consequently, GLM models can be significantly vulnerable to inaccurate specifications and their fundamental complexity makes the practical incorporation of credibility adjustments, such as including random effects or fitting parameters through Gibbs sampling, very complex.

Philosophically, simpler modeling is desirable. In practice, simpler models are beneficial in many ways, such as lower skill requirements for operational personnel and greater lucidity to a much wider audience. Some previous papers, such as Brown in [5] and Mildenhall in [10], have highlighted the sense in which minimum bias iteration is a special case of GLM and encouraged – at least implicitly – minimum bias practitioners to switch to GLM as a richer framework. There is some irony that with the advent of the predictive framework minimum bias may often be somewhat more advantageous, in principle and practice.

While GLM models are powerful and belong in the set of tools applied by actuaries, consideration should also be given to multiplicative minimum bias models and the traditional actuarial concept of partial credibility. Ultimately the test of any predictive model should be how it performs on out-of-sample data.

Acknowledgment

The authors acknowledge are thankful to Jose Couret, Louise Francis, Chris Laws and Frank Schmid for answering some questions that arose in the course of writing this paper.

Appendix A Details of Empirical Case Study

The empirical data used in this case study consists of 371,123 records of medical malpractice payments obtained from the National Practitioner Data Bank. Three explanatory variables were used for modeling payment amounts: *Original Year*, *Allegation Group* and *License Field*. The following Tables A.1 through A.3 display record counts by each of the explanatory variables overall and for the individual predictive modeling splits.

Table A.1 Counts of Records by License Field

License Field	Total	Large Split		Smaller Sample		
		Fit	Test	5,000 Fit	5,000 Test	Full Test
Allopathic Physician (MD)	271,443	135,514	135,929	3644	3661	267,799
Phys. Intern/Resident (MD)	2,113	1,063	1,050	34	28	2,079
Osteopathic Physician (DO)	17,612	8,829	8,783	237	244	17,375
Osteo. Phys. Intern/Resident (DO)	324	161	163	8	6	316
Dentist	46,516	23,425	23,091	623	596	45,893
Dental Resident	145	64	81	4	3	141
Pharmacist	1,890	952	938	24	20	1,866
Pharmacy Intern [available 9/9/2002]	2	1	1	0	0	2
Pharmacist, Nuclear	6	4	2	0	0	6
Pharmacy Assistant	19	12	7	0	0	19
Pharmacy Technician [available 9/9/2002]	12	7	5	0	1	12
Registered (RN) Nurse	5,715	2,885	2,830	91	80	5,624
Nurse Anesthetist	1,568	777	791	19	19	1,549
Nurse Midwife	873	431	442	18	8	855
Nurse Practitioner	1,288	598	690	19	24	1,269
Doctor of Nursing Practice [available 11/8/2010]	1	-	1	0	0	1
Advanced Nurse Practitioner [3/5/02 - 9/9/02]	4	3	1	0	0	4
LPN or Vocational Nurse	692	345	347	9	9	683
Clinical Nurse Specialist [available 9/9/02]	18	12	6	1	0	17
Certified Nurse Aide/Nursing Assistant [available 10/17/05]	36	18	18	0	1	36
Nurses Aide	78	39	39	2	2	76
Home Health Aide (Homemaker)	22	10	12	0	0	22
Health Care Aide/Direct Care Worker [available 10/17/05]	3	1	2	0	0	3
Psychiatric Technician	15	10	5	0	0	15
Dietician	22	11	11	0	1	22
Nutritionist	1	1	-	0	0	1
EMT, Basic	200	106	94	3	2	197
EMT, Cardiac/Critical Care	28	17	11	0	0	28
EMT, Intermediate	26	13	13	1	2	25
EMT, Paramedic	59	32	27	0	1	59
Clinical Social Worker	206	107	99	2	0	204
Podiatrist	7,654	3,809	3,845	92	113	7,562
Clinical Psychologist [last use 9/9/02]	875	436	439	15	15	860
Psychologist [available 9/9/02]	352	174	178	2	5	350
School Psychologist [available 9/9/02]	1	-	1	0	0	1
Audiologist	39	23	16	2	1	37
Art/Recreation Therapist	2	1	1	0	0	2
Massage Therapist	82	54	28	3	1	79
Occupational Therapist	85	43	42	0	0	85
Occup. Therapy Assistant	11	7	4	0	0	11
Physical Therapist	1,094	545	549	14	14	1,080

Minimum Bias, GLMs, and Credibility in the Context of Predictive Modeling

Table A.1 Counts of Records by License Field (continued)

License Field	Total	Large Split		Smaller Sample		
		Fit	Test	5,000 Fit	5,000 Test	Full Test
Phys. Therapy Assistant	94	48	46	0	3	94
Rehabilitation Therapist	9	3	6	0	0	9
Speech/Language Pathologist	14	9	5	0	0	14
Hearing Aid/Instrument Specialist [available 10/17/05]	2	1	1	0	0	2
Medical Technologist [changed to 501(6/15/09)]	64	28	36	0	0	64
Medical/Clinical Lab Technologist [available 6/15/09]	1	1	-	0	0	1
Medical/Clinical Lab Technician [available 6/15/09]	2	-	2	0	0	2
Surgical Technologist [available 6/15/09]	7	4	3	0	0	7
Surgical Assistant [available 6/15/09]	1	-	1	0	0	1
Cytotechnologist [available 11/22/99]	11	7	4	0	0	11
Nuclear Med. Technologist	14	5	9	0	0	14
Rad. Therapy Technologist	12	5	7	0	0	12
Radiologic Technologist	169	89	80	1	0	168
X-Ray Technician or Operator [available 6/15/09]	5	2	3	0	0	5
Acupuncturist	58	22	36	0	0	58
Athletic Trainer [available 11/22/99]	6	3	3	1	0	5
Chiropractor	5,834	2,928	2,906	78	87	5,756
Dental Assistant	15	8	7	1	1	14
Dental Hygienist	41	22	19	1	2	40
Denturist	27	8	19	0	0	27
Homeopath	6	5	1	1	0	5
Medical Assistant	33	14	19	1	0	32
Counselor, Mental Health	167	84	83	1	2	166
Midwife, Lay (Non-Nurse)	22	14	8	0	0	22
Naturopath	17	9	8	0	0	17
Ocularist	25	12	13	0	1	25
Optician	17	10	7	0	0	17
Optometrist	715	367	348	6	11	709
Orthotics/Prosthetics Fitter	9	5	4	1	0	8
Phys. Asst., Allopathic	1,713	847	866	26	22	1,687
Phys. Asst., Osteopathic	137	71	66	3	3	134
Perfusionist [available 11/22/99]	8	2	6	1	0	7
Podiatric Assistant	14	9	5	0	0	14
Prof. Counselor	209	109	100	4	3	205
Prof. Cnslr., Alcohol	9	2	7	0	1	9
Prof. Cnslr., Family/Marriage	177	96	81	4	5	173
Prof. Cnslr, Substance Abuse	23	13	10	0	0	23
Marriage and Family Therapist [available 9/9/02]	27	15	12	1	0	26
Respiratory Therapist	48	24	24	1	0	47
Resp. Therapy Technician	14	4	10	0	0	14
Other Health Care Pract, Not Classified [available 11/22/99]	45	31	14	0	0	45
Unspecified or Unknown	170	86	84	1	2	169
Total	371,123	185,562	185,561	5,000	5,000	366,123

Table A.2 Counts of Records by Allegation Nature

Allegation Nature	Total	Large Split		Smaller Sample		
		Fit	Test	5,000 Fit	5,000 Test	Full Test
Diagnosis Related	105,674	52,516	53,158	1,409	1,388	104,265
Anesthesia Related	10,974	5,421	5,553	127	153	10,847
Surgery Related	88,763	44,538	44,225	1,176	1,211	87,587
Medication Related	20,197	10,047	10,150	259	268	19,938
IV & Blood Products Related	1,259	625	634	14	16	1,245
Obstetrics Related	25,988	13,081	12,907	384	345	25,604
Treatment Related	100,666	50,517	50,149	1,380	1,372	99,286
Monitoring Related	7,313	3,594	3,719	103	106	7,210
Equipment/Product Related	2,037	989	1,048	32	24	2,005
Other Miscellaneous	7,404	3,791	3,613	106	106	7,298
Behavioral Health Related	677	361	316	7	9	670
blank	171	82	89	3	2	168
Total	371,123	185,562	185,561	5,000	5,000	366,123

Table A.3 Counts of Records by Origination Year Group

Origination Year	Total	Large Split		Smaller Sample		
		Fit	Test	5,000 Fit	5,000 Test	Full Test
1990-1992	40,574	20,306	20,268	568	515	40,006
1993-1994	39,016	19,480	19,536	570	529	38,446
1995-1996	37,048	18,557	18,491	516	509	36,532
1997-1998	35,689	17,838	17,851	490	493	35,199
1999-2000	38,036	19,045	18,991	469	516	37,567
2001-2002	39,277	19,650	19,627	491	533	38,786
2003-2004	36,565	18,256	18,309	472	508	36,093
2005-2007	47,519	23,756	23,763	659	646	46,860
2008-2012	57,399	28,674	28,725	765	751	56,634
Total	371,123	185,562	185,561	5,000	5,000	366,123

Appendix B Gibbs Sampling Model Code

With Poisson Over-dispersion

```
model
{
  U[1,4]~dunif(0,20)
  U[1,1]<-0
  U[1,2]<-0
  U[1,3]<-0
  Tau[1] ~ dlnorm(0,0.1)
  Mu<- -pow(Tau[1],-1)/2
  Tau[2] ~ dlnorm(0,0.1)
  Mu2<- -pow(Tau[2],-1)/2
  Tau[3]<-Tau[1]/Tau[2]
  for(i in 2:N1) { U[i,1]~dnorm(Mu,Tau[1]) }
  for(i in 2:N2) { U[i,2]~dnorm(Mu,Tau[1]) }
  for(i in 2:N3) { U[i,3]~dnorm(Mu,Tau[1]) }
  for(i in 1:N) {
    ProcError[i]~dnorm(Mu2,Tau[2])
    lambda1[i]<-exp(min(20,ProcError[i]+U[1,4]+U[X[i,1],1]+U[X[i,2],2]+U[X[i,3],3]))
    Y[i]~dpois(lambda1[i])
  }
}
```

Without Poisson Over-dispersion

```
model
{
  U[1,4]~dunif(0,20)
  U[1,1]<-0
  U[1,2]<-0
  U[1,3]<-0
  Tau[1] ~ dlnorm(0,0.1)
  Mu<- -pow(Tau[1],-1)/2

  for(i in 2:N1) { U[i,1]~dnorm(Mu,Tau[1]) }
  for(i in 2:N2) { U[i,2]~dnorm(Mu,Tau[1]) }
  for(i in 2:N3) { U[i,3]~dnorm(Mu,Tau[1]) }
  for(i in 1:N) {
    lambda1[i]<-exp(min(20,U[1,4]+U[X[i,1],1]+U[X[i,2],2]+U[X[i,3],3]))
    Y[i]~dpois(lambda1[i])
  }
}
```


Minimum Bias, GLMs, and Credibility in the Context of Predictive Modeling

REFERENCES

- [1] Anderson, D., et al., “A Practitioner's Guide to Generalized Linear Models,” *CAS Exam Study Note* **2007**, 3rd ed.
- [2] Bailey, Robert A., “Insurance Rates with Minimum Bias”, *PCAS* **1963**, Vol. L, 4–13.
- [3] Bailey, Robert A. and LeRoy J. Simon, “Two Studies in Automobile Insurance Ratemaking”, *PCAS* **1960**, Vol. XLVII, 1–19.
- [4] Brosius, E. and Feldblum, S., “The Minimum Bias Procedure: A Practitioner's Guide,” *PCAS* **2003**, Vol. XC, 196-273.
- [5] Brown, Robert L., “Minimum Bias with Generalized Linear Models”, *PCAS* **1988**, Vol. LXXV, 187–217.
- [6] Evans, J. and Dean, C., “The Optimal Number of Quantiles for Predictive Performance Testing of the NCCI Experience Rating Plan,” *Variance* **2014**, Vol. 8 No. 2, 89-104.
- [7] Frees, E., et al., *Predictive Modeling Applications in Actuarial Science: Volume 1, Predictive Modeling Techniques*, Cambridge University Press, **2014**.
- [8] Klinker, F., “Generalized Linear Mixed Models for Ratemaking: A Means of Introducing Credibility into a Generalized Linear Model Setting,” *CAS Forum* **2001**, Winter Vol. 2, 1-25.
- [9] Klugman, S., et al., *Loss Models: From Data to Decisions 4th Edition*, Wiley, **2012**.
- [10] Mildenhall, S., “Minimum Bias and Generalized Linear Models,” *PCAS* **1999**, Vol. LXXVI, 393-487.
- [11] Nelder, J. and Verrall, R., “Credibility Theory and Generalized Linear Models,” *ASTIN Bulletin* **1997**, Vol. 27:1, 71-82.
- [12] Scollnik, D., “An Introduction to Markov Chain Monte Carlo Methods and their Actuarial Applications,” *PCAS* **1996**, Vol. LXXXIII, 114–165.
- [13] Venter, Gary G., “Discussion of Minimum Bias with Generalized Linear Models,” *PCAS* **1990**, Vol. LXXVII, 337–349.

Biographies of the Authors

Christopher Gross, ACAS, MAAA is a property and casualty consulting actuary with 20 years of experience. He is the founder and president of Gross Consulting, which provides actuarial services and the Cognalysis™ line of actuarial software: Claim Life Cycle Model (CLCM™), MultiRate™, and Cognalysis Reserving System™.

Jonathan Evans, FCAS, FSA, FCA, CERA, MAAA, WCP is a property and casualty consulting actuary with 20 years of experience. He is currently President of Convergent Actuarial Services, Inc.