

1999

Credibility Calculations Using Analysis of Variance Computer Routines

Dennis H. Tolley

Brigham Young University, toolley@byu.edu


Michael D. Nielsen

University of Pennsylvania's Wharton School of Business, mdn2@wharton.upenn.edu

Robert Bachler

Educators Mutual Insurance Association, bachlero@educatorsmutual.com

Follow this and additional works at: <http://digitalcommons.unl.edu/joap>

 Part of the [Accounting Commons](#), [Business Administration, Management, and Operations Commons](#), [Corporate Finance Commons](#), [Finance and Financial Management Commons](#), [Insurance Commons](#), and the [Management Sciences and Quantitative Methods Commons](#)

Tolley, Dennis H.; Nielsen, Michael D.; and Bachler, Robert, "Credibility Calculations Using Analysis of Variance Computer Routines" (1999). *Journal of Actuarial Practice 1993-2006*. 87.
<http://digitalcommons.unl.edu/joap/87>

This Article is brought to you for free and open access by the Finance Department at DigitalCommons@University of Nebraska - Lincoln. It has been accepted for inclusion in Journal of Actuarial Practice 1993-2006 by an authorized administrator of DigitalCommons@University of Nebraska - Lincoln.

Credibility Calculations Using Analysis of Variance Computer Routines

H. Dennis Tolley,* Michael D. Nielsen,[†] and Robert Bachler[‡]

Abstract

In this paper we present a method of calculating Bühlmann-Straub credibility factors using standard statistical techniques developed for the analysis of variance. Emphasis is placed on using readily available statistical packages such as SAS and SPSS. Additionally many other computational tools such as EXCEL can be programmed to make such calculations. An example and some sample SAS programs are provided.

Key words and phrases: Bühlmann-Straub credibility factors, empirical Bayes, borrowing strength, random ANOVA model

*Dennis Tolley, Ph.D., A.S.A., is professor of statistics at Brigham Young University. He received his Ph.D. in biostatistics from the University of North Carolina in 1981. He has taught statistics and actuarial methods at Duke University and at Texas A&M University. His research interests are in health and mortality statistics, especially as these regard health care costs issues. He is currently active in health care cost research and in models of health care needs.

Dr. Tolley's address is: Department of Statistics, Brigham Young University, Provo UT 84602, USA. Internet address: tolley@byu.edu

[†]Michael D. Nielsen, A.C.A.S., is currently a doctoral student studying insurance and risk management at the University of Pennsylvania's Wharton School of Business. Before returning to school, he worked as an actuary for both Fireman's Fund Insurance Company and the Allstate Research and Planning Center. He received his M.S. in statistics from Brigham Young University.

Mr. Nielsen's address is: University of Pennsylvania, Wharton School of Business, 3641 Locust Walk, Philadelphia PA 19104, USA. Internet address: mdn2@wharton.upenn.edu

[‡]Robert Bachler, A.S.A., A.C.A.S., M.A.A.A., is a vice president at Educators Mutual Insurance Association in Salt Lake City, Utah. His main practice areas are group health, group disability, and group and individual life. He graduated with a B.S. in Statistics from Brigham Young University.

Mr. Bachler's address is: Educators Mutual Insurance Association, 852 East Arrowhead Lane, Murray UT 84107, USA. Internet address: bachlero@educatorsmutual.com

1 Introduction

Casualty actuaries long have recognized the use of the methods of credibility theory as important in assisting them when setting premiums for (i) renewing business, (ii) blocks of new business, and (iii) determining experience-based refunds. The value of these methods also is gaining recognition among health actuaries.¹ Implementation of these credibility methods, however, is varied. Although formal methods of calculating credibility rates are well established, their implementation varies mathematically from ad hoc computations to simple approximations to detailed estimation of the model parameters. One of the reasons for this is the differences in computational complexity. Despite the fact that company experience is maintained in well-documented databases, use of computer programs on these databases to form credibility estimates is far from seamless and may be too complex to warrant the effort.

We present a method of calculating credibility factors under the Bühlmann-Straub (1970) model using readily available statistical software.² The Bühlmann-Straub model is one of a variety of credibility models and is based on a least squares argument. Though the least squares basis for credibility is adequate justification for the procedure, it has been shown that the Bühlmann-Straub method of calculating credibility is identical to the empirical Bayes method when the distribution of losses is a member of the linear exponential family, the loss is quadratic, and when the Bayesian prior used is the conjugate prior for this distribution (Ericson, 1970). Although software programs do not explicitly identify the credibility factors in the software documentation and are not part of the traditional statistical reports generated by these packages, Bühlmann-Straub credibility factors can be calculated from such packages with minimal effort. This paper illustrates these procedures.

A credibility premium uses data from two sources: the estimate of the pure premium based only on the data from a specific group of interest at a specific time and an estimate of the pure premium based on the other data sources and/or prior information. This second estimate may be the overall average of observed rates taken from samples of other groups of policies or the historical average of the group of policies of interest.

¹There is an extensive literature on credibility in general (see, e.g., Longley-Cook, 1962; Norberg, 1979; Hossack et al., 1983; Herzog, 1996; Goulet, 1998).

²For other papers on the Bühlmann-Straub model see, for example, Morris and Slyke, (1978), and Venter (1985, 1990), and Klugman (1987).

The credibility premium classically takes the form

$$C = ZR + (1 - Z)H, \quad 0 \leq Z \leq 1, \quad (1)$$

where C is the credibility premium; R is the estimate of pure premium using the data from the group of interest; H is a global premium (i.e., an exogenous estimate or assumed value of the average of observations); and Z is the credibility factor and denotes the weight assigned to R . If $Z = 1$ then the data are said to be fully credible, and no compromise estimate is needed.

Although the simple form given in equation (1) is found in most of the literature, there are many different approaches to calculate the credibility factor.³ Bühlmann (1967) arrives at a credibility premium by finding the linear estimator that minimizes the expected squared error. The resulting credibility premium follows the form of the model shown in equation (1), with the credibility factor, Z , given as

$$Z = \frac{n \times \text{VHM}}{n \times \text{VHM} + \text{EPV}} \quad (2)$$

where EPV is the expected value of the process variance and refers to the value of the variance of the pure premium within each group, averaged across all groups; and VHM is the variance of the hypothetical means, which is the mean square distance between the mean of the pure premium in each group and the mean over all groups. Bühlmann (1967) proposes this estimate of credibility for cases when the n_i are equal. The extension to the case where the n_i are not equal is presented by Bühlmann-Straub (1970).

2 The Analysis of Variance (ANOVA) Approach

The connection between credibility methods and analysis of variance (ANOVA)⁴ has been alluded to in several papers. For example, both Venter (1990) and Morris and Van Slyke (1978) describe a model similar

³Morris and Van Slyke (1978) determine Z using a Bayesian framework to obtain a form of equation (1). Bühlmann (1970) suggests an alternative method that is also related to the empirical Bayes approach. Herzog (1996), Philbrick (1981), and Venter (1990) also describe this method.

⁴Analysis of variance is a standard statistical technique in the design and analysis of experiments. For more on analysis of variance, see, for example, Scheffé (1959) and Neter, Wasserman, and Craig (1990, Part 3.)

to the random one-way analysis of variance model. Dannenburg (1995) uses a one-way random effects model in a cross-classification credibility model that determines the credibility score using estimated variance components. Dannenburg et al. (1996) use the general variance components models of which this is a special case. (See also Goulet, 1998.)

Analysis of variance can be put into the context of the insurance model as follows: Consider an insurance company with I groups of policies. Suppose further that there are n_i individuals from group i who have a claim within a single period (a month, quarter, or year, say). For $i = 1, 2, \dots, I$, the claim amount, Y_{iu} , associated with individual u in group i , is modeled as

$$Y_{iu} = \mu + \alpha_i + e_{iu}, \quad u = 1, \dots, n_i, \quad (3)$$

where μ represents the mean over all groups and α_i represents the amount that the mean of the i th group varies from this overall mean, α_i s are mutually independent random variables mean zero and variance σ_1^2 , and the e_{iu} s are mutually independent random variables mean zero and variance σ_0^2 . We also assume that α_i and e_{iu} are mutually independent.

If an assumption of normality of the distribution of α_i and e_{iu} were added to equation (3), this would be the standard formulation of the random one-way ANOVA model. This assumption is unnecessary to form the Bühlmann-Straub credibility premium.

Equation (3) implies that the unconditional expected value of Y_{iu} is μ . Conditional on α_i , however, the expected value of Y_{iu} is $\mu + \alpha_i$. It is the past experience that provides the basis for improving our estimate of the expected value of Y_{iu} , for each group by providing information regarding α .

In the ANOVA model of equation (3), the credibility factor is easy to estimate if we use the method of moments estimate of the variance components as suggested by Venter (1990). The method of moments estimate of σ_1^2 is referred to in the European literature as \hat{a} . Other than simplicity and unbiasedness, this method of estimation has no known optimality properties. Other estimates of σ_1^2 exist with optimality properties, however (see Goulet, 1998; and DeVyllder and Goovaerts, 1992). We will use the simple method of moments estimator.

The following notation is used:

$$N = \sum_{i=1}^t n_i; \quad (4)$$

$$\begin{aligned} \bar{Y}_i &= \text{Average of all observations in group } i; \\ &= \frac{\sum_{u=1}^{n_i} Y_{iu}}{n_i}; \end{aligned} \quad (5)$$

$$\begin{aligned} \bar{Y}_{..} &= \text{Average of all observations, across all groups;} \\ &= \frac{1}{N} \sum_{i=1}^t \sum_{u=1}^{n_i} Y_{iu}; \end{aligned} \quad (6)$$

$$s_i^2 = \frac{1}{n_i - 1} \sum_{u=1}^{n_i} (Y_{iu} - \bar{Y}_i)^2 \quad (7)$$

$$\text{MSE} = \frac{1}{N - t} \sum_{i=1}^t (n_i - 1) s_i^2, \quad (8)$$

$$\text{MS}(\alpha) = \frac{1}{t - 1} \sum_{i=1}^t n_i (\bar{Y}_i - \bar{Y}_{..})^2. \quad (9)$$

The last two expressions are referred to as the mean square for error (MSE) and the mean square for groups ($\text{MS}(\alpha)$), respectively. The expected values of these mean squares are:⁵

$$E[\text{MSE}] = \sigma_0^2$$

and

$$E[\text{MS}(\alpha)] = \sigma_0^2 + n_0 \sigma_1^2,$$

where

$$n_0 = \frac{N^2}{t - 1} \left(1 - \sum_{i=1}^t \frac{n_i^2}{N^2} \right). \quad (10)$$

In Bühlmann' notation, σ_0^2 is the expected value of the process variance and σ_1^2 is the variance of the hypothetical means. Thus, Bühlmann's k is given as

⁵For a derivation of $E[\text{MSE}]$ and $E[\text{MS}(\alpha)]$ see Scheffé (1959, Chapter 3) or Neter, Wasserman, and Craig (1990, Chapters 14, pages 543-546).

$$k = \frac{n_0 \times \text{MSE}}{\text{MS}(\alpha) - \text{MSE}}.$$

From these expectations we can calculate the following method of moments estimators for the variance components:

$$\hat{\sigma}_0^2 = \text{MSE},$$

and

$$\hat{\sigma}_1^2 = \frac{\text{MS}(\alpha) - \text{MSE}}{n_0}. \quad (11)$$

Thus, for the simple one-way model in equation (3), the Bühlmann-Straub credibility factor, Z , given in equation (2) becomes

$$\begin{aligned} Z_i &= \frac{n_i}{n_i + k}, \\ &= \frac{n_i}{n_i + \hat{\sigma}_0^2 / \hat{\sigma}_1^2} \\ &= \frac{n_i \hat{\sigma}_1^2}{n_i \hat{\sigma}_1^2 + \hat{\sigma}_0^2}, \end{aligned} \quad (12)$$

which can be rewritten as

$$Z_i = \frac{\text{MS}(\alpha) - \text{MSE}}{\text{MS}(\alpha) + \left(\frac{n_0}{n_i} - 1\right) \times \text{MSE}}. \quad (13)$$

Most analysis of variance routines calculate MSE and $\text{MS}(\alpha)$. Only the number of observations in the i th group, n_i , and the value of n_0 need to be determined.

The credibility factor is different for each group depending on the value of n_i . As n_i increases, Z_i goes to unity and the group becomes fully credible. On the other hand, as $\hat{\sigma}_1^2$ increases, indicating a high degree of variability from group to group, Z_i approaches unity and the group becomes fully credible. When $\hat{\sigma}_1^2$ is small relative to $\hat{\sigma}_0^2$ and/or n_i is small relative to n_0 , Z_i drops below unity and the group experience is less credible. In this case the compromise estimate borrows more strength from the experience of other groups.

Equation (13) provides a simple calculation of the credibility factor using output from ANOVA routines. Many times, however, the data have been summarized so that for each group i only the observed pure premium, say \bar{Y}_i , the number insured, n_i , and the standard deviation, s_i , are known. In this case the formulas can be used by first observing that

$$\bar{Y}_{..} = \sum_{i=1}^t \frac{n_i}{N} \bar{Y}_i. \quad (14)$$

Thus, $MS(\alpha)$ is calculated as given in equation (9) using $\bar{Y}_{..}$ as given in equation (14). Rearranging the terms in equation (9) yields a formula that is often easier to use. Explicitly,

$$MS(\alpha) = \frac{1}{t-1} \left(\sum_{i=1}^t n_i \bar{Y}_i^2 - N \bar{Y}_{..}^2 \right). \quad (15)$$

Second, MSE is calculated as in equation (8).

The credibility factors Z_i can be calculated using equation (13) where the MSE is given by equation (8) and $MS(\alpha)$ is calculated using equation (15) with $\bar{Y}_{..}$ as defined in equation (14).

3 Calculation of Z via Computer Programs

3.1 Individual Data Case

To illustrate the formulas and computer programs we consider the hypothetical data given in Table 1. The data sets are small and would not be seriously considered as reliable insurance experience. With such small data sets, however, the details of calculations are more apparent. The data in Table 1 represent four hypothetical groups with claims for each group. We wish to determine the credibility factors for each group assuming that the four groups represent the entire experience of interest for the insurer.

Table 2 gives the EXCEL⁶ output for a one-way analysis of variance of the data in Table 1. To obtain this analysis we perform the following steps:

⁶EXCEL is a registered trademark of: Microsoft Corporation, One Microsoft Way, Redmond WA 98052-6399, USA.

Table 1
Hypothetical Individual Cost Data
For Four Groups of Insureds for a Single Year

Groups			
1	2	3	4
1550	1879	1440	1014
1325	2028	1601	1231
1417	2150	1790	1487
1824	2245	1852	1491
2138	2516	1998	
	2918	2081	
		2171	

Step 1: Click the *Data Analysis* menu selection under *Tools*;

Step 2: We then click *One-Way*;

Step 3: As each column represents a different group, we indicate the *Grouped by Columns* option and then proceed.

The output consists of one table (Table 2) with two panels, Panel A and Panel B. The first column in Panel A lists the group name. The second gives the value of n_i for group i , where i indicates the column of the group data. The fourth column gives \bar{Y}_i for group i as given by equation (5). The fourth column of Panel B lists the $MS(\alpha)$ in the first row and the MSE in the second row.

Using the second column of Table 2, Panel A we calculate n_0 using equation (10). For this equation $t - 1 = 4 - 1 = 3$. The other components of the equation are given as:

$$N = 22$$

$$\sum n_i^2 = 126, \quad \text{and}$$

$$n_0 = (22^2 - 126)/(22 * 3) = 5.4242.$$

Table 2
Output from Excel Program of the
One-Way ANOVA Analysis of the Data in Table 1

Panel A: ANOVA Single Factor (Summary)						
Groups	Count (n_i)	Sum	Average (\bar{Y}_i)	Variance		
Group 1	5	8254	1650.800	109582.70		
Group 2	6	13736	2289.333	140929.50		
Group 3	7	12933	1847.571	68661.60		
Group 4	4	5223	1305.750	52624.92		
Panel B: ANOVA						
Source of Variation	SS	df	MSE	F-Value	P-Value	F-Crit
Between Groups	2527409	3	842469.6	8.853487	0.000805	3.159911
Within Groups	1712823	18	95156.81			
Total	4240231	21				

Notes: SS = Sum of Squares; *MSE(α) = Between Groups MSE; F-value = Test statistic to test whether mean costs are the same across groups under the ANOVA assumptions; P-value = Probability of a value greater than or equal to the F-value assuming the means are the same; F-Crit = The value which, if it is exceeded by the F-value, there is statistical evidence that the mean costs differ from between groups.

Using these values we calculate the Z_i for each group using equation (13). Explicitly, for group 1 we have

$$\begin{aligned} Z_1 &= \frac{842469.6 - 95156.81}{842469.6 + \left(\frac{5.4242}{5} - 1\right) \times 95156.81} \\ &= 0.878631 \end{aligned}$$

Thus, the credibility score for group 1 is about 87.9 percent. Relative to the complete set of data available, the data on group 1 are relatively credible—there is little difference between the compromise estimate of the group pure premium and the estimate using the observed average of the group.

3.2 Grouped Data Case

Suppose that only the summary data consisting of n_i , \bar{Y}_i , and s_i^2 for each group are available (columns (2), (4), and (5) of Table 2, Panel A). In this case we can use equations (15) and (8) to calculate the components of equation (13). Explicitly we make the following calculations. First from equation (14) we have

$$\begin{aligned} \bar{Y} &= (5 \times 1650.8 + 6 \times 2289.333 + 7 \times 1847.571 + 4 \times 1305.75) / 22 \\ &= \frac{40146}{22} \\ &= 1824.818182. \end{aligned}$$

Using these in equation (15) we obtain

$$\begin{aligned} MS(\alpha) &= \frac{75786559.41 - 73259150.73}{3} \\ &= \frac{2527408.68}{3} \\ &= 842469.56 \end{aligned}$$

This is close to the value given in Table 2, Panel B (row (1), column (4)). The difference is due to roundoff error.

Calculation of MSE follows similarly using equation (8). Explicitly, we get

$$\begin{aligned} MSE &= \frac{1712822.66}{18} \\ &= 95156.81 \end{aligned}$$

These results can be used to calculate the credibility scores as before.

Computer code for the same calculations using SAS are given in the appendix; no code is provided for SPSS.⁷

4 Discussion

We have illustrated how the Bühlmann-Straub credibility factors can be calculated using one-way ANOVA statistical routines common in many computer programs. In order to form such scores the mean squares reported in the ANOVA tables must be used as given in equation (13). Under certain situations estimated $MS(\alpha)$ can be negative. In this case the value of $Z_i = 0$ is used. This reduces the bias of the compromise estimate as shown by Morris (1983).

References

- Bühlmann, H. "Experience Rating and Credibility." *ASTIN Bulletin* 4 (1967): 199-207.
- Bühlmann, H. *Mathematical Methods in Risk Theory*. New York, N.Y.: Springer-Verlag, 1970.
- Bühlmann, H. and Straub, E. "Glaubwürdigkeit Für Schadensätze." *Mitteilungen der Vereinigung Schweizerischer Versicherungsmathematiker* 70 (1970): 111-133.
- Dannenburg, D. "Cross Classification Credibility Models." *Transactions of the 25th International Congress of Actuaries* 4 (1995): 1-35.
- Dannenburg, D.R., Kaas, R. and Goovaerts, M.J. *Practical Actuarial Credibility Models*. Amsterdam, Holland: Institute of Actuarial Science and Econometrics, University of Amsterdam, 1996.
- DeVylder, F. and Goovaerts, M. "Optimal Parameter Estimation Under Zero-Excess Assumptions in the Bühlmann-Straub Model." *Insurance: Mathematics and Economics* 11 (1992): 167-171.
- Ericson, W.A. "On the Posterior Mean and Variance of a Population Mean." *Journal of the American Statistical Association* 65 (1970): 649-652.

⁷SAS is a registered trademark of: SAS Institute Inc., Cary, NC 27512-8000, USA; and SPSS is a registered trademark of: SPSS Inc., 444 North Michigan Avenue, Chicago IL 60611, USA.

- Goulet, V. "Principles and Applications of Credibility Theory." *Journal of Actuarial Practice* 6 (1998): 5-62.
- Herzog, T.N. *Credibility Theory*. Winsted, Conn.: ACTEX Publications, 1996.
- Hossack, I.B., Pollard, J.H. and Zehnwirth, B. *Introductory Statistics with Applications in General Insurance*. Cambridge, England: Cambridge University Press, 1983.
- Klugman, S. "Credibility for Classification Ratemaking via the Hierarchical Normal Linear Model." *Proceedings of the Casualty Actuarial Society* 74 (1987): 272-321.
- Longley-Cook, L.H. "An Introduction to Credibility Theory." *Proceedings of the Casualty Actuarial Society* 49 (1962): 194-221.
- Morris, C. and van Slyke, L. "Empirical Bayes Methods for Pricing Insurance Classes." *Proceedings of the Business and Economics Statistics Section, American Statistical Association* (1978): 579-582.
- Morris, C.N. "Parametric Empirical Bayes Inference: Theory and Applications." *Journal of the American Statistical Association* 78 (1983): 47-65.
- Mosteller, F. and Tukey, J.W. "Purposes of Analyzing Data That Come in a Form Inviting Us to Apply Tools From the Analysis of Variance." In *Fundamentals of Exploratory Analysis of Variance*. New York, N.Y.: John Wiley & Sons, 1991.
- Neter, J., Wasserman, W. and Craig, A.T. *Applied Linear Statistical Models*. 3rd Edition. Boston, Mass.: Irwin, 1990.
- Norberg, R. "The Credibility Approach to Ratemaking." *Scandinavian Actuarial Journal* (1979): 181-221.
- Philbrick, S.W. "An Examination of Credibility Concepts." *Proceedings of the Casualty Actuarial Society* 68 (1981): 195-219.
- Rubin, D.B. "Using Empirical Bayes Techniques in the Law School Validity Studies." *Journal of the American Statistical Association* 75 (1980): 801-827.
- Scheffé, H. *The Analysis of Variance*. New York, N.Y.: John Wiley & Sons, 1959.
- Tsutakawa, R.K. "Mixed Model for Analyzing Geographic Variability in Mortality Rates." *Journal of the American Statistical Association* 83 (1988): 37-42.
- Venter, G.G. "Structured Credibility in Application—Hierarchical, Multi-dimensional and Multivariate Models." *Actuarial Research Clearing*

House (ARCH) 2 (1985): 267-308. (ARCH is a publication of the Society of Actuaries, Schaumburg, Ill.)

Venter, G.G. "Credibility." In *Foundations of Casualty Actuarial Science*. Arlington, Va.: Casualty Actuarial Society, 1990.

Appendix

The codes for making credibility calculations using SAS for the data in Table 1 are given below. First we use the individual data. We have used the *cards* option. In practice one would read a data file. Below we give the code for grouped data. In both cases the amount of work to get the SAS code seems long relative to the simple problem considered. For longer, more practical problems, however, the benefits of SAS routines are more apparent.

```
DATA costs;
  INFILE cards;
  INPUT cost group;
  CARDS;
1550 1
1325 1
1417 1
1824 1
2138 1
1879 2
2028 2
2150 2
2245 2
2516 2
2918 2
1440 3
1601 3
1790 3
1852 3
1998 3
2081 3
2171 3
```

```

1014 4
1231 4
1487 4
1491 4
;
RUN;

/** Getting number of individuals per group ***/
PROC SQL;
  CREATE TABLE counts AS
  SELECT DISTINCT group,count(group) AS number
  FROM costs
  GROUP BY group;

/** Calculating n_not ***/
PROC SQL;
  SELECT (sum(number)-(sum(number**2)/sum(number)))
  /(count(number)-1)
  INTO :n_not
  FROM counts;

/** Calculating MSE, MSA ***/
PROC ANOVA DATA=costs OUTSTAT=results NOPRINT;
  CLASS group;
  MODEL cost=group;
RUN;

DATA _null_;
  SET results;
  mean_sqr=ss/df;
  SELECT (_source_);
    WHEN ("ERROR") CALL SYMPUT("MSE",mean_sqr);
    WHEN ("GROUP") CALL SYMPUT("MSA",mean_sqr);
  END;
RUN;

/** Calculating credibilities ***/
DATA creds;
  SET counts;
  cred=(&MSA-&MSE)/(&MSA+(&n_not/number-1)*&MSE);
  KEEP group cred;
RUN;

```

```

PROC PRINT NOOBS DATA=creds;
  TITLE 'Credibility Factors for Individual Data';
RUN;

/*****
  USING GROUPED DATA
*****/

DATA grouped;
  INFILE cards;
  INPUT group number avg_cost var_cost;
  CARDS;
1 58 1666 49597893
2 115 5051 216276545
3 81 4670 193990984
4 108 8966 757144094
;
RUN;

/*** Calculating n_not and the overall mean ***/
PROC SQL;
  SELECT (sum(number)-(sum(number**2)/sum(number)))
         /((count(number)-1),
         sum(avg_cost*number)/sum(number))
  INTO :n_not,:y_bar2
  FROM grouped;

/*** Calculating MSE, MSA ***/
PROC SQL;
  SELECT 1/((count(group)-1)*(sum(number*avg_cost**2)
  -sum(number)*&y_bar2**2),
  1/((sum(number)-count(group))*sum((number-1)
  *var_cost))
  INTO :msa,:mse
  FROM grouped;

/*** Calculating credibilities ***/
DATA creds;
  SET grouped;
  cred=(&msa-&mse)/(&msa+(&n_not/number-1)*&mse);
  KEEP group cred;

```


time period. Recent examples of collective risk modeling in insurance include Butler, Gardner, and Gardner (1998); Butler and Worall (1991); and Cummins and Tennyson (1996).

The stochastic structure is two-pronged: both the size of the individual claims and the number of claims are considered random variables. Specifically, let S denote the aggregate claims random variable, i.e.,

$$S = \sum_{i=1}^N X_i \quad (1)$$

where N is the number of claims and X_i is the size of the i th individual claim. The X_i s are assumed to be mutually independent and identically distributed (i.i.d.) and are mutually independent of N . In the literature equation (1) is referred to as a compound random variable; see, for example, Bowers et al. (1997, Chapter 12).

Theoretically, the distribution of S can be obtained from equation (1) as follows:

$$\Pr[S \leq s] = \sum_{n=0}^{\infty} p_n F^{*n}(s) \quad (2)$$

where $p_n = \Pr[N = n]$ and $F^{*n}(s) = \Pr[X_1 + \dots + X_n \leq s]$, i.e., $F^{*n}(s)$ is the n th convolution of the X_i s, with $F(x) = F^{*1}(x)$ being the cumulative distribution function of X_1 .

The difficulty in using equation (2), however, often lies in calculating $F^{*n}(s)$. Thus, approximations are frequently used. There are several approximations used by actuaries, including discretizing the claim size distribution (Panjer 1981); using the Wilson-Hilferty approximation or Haldane Type A approximation (Pentikäinen, 1987); and, of course, the normal approximation. See Panjer and Willmot (1992, Chapter 6) and Bowers et al. (1997, Chapters 2 and 12) for a discussion of the actuarial approaches. Other methods such as the Edgeworth expansion (Feller, 1971) or the conjugate density method (Esscher, 1932) have been applied.

The methods mentioned above provide good approximations near the center of the distribution but can be slow or inaccurate for calculating tail probabilities of the form $\Pr[S > s]$ (for large values of s). For a discussion of the tail behavior of aggregate distributions; see Panjer and Willmot (1992, Chapter 10). A fairly quick and accurate method of calculating tail probabilities is the so-called saddlepoint approximation.

Since their introduction by Daniels (1954) saddlepoint approximations have been utilized to approximate tail probabilities in a variety of situations; see, for example, Goutis and Casella (1999), Huzurbazar (1999), Butler and Sutton (1998), Tsuchiya and Konishi (1997), and Wood, Booth, and Butler (1993). Field and Ronchetti (1990) document the accuracy of these procedures for small sample sizes (even of sample size one). In this paper a saddlepoint approximation is developed for $\Pr[S > s]$ and is applied to specific examples.

2 The Saddlepoint Approximation

The key assumption in the saddlepoint approximation is the assumption of the existence of the moment-generating functions corresponding to X_i and N , which are denoted by $M_X(\theta)$ and $M_N(\theta)$, respectively, where θ is a real valued parameter.¹ The moment-generating function of S , $M_S(\theta)$, is then given by

$$\begin{aligned} M_S(\theta) &= E[e^{\theta S}] \\ &= E[E[e^{\theta S} | N]] \\ &= M_N(\log(M_X(\theta))). \end{aligned} \quad (3)$$

Equation (3) can be used to derive the well-known results on the moments of compound sums of i.i.d. random variables:

$$\mu_S = E[S] = E[N]E[X_1] \quad (4)$$

$$\sigma_S^2 = \text{Var}[S] = \text{Var}[N](E[X_1])^2 + E[N]\text{Var}[X_1]. \quad (5)$$

The saddlepoint approximation for the tail probability $\Pr[S > s]$ is adapted from Field and Ronchetti (1990) for sample size one. First let T denote the standardized random variable

$$T = \frac{S - \mu_S}{\sigma_S}$$

¹The moment-generating function of a random variable Z is defined as

$$M_Z(\theta) = E[e^{\theta Z}], \quad \theta > 0.$$

where μ_S and σ_S and the mean and standard deviation of S respectively (which can be obtained from equations (4) and (5)). The moment-generating function for T is easily seen to be:

$$M_T(\theta) = e^{-\mu_S/\sigma_S} M_S(\theta/\sigma_S). \quad (6)$$

For a fixed value of s , let $t = (s - \mu_S)/\sigma_S$ and let β be the solution to the equation

$$M_T'(\beta) = t M_T(\beta) \quad (7)$$

where the ' denotes differentiation with respect to θ . Note that β is a function of t . Further, let

$$c = \frac{e^{\beta t}}{M_T(\beta)} \quad (8)$$

and

$$\sigma^2 = \frac{M_T''(\beta)}{M_T(\beta)} - t^2. \quad (9)$$

The saddlepoint approximation for $P(S > s)$ is:

$$\Pr(S > s) \approx 1 - \Phi(\sqrt{2 \ln(c)}) + \frac{1}{c\sqrt{2\pi}} \left[\frac{1}{\beta\sigma} - \frac{1}{\sqrt{2 \ln(c)}} \right] \quad (10)$$

where $\Phi(\cdot)$ is the standard normal distribution function, and c and σ are defined in equations (8) and (9).

In practice, once s is chosen and t is computed, equation (7) is solved numerically using a technique such as Newton's method or the secant method; see, for example, Burden and Faires (1997, Chapter 2).

3 Examples

The saddlepoint approximations of tail probabilities are now applied to several specific collective risk models. These saddlepoint approximations are compared to the Haldane Type A and the normal approximations, and the exact probabilities. The exact calculations are found

by simulation using 10,000 repetitions, which gives accuracy to four decimal places.

If X has mean μ_X , standard deviation σ_X , and coefficient of skewness γ_X , then the Haldane Type A approximation is as follows:

$$Pr[X \leq x_0] \approx \Phi \left[\frac{\left((1 + r\tilde{x}_0)^h - \mu(h, r) \right)}{\sigma(h, r)} \right] \tag{11}$$

where

$$\tilde{x}_0 = \frac{(x_0 - \mu_X)}{\sigma_X} \tag{12}$$

$$r = \frac{\sigma_X}{\mu_X} \tag{12}$$

$$h = 1 - \frac{\gamma_X}{3r} \tag{13}$$

$$\mu(h, r) = 1 - \frac{1}{2}h(1-h) \left[1 - \frac{1}{4}(2-h)(1-3h)r^2 \right] r^2 \tag{14}$$

$$\sigma(h, r) = hr \sqrt{1 - \frac{1}{2}(1-h)(1-3h)r^2} \tag{15}$$

The Haldane approximation is chosen because Pentikäinen's (1987) results show it to be, under certain circumstances, an accurate approximation. Recall that the normal approximation is

$$Pr[X \leq x_0] \approx \Phi[\tilde{x}_0]. \tag{16}$$

The relative errors shown in the tables are calculated as:

$$\text{Relative Error} = \left| \frac{\text{Approximation} - \text{Exact}}{\text{Exact}} \right|.$$

3.1 Light and Medium Tailed Claim Size Distributions

Example 1: X_1 is normal random variables with mean $\mu_X = 100$ and standard deviation $\sigma_X = 10$ while N is Poisson with mean $\lambda = 10$. From equation (3)

$$M_S(\theta) = \exp \left[\lambda \left(\exp \left(\mu_X \theta + \frac{1}{2} \sigma_X^2 \theta^2 \right) - 1 \right) \right]. \tag{17}$$

Table 1
Approximating Tail Probabilities for
The Compound Normal-Poisson Model

t	β	Exact	Relative Error		
			Normal	HALD A	SADP
0.5	0.4637	0.2964	0.0411	0.0039	0.0034
1.0	0.8672	0.1575	0.0074	0.0077	0.0070
1.5	1.2243	0.0750	0.1089	0.0062	0.0087
2.0	1.5445	0.0303	0.2498	0.0125	0.0082
2.5	1.8347	0.0112	0.4469	0.0019	0.0089
3.0	2.1001	0.0036	0.6351	0.0091	0.0084

In this setting the central limit theorem is known to hold for large λ .

Example 2: X_1 is a gamma random variable with a mean of $\mu_X = 100$ and standard deviation $\sigma_X = 10$. N is a negative binomial random variable with mean of $\alpha = 10$ and standard deviation $y = 20$. Here

$$M_S(\theta) = \left[\frac{1-q}{1-q(1-\beta)^{-\delta}} \right]^r \quad (18)$$

where $q = 0.5$, $\mu_X = \beta\delta$, $\sigma_X = \beta\sqrt{\delta}$, $\alpha = rq/(1-q)$ and $y = rq/(1-q)^2$.

Table 2
Approximating Tail Probabilities for
The Compound Gamma-Negative Binomial

t	β	Exact	Relative Error		
			Normal	HALD A	SADP
0.5	0.4284	0.2684	0.1494	0.0961	0.0417
1.0	0.7502	0.1548	0.0252	0.0284	0.0032
1.5	1.001	0.0796	0.1608	0.0515	0.0050
2.0	1.203	0.0375	0.3920	0.6907	0.0027
2.5	1.369	0.0166	0.6265	0.3012	0.0084
3.0	1.508	0.0070	0.8086	0.4571	0.0100

Example 3: X_1 is an inverse Gaussian random variable with mean $\mu_X = 100$ and standard deviation $\sigma_X = 10$. N is Poisson with mean $\lambda = 10$. The moment-generating function for the inverse Gaussian distribution is

$$M_X(\theta) = \exp \left[\left(\frac{\mu_X}{\sigma_X} \right)^2 \left(1 - \left(1 - \frac{2\sigma_X^2\theta}{\mu_X} \right) \right) \right],$$

see Johnson and Kotz (1970, Chapter 15). Hence

$$M_S(\theta) = \exp \left[\lambda \left(\left(\left(\frac{\mu_X}{\sigma_X} \right)^2 \left(1 - \left(1 - \frac{2\sigma_X^2\theta}{\mu_X} \right) \right) \right) - 1 \right) \right]. \quad (19)$$

Table 3
Approximating Tail Probabilities for
The Compound Inverse Gaussian–Poisson Model

t	β	Exact	Relative Error		
			Normal	HALD A	SADP
0.5	0.4537	0.2998	0.0290	0.0153	0.0147
1.0	0.8671	0.1629	0.0258	0.0258	0.0264
1.5	1.2242	0.0775	0.1381	0.0387	0.0413
2.0	1.5444	0.0316	0.2785	0.0285	0.0348
2.5	1.8345	0.0119	0.4790	0.0588	0.0672
3.0	2.0998	0.0038	0.6474	0.0526	0.0526

These examples show that the saddlepoint approximation is superior to the central limit theorem, but seems to be on par with the Haldane approximation in calculating tail probabilities. Next we consider a more difficult setting involving heavy tailed distributions.

4 Heavy Tailed Claim Size Distributions

The saddlepoint approximation requires the existence of the moment-generating function of the claim variable. For heavy tailed distributions, such as the Pareto (the moment-generating function does not exist)

and lognormal (the moment-generating function is not in convenient a closed form), an approximation is required. For these problem cases a censoring limit is imposed on the claim distribution.

For cases where the moment-generating function does not exist, the distribution of the claim variable is approximated utilizing an upper tail censoring limit. For small ϵ the censoring limit, L , satisfies $\Pr[X_1 > L] = \epsilon$. Let us define the censored claim random variable as

$$Y_i = \begin{cases} X_i & \text{if } X_i \leq L \\ L & \text{if } X_i > L. \end{cases}$$

The distribution function for the Y_i s is now

$$F_Y(x) = \begin{cases} F(x) & \text{if } x < L \\ 1 & \text{if } x \geq L. \end{cases}$$

The corresponding moment-generating function is

$$M_Y(\theta) = \int_{x=0}^L e^{\theta x} dF(x) + \epsilon e^{\theta L}. \quad (20)$$

The saddlepoint approximation is applied using the censoring moment-generating function in equation (20). This technique is now demonstrated on two examples of heavy tailed claim distributions. In both cases the number of claims is assumed to be Poisson with mean 5.

Example 4: Claims are assumed to follow a lognormal distributed with probability density function (pdf) of X_1 is

$$f(x) = \frac{1}{\sqrt{2\pi\sigma}} \exp\left[-\frac{1}{2} \left(\frac{\ln(x) - \mu}{\sigma}\right)^2\right] \quad -\infty < x < \infty. \quad (21)$$

where $\mu = 0$ and $\sigma = 1$. We assume that $\epsilon = 0.001$, which produces a censoring limit of $L = 59.7697$.

Example 5: Here we assume the claim size follows a Pareto distribution with distribution function given by

$$F(x) = 1 - \frac{1}{(1+x)^3}.$$

Table 4
Approximating Tail Probabilities for
The Compound Lognormal-Poisson Model

<i>t</i>	β	Exact	Relative Error		
			Normal	HALD A	SADP
0.5	0.7251	0.1628	0.8950	0.5565	0.0498
1.0	0.9501	0.0630	1.5190	1.3016	0.0825
1.5	1.0512	0.0241	1.7718	2.3361	0.0622
2.0	1.2001	0.0108	1.1111	1.0463	0.1574
2.5	1.4211	0.0047	0.3191	5.5319	0.3404

Again, $\epsilon = 0.001$, and this produces a censoring limit of $L = 9.0$.

As in the previous section, normalized tail probabilities and the saddlepoint approximations are compared to the exact values as obtained by simulation. These computations are listed in Tables 4 and 5.

Table 5
Approximating Tail Probabilities for
The Compound Pareto-Poisson Model

<i>t</i>	β	Exact	Relative Error		
			Normal	HALD A	SADP
0.5	0.6959	0.1664	0.8540	0.6280	0.0313
1.0	0.9880	0.0688	1.3067	1.2456	0.1933
1.5	1.1623	0.0327	1.0428	1.5199	0.1804
2.0	1.2842	0.0165	0.3818	1.5091	0.0727
2.5	1.3772	0.0094	0.3404	1.1064	0.1596

For the heavy tailed distributions, the saddlepoint approximation is superior to the central limit theorem and the Haldane approximation in calculating tail probabilities.

References

- Bowers, N.L., Gerber, H.U., Hickman, J.C., Jones, D.A. and Nesbitt, C.J. *Actuarial Mathematics*, (2nd edition). Schaumburg, Ill.: Society of Actuaries, 1997.
- Burden, R.L. and Faires, J.D. *Numerical Analysis*, (6th edition). New York, N.Y.: Brooks/Cole Publishing Company, 1997.
- Butler, R.J., Gardner, H. and Gardner, H. (1998). "Workers Compensation Costs When Maximum Benefits Change." *Journal of Risk and Uncertainty* 15 (1998): 259-269.
- Butler, R. and Sutton, R. "Saddlepoint Approximation for Multivariate Cumulative Distribution Functions and Probability Computations in Sampling Theory and Outlier Testing." *Journal of the American Statistical Association* 19, no. 442 (1998): 596-604.
- Butler, R.J. and Worall, J.D. "Claim Reporting and Risk Bearing Moral Hazard in Workers Compensation." *Journal of Risk and Insurance* 53 (1991): 191-204.
- Cummins, J.D. and Tennyson, S. "Moral Hazard in Insurance Claiming: Evidence from Automobile Insurance." *Journal of Risk and Uncertainty* 12 (1996): 29-50.
- Daniels, H.E. "Saddlepoint Approximations in Statistics." *Annals of Mathematical Statistics* 25 (1954): 631-650.
- Esscher, F. "On the Probability Function in Collective Risk Theory." *Scandinavian Actuarial Journal* 15 (1932): 175-195.
- Feller, W. *An Introduction to Probability and Its Application, Volume 2*, (2nd edition). New York, N.Y.: Wiley and Sons, 1971.
- Field, C. and Ronchetti, E. *Small Sample Asymptotics*. IMS Lecture Notes-Monograph Series 13. Hayward, Calif.: Institute of Mathematical Statistics, 1990.
- Goutis, C. and Casella, G. "Explaining the Saddlepoint Approximation." *The American Statistician* 53, no. 3 (1999): 216-224.
- Huzurbazar, S. "Practical Saddlepoint Approximations." *The American Statistician* 53, no. 3 (1999): 225-232.
- Johnson, N.L. and Kotz, S. *Distributions in Statistics: Continuous Univariate Distributions, Volume 1*. Boston, Mass.: Houghton Mifflin Company, 1970.
- Panjer, H.H. "Recursive Evaluation of a Family of Compound Distributions." *ASTIN Bulletin* 12 (1981): 22-26.

- Panjer, H.H. and Willmot, G.E. *Insurance Risk Models*. Schaumburg, Ill.: Society of Actuaries, 1992.
- Pentikäinen, T. "Approximate Evaluation of the Distribution Function of Aggregate Claims." *ASTIN Bulletin* 17 (1987): 15-40.
- Tsuchiya, T. and Konishi, S. "General Saddlepoint Approximation and Normalizing Transformations for Multivariate Statistics." *Communications in Statistics, Part A—Theory and Methods* 26, no. 11 (1997): 2541-2563.
- Wood, A.T., Booth, J.G. and Butler, R.W. "Saddlepoint Approximation to the CDF of Some Statistics with Nonnormal Limit Distributions." *Journal of the American Statistical Association* 88, no. 442 (1993): 680-686.

