

Modeling Loss Emergence and Settlement Processes

CAS Loss Simulation Model Working Party Summary Report

Abstract: The Loss Simulation Model Working Party (LSMWP) has been charged by the Committee on Dynamic Risk Modeling Committee with creating a simulation model of the processes of loss emergence and settlement, commonly known as loss development, that underlie the loss "triangles" and other statistics used to estimate loss reserves. The goal was to create a tool that researchers could use to generate claims that would be summarized into loss development triangles and complete rectangles which could then be used to test loss reserving methods and models.

Motivation. Actuaries need tools that will enable them to better understand the underlying loss development process and will aid them in determining what methods and models work best in different reserving situations.

Method. The LSMWP first developed a prototype model that met its basic objectives. It then engaged a consultant to develop an open source model that could be applied and further developed by the CAS membership.

Results. The open source model developed by the consultant is documented in this paper, along with the testing which validated the model.

Conclusions. A valuable tool has been created for use by the actuarial community in its research work on reserving methods and models. We encourage actuaries to perform additional validation tests of the model and contribute model enhancements.

Availability. This report and the associated model are available by clicking on the LSMWP link on the Committee on Dynamic Risk Modeling (DRM) page on the CAS web site (www.casact.org/research/drm).

Keywords. loss reserving; simulation model; reserve variability; reporting pattern; payment pattern; open source model.

TABLE OF CONTENTS

1. INTRODUCTION.....	4
1.1 RESEARCH CONTEXT.....	4
1.2 OBJECTIVE.....	4
1.3 DISCLAIMER.....	5
1.4 OUTLINE.....	5
2. SURVEY OF EXISTING LITERATURE.....	5
3. STATISTICAL TESTS OF SIMULATED MODEL OUTPUT.....	6
3.1 PP-PLOTS.....	7
3.2 THE MEYERS APPROACH.....	8
3.3 ADDITIONAL COMMENTS.....	9
4. BASIC FEATURES IN THE PROTOTYPE MODEL.....	9
5. DOCUMENTATION OF OPEN SOURCE MODEL.....	11
6. TESTING DETAILED OUTPUT AND FITTING THE MODEL.....	11
6.1 DATA FORMATS.....	12
6.2 DESCRIPTION OF TESTING.....	13
6.2.1 Test of Elementary Frequencies, Trend, and Zero-modification.....	13
6.2.2 Testing Size of Loss Distributions.....	19
6.2.3 Testing Correlated Frequencies.....	30
7. POTENTIAL APPLICATIONS AND MODEL ENHANCEMENTS.....	34
8. CONCLUSIONS.....	35
ACKNOWLEDGMENT.....	43
SUPPLEMENTARY MATERIAL.....	43
WORKING PARTY OVERSIGHT, MODEL BUILDING AND TESTING TEAM MEMBERS.....	43
APPENDIX A.....	47
1) HOW TO INSTALL THE PUBLIC LOSS SIMULATION MODEL.....	47
2) HOW TO RUN THE PUBLIC LOSS SIMULATION MODEL.....	50
2.1 System Overview.....	50
2.2 Simulation Project.....	54
2.3 Overall Simulation Properties.....	55
2.4 Line Level Properties.....	56
2.5 Type Level Properties.....	61
2.6 Run Simulation.....	69
3) SIMULATION EXAMPLE.....	75

Modeling Loss Emergence and Settlement Processes

<u>3.1 Test Objective: Test univariate ultimate severities.....</u>	<u>76</u>
<u>3.2 Simulation Project Level Parameters Setup.</u>	<u>76</u>
<u>3.3 Line Level Parameters Setup.</u>	<u>77</u>
<u>3.4 Coverage Level Parameters Setup.</u>	<u>77</u>
<u>3.5 Coverage Level Severity Distributions.....</u>	<u>79</u>
<u>3.6 Number of Iterations: 100.....</u>	<u>80</u>
4) DISTRIBUTIONS USED IN THE PUBLIC LOSS SIMULATOR AND THEIR PARAMETERIZATIONS.....	81
<u>BETA.....</u>	<u>81</u>
<u>EXPONENTIAL.....</u>	<u>82</u>
<u>GAMMA.....</u>	<u>83</u>
<u>GAUSSIAN (NORMAL).....</u>	<u>84</u>
<u>GEOMETRIC.....</u>	<u>85</u>
<u>LOGNORMAL.....</u>	<u>86</u>
<u>NEGATIVEBINOMIAL.....</u>	<u>87</u>
<u>PARETO.....</u>	<u>88</u>
<u>POISSON.....</u>	<u>89</u>
<u>WEIBULL.....</u>	<u>90</u>
<u>APPENDIX B.....</u>	<u>91</u>
<u>6.2.1 TEST OF ELEMENTARY FREQUENCIES, TREND, AND ZERO-MODIFICATION.....</u>	<u>91</u>
<u>APPENDIX TO SECTION 6.2.2 TEST OF SEVERITY DISTRIBUTIONS.....</u>	<u>95</u>
<u>APPENDIX TO SECTION 6.2.3 -- TESTING CORRELATED FREQUENCIES.....</u>	<u>109</u>
<u>APPENDIX C.....</u>	<u>111</u>

1. INTRODUCTION

The work of the LSMWP evolved in the following stages:

- Survey of existing literature and preparation of an appropriate bibliography. This bibliography includes papers that provide guidance concerning how alternative reserving methods and models may be tested.
- Develop methods of testing simulated data to assure the data generated represent “real” results that could be produced by a company and, therefore, used in testing the relative worth of various reserve models.
- Develop prototype model that has key features desired in final open source model.
- Develop and document open source model through work with a consultant.
- Test open source model by evaluating output to determine if it could be distinguished from real data and if it is consistent with actuarial assumptions.
- Develop procedures to enable a user to develop model parameters that fit his/her own data

The model testing and the procedures for fitting the model to a user’s own data are contained in two sections. The model generates both detailed claim transaction data and aggregate data in the form of “triangles.” Section 3 discusses testing of the aggregate data features. Section 6 tests the detailed data and also describes procedures for fitting the model to actual detailed data.

The sections below elaborate on the above developmental stages.

1.1 Research Context

Actuaries need tools that will enable them to better understand the underlying loss development process and will aid them in determining what methods and models work best in different reserving situations.

1.2 Objective

Our objective is to make a valuable tool available for use by the actuarial community in its research work on reserving methods and models. We encourage actuaries to perform additional validation tests of the model and contribute model enhancements.

1.3 Disclaimer

While this paper is the product of a CAS working party, its findings do not represent the official view of the Casualty Actuarial Society. Moreover, while we believe the approaches we describe are very good examples of how to address the issue of loss development, we do not claim they are the only acceptable ones.

While we made a reasonable effort to test and validate the new open source model, users should do their own independent testing and validation. We cannot assure users that the model is completely valid or free from error. We look forward to working with users to correct any errors and enhance it so that it is more useful to the actuarial community.

1.4 Outline

The remainder of the paper proceeds as follows. Section 2 will discuss our survey of existing literature and present our bibliography. Section 3 will discuss our recommended method of testing simulated data to assure the data generated represent “real” results that could be produced by a company and, therefore, used in testing the relative worth of various reserve models. Section 4 will present the model features that were developed in the prototype model and ultimately programmed in the new open source model. Section 5 will present summary features and documentation of the new open source model. Finally, Section 6 will present test results of the new open source model, both evaluating output to determine if it is consistent with actuarial assumptions and fitting data to common distributions. Section 7 will present a list of potential applications of the model as well as model enhancements that we hope the actuarial community will undertake, and Section 8 will present our conclusions.

2. SURVEY OF EXISTING LITERATURE

As we began our survey work, we decided to associate our findings with four basic categories of articles that might prove of value to those who ultimately use the model developed from the LSMWP efforts:

General Interest

Building Models

Testing Simulated Data

Testing Reserve Methods.

Modeling Loss Emergence and Settlement Processes

Our final list of readings for consideration is provided as the Bibliography at the end of this paper. While we do not represent this to be an exhaustive list, we think we have a good balance of the four areas.

In most cases, the articles suggested reference additional publications for consideration of the interested reader as she/he delves deeper into one or more aspects of the issues associated with development of a model generating loss data and testing the simulated output from that model.

3. STATISTICAL TESTS OF SIMULATED MODEL OUTPUT

We began by asking questions that one would ask when inspecting the output from any model supposedly representing potential loss results for a line of business:

- Do aggregate results appear reasonable in terms of losses (loss ratios) generated in a given year and the distribution of losses (loss ratios) over multiple years?
- Do incremental changes appear “logical” with respect to the line of business being modeled?
- Do the distributions (claim count, severity, etc.) produced appear “reasonable”?

While these questions would properly be asked by someone with reserving experience when presented with real data by line of business, this type of review lacks the level of objective analysis required to use data generated from a model for the purpose of evaluating the value of different reserving models (methods).

Depending on the amount of real data available, many modelers suggest holding back data for use in testing a model. In this approach, a model is parameterized based on a subset of the total actual data, then tested against the unused portion of data. This “control” data is compared with what comes out of the model to determine if the model “accurately” represents the real world.

The ability to apply the “control data” approach depends on the amount and granularity of data being modeled. Unfortunately, it is often difficult to amass enough data to apply this approach. Therefore, modelers have, on occasion “created” more control data. Fundamentally, this is done by building a set of data or supplementing it with randomly selected data elements from the initial set of data available to the modeler. There are a host of potential pitfalls with this approach, not the

least of which is making certain you have in fact generated a randomly selected set of data, but it does have its merits.

After reviewing the discussion of statistical testing methods in our bibliography, we chose to recommend a testing approach brought to our attention by Glenn Meyers. Previously, we had received input from experts in model building suggesting we test the simulated data against actual data by comparing the differences at various incremental points of development. If these differences were randomly distributed with a mean of zero, then one could conclude the simulated data would serve as a good proxy for real data.

Mr. Meyers has authored two papers documented in the literature search results that provide detailed explanation of a way to assess the ability of data from a model to represent “real” data:

Estimating Predictive Distributions for Loss Reserve Models, and
Thinking Outside the Triangle.

In the following, we will provide both general and specific insights gained while working with the Meyers approach, along with some more general thoughts related to testing of data.

3.1 PP-Plots

A commonly used approach to testing data output from a model is to use the Kolmogorov-Smirnov (K-S) test. This test is described in the “Loss Models: From Data to Decisions” book by Klugman, Panjer and Willmot ¹ and in Wikipedia.

An approach that is similar to graphically analyzing a sample of data output from the model is to construct PP-Plots. The points $(\frac{i}{n+1}, F_i)$ are plotted for $i=1, \dots, n$. One evaluates whether these points are within a fixed percentage of the 45° line. This is the approach recommended and

¹ [2] in bibliography

documented by Meyers in his papers and subsequently adopted by the LSMWP in its tests of the open source loss simulation model developed for the CAS. We recommend the Meyers papers cited above for the insights shared on interpreting PP-Plots.

3.2 The Meyers Approach

The approach taken by Meyers was as follows:

- Starting with a real set of data, a model is constructed which can be used to generate accident year payments by settlement lags. It is used to construct “typical” accident year (AY) triangles and rectangles (all AYs fully developed) we face during normal reserve reviews.
- Using the model, generate a significant set of simulated AY rectangles (500 in this example) for use in testing output from the model and how well it represents real data.
- Produce as additional output a set of simulated triangles (9 for this example) and randomly insert a real triangle to be tested with the others.
- For this example, parameters for the model were selected to purposely NOT fit the real data. Therefore, we’d expect the real triangle to not pass a goodness of fit test, while the simulated triangles would do so.

The PP-Plot tests were performed as follows:

- A set of 500 simulated rectangles was generated with each simulation including 10 AYs. For each AY, premium and incremental paid losses for each of 10 settlement lags were calculated. (While the premium is not needed for our example, it was an integral part of the model, as it used an expected loss ratio to determine the total losses ultimately paid for an AY.)
- We used Excel commands to reformat the output into a summary of paid losses by settlement lag for the 5000 AYs provided.
- We are given the 10 candidate triangles (9 simulated and one real triangle) to test.
- Candidate data was reformatted, percentiles determined and then sorted for ease of output as PP-Plot graphs. (Note: For a given candidate, the percentile of the accident year emergence for each lag is calculated from the distribution of simulated loss

emergences for the same accident year and lag. Also, for data points outside the percentile range determined by the 500 simulations, a value of zero is substituted.)

- Appendix C, Sheet 1-10: Candidate PP-Plots are provided. We think it's fairly obvious which of the candidate triangles represents the real data.

3.3 Additional Comments

An underlying tenet of the highlighted approach to testing the ability of a model to represent real data is: If real data passes the PP-Plot with the Kolmogorov-Smirnov ("KS") test using data simulated by the model, then we can infer that a formula/method that works well (passes appropriate tests) for data from the loss simulation model will work well on real data. We are comfortable with this assumption but it should not be accepted without some reflection.

4. BASIC FEATURES IN THE PROTOTYPE MODEL

The prototype model contained the following basic features that were ultimately programmed in the new open source model:

- (1) Observation period: We assume that the relevant loss process involves accidents or occurrences between dates t_0 and t_1 . The simulator tracks transactions until accidents are settled.
- (2) Time intervals: We assume that parameters are constant throughout calendar months but may change from one month to next. Lags are measured in days.
- (3) Exposures: The user may specify a measure of exposure for each month. By default, the system assumes constant unit exposure. The purpose of the exposure parameter is to allow the user to account for a principal source of variation in monthly frequencies.
- (4) Events: Each claim may be described by the dates and amounts of the events it triggers: the accident date, the report date and an initial case reserve, zero or more subsequent valuation dates and case reserves changes, zero or one payment date and amount, and zero or one recovery date and amount.
- (5) Distributions: For most variables, the user may specify a distribution and associated parameters.
- (6) Frequency: The user may specify monthly claim frequency as a Poisson distribution with mean proportional to earned exposure, or as a Negative Binomial distribution with mean proportional to earned exposure and variance proportional to the mean frequency (which implies that the variance is also proportional to the earned exposure..)

Modeling Loss Emergence and Settlement Processes

- (7) Report lag: The lag between occurrence and reporting is assumed to be distributed Exponential, Lognormal, Weibull, or Multinomial. The Multinomial distribution allows the user to define proportions of claims reporting within one month, two months, and so on.
- (8) The lags between reporting and payment, between one valuation date and the next, and between payment *and recovery or adjustment*, are also assumed to be distributed Exponential, Lognormal, Weibull, or Multinomial.
- (9) Size of loss: The actual size of the loss to the insured, independent of responsibility for payment, is distributed Lognormal, Pareto, or Weibull.
- (10) Case reserve factor: Case reserves are assumed to equal the actual size of loss, adjusted for the minimum, the maximum, the deductible, and the probability of closure without payment, all multiplied by an adequacy factor. This factor is assumed to be distributed Lognormal. The user may specify the mean factor at four points in time between the report and payment dates.
- (11) Fast-track reserve: A value may be assigned to each loss at first valuation, independent of regular case reserves and case reserve factor.
- (12) Initial payment factor: The initial payment of each loss not closed without payment is assumed to equal the actual size of loss, adjusted for the minimum, the maximum, the deductible, multiplied by a payment adequacy factor (PAF). The PAF determines the size of any subsequent adjustment or recovery.
- (13) Second-level distributions: The LSMWP models the drift in parameter values that may take place for many reasons but chiefly because of business turnover. It has developed an autoregressive model to reflect parameter drift.
- (14) Monthly vectors of parameters: For nearly all distributional parameters, the user may specify a single value or a vector of values.
- (15) Frequency Trend and Seasonality: The user may specify monthly trend and seasonality factors for frequency that are applied to means.
- (16) Severity Trend: The user may specify monthly trend factors for severity.
 - The “main” trend is allowed to operate up to the accident date and a fraction of this trend, defined by Butsic’s “alpha” parameter, is allowed to operate between accident and payment dates.
 - Case reserves before the adequacy factor are centered around the severity trended to the payment date.
- (17) Lines and Loss Types: The prototype model recognizes that loss data often involves a mixture of coverages and/or loss types with quite different frequencies, lags, and severities. Therefore, it allows the user to specify a two-level nested hierarchy of simulation specifications, with one or more “Lines” each containing one or more “Types”.
- A typical Line might be “Auto,” typical Types within that Line might be “APD,” “AL-BI,” and “AL-PD.”
- This hierarchy allows the user to set up any reasonable one or two level classification scheme.

- Accident frequencies are modeled at the Line level and loss counts per accident are distributed among Types using a discrete distribution.
- (18) Lines and Loss Types Example: An Automobile occurrence might give rise to a single Auto Physical Damage (APD) claim with probability 0.4, to a single Auto Property Damage Liability (AL-PD) claim with probability 0.2, to a single APD and a single AL-PD claim with probability 0.2, to a single Auto Bodily Injury Liability (AL-BI) claim with probability 0.1, to two AL-BI claims with probability 0.05, etc.
- (19) Correlations: The prototype model makes it possible to request correlated samples of certain variables without fully specifying their joint distribution. These variables are (a) the mean frequencies across Lines and (b) the size of loss and report lag within a Type. To specify correlated frequencies among lines, the user specifies the marginal frequency distribution for each Line and then specifies a correlation matrix and copula to determine the joint distribution.
- (20) Clustering: The prototype simulator allows a selectable fraction of loss sizes and a selectable fraction of case reserves to be rounded to two significant digits, imitating clustering around round numbers frequently observed.
- (21) Output: The prototype simulator produces output as comma-delimited (“csv”) text files or by launching an instance of Excel and populating it with worksheets. In both cases, the possible output tables include claim and transaction files (together displaying the complete loss history), all the usual triangles, a table of large losses, a summary of the simulation specifications, and a summary of the frequency derivation by month.

5. DOCUMENTATION OF OPEN SOURCE MODEL

The help files in the new open source model fully document the implementation of all of the features that were included in the prototype model. These help files are organized into comprehensive model documentation on the CAS web site at <http://www.casact.org/research/lsmwp/losshelp/index.cfm?fa=main>. The documentation within the model and on the CAS web site will be kept up to date as this open source software is enhanced.

Comprehensive program instructions are provided in Appendix A, and are also included within the model and in the CAS web site model documentation.

6. TESTING DETAILED OUTPUT AND FITTING THE MODEL

This section describes tests performed on the Simulator’s output. The basic procedure for these tests was to run the Simulator with parameters that generate random claims with their transactions, from distributions whose densities and/or distribution functions are easily computed. Then we test the output files to see how closely the empirical distributions match the theoretical distributions.

Modeling Loss Emergence and Settlement Processes

An important point is that properly designed and documented statistical tests on simulator output also enable the user to fit models to real data. The user would transform the real data into the same format as the “claims” and “transactions” files output from the simulator. Then the user can run the same tests that we have developed for testing simulator output. Many of these tests include maximum likelihood estimates of the parameters. The user can generate the same maximum likelihood estimates on his/her own data. To enable this, we are including much of the source code used to perform the tests.

This section deals with the claims at their ultimate values, not with age-to-age loss development. Chapter 3 discusses the aggregate reserve triangles generated by summarizing the reserve transactions as of specific accident periods and development dates.

6.1 Data Formats

The simulator generates a Claims file and a Transaction file. The Claims file contains one record per claim. A sample of the format follows:

Simulation No	Occurrence No	Claim No	Accident Date	Report Date	Line	Type
1	1	1	20000104	20000227	1	1
1	2	1	20000105	20000818	1	1
.....						
1000	324	1	20001222	20010120	3	3
1000	325	1	20001215	20010224	3	3

The Transaction file may have multiple transactions per claim, and has format:

Simulation No	Occurrence No	Claim No	Date	Transaction	Case Reserve	Payment
1	1	1	20000227	REP	2000	0
1	1	1	20000413	RES	89412	0
1	1	1	20000417	CLS	-91412	141531
1	2	1	20000818	REP	2000	0
.....			

In this case, the first three records in the Transaction file correspond to the first record in the Claims file. The amounts in the “Case Reserve” column are incremental reserve changes, while the “Payment” column shows the payment made on each date. Thus, on 4/17/2000, Occurrence 1 –

Modeling Loss Emergence and Settlement Processes

Claim No 1 closed, had zero case reserve (i.e., the sum of the Case Reserve column) and a cumulative paid amount of \$141,531.

For purposes of testing ultimate severity and number of claims, the two files were merged into an “Ultimate Loss” file whose fields are listed:

Simulation.	Occurrence.	Claim.	Accident.	Report.	Line	Type	Case.	Payment
No	No	No	Date	Date			Reserve	

This file has the same number of records as the Claims file. For each Simulation / Occurrence / Claim combination, the Case Reserve is zero, since all records are at ultimate value, and the Payment is the sum of all the Payment amounts from the Transaction file.

Many of the tests were performed on this “Ultimate Loss” file. **If the Modeler can format his/her own data into this format, then he/she can use the LSMWP source code to parameterize the Simulator.**

6.2 Description of Testing

The LSMWP Testing Group ran tests on the model as it was being developed. One reason for doing this is that these tests helped debug the model’s calculations. The model simulates claims and transactions using the “R” language. Part of the “debugging” consisted of making sure the R commands work the way the LSMWP members think they work.. The subsections of section 6.2 each describe specific groups of tests that were performed.

6.2.1 Test of Elementary Frequencies, Trend, and Zero-modification

The model allows the user to input annual frequency by line and year, to specify the probability distribution for the types of claims within a line, and to specify “P(0)”, the probability that a claim closes without payment. This section describes the final test among a series of tests of these parameters. The size of loss distribution was not used in this test except that a status indicator was generated for each claim, with “CNP” (resp. “CWP”) meaning that the claim was closed without (resp. with) payment.

A partial description of this test follows. The full list of parameters is in the Appendix B.

Test was run 10/27/2009. Project name: Frequency Test

Modeling Loss Emergence and Settlement Processes

Purpose: Test frequency with trend. Two types within one line.

- One Line with annual frequency distributed Poisson with mean 120
- Set claim/acc distribution matrix as follows:
 - Prob =75% that one Type 1 claim is generated.
 - Prob =25% that one Type 2 claim is generated.
- Freq Trend: 1.02 constant throughout
- $P(0) = 0.4$, $EstP(0) = 0.4$ for each Type.
- Accident Years: 2000-2002
- Random Seed: 16807
- Frequency correlation copula: normal $Correlation=c()$ Dim = 1

Run: 1,000 simulations.

What would we expect from this run? If we think of the months k in the three accident years as numbered 1 through 36, the Poisson λ_k for each month² fits the formula:

$$\lambda_k = 10 (1.02)^{k/12}, \text{ where } 10 = 120/12.$$

The total number of expected claims in each simulation is $\lambda = 10 \sum_{k=1}^{36} (1.02)^{k/12} = 371.2144835$.

With 1,000 simulations, the total expected claim count equals $EN = 371,214$. The actual number generated was $N = 371,198$. The closeness of the actual and expected counts indicates that the Poisson frequency and trend parameters are operating in the manner expected.

The parameters were chosen so that a claim has probability of 0.75 (0.25) of being Type 1 (Type 2) and probability 0.60 (0.40) of having status CWP (CNP). A very important property of Poisson distributions states that the number of claims in 1,000 simulations for each Type(x) x Status(y) category are mutually independent Poisson random variables with mean $EN * Pr[Type=x]*Pr[Status=y]$.³

Chi-Square test for independence of Type and Status on claim counts.

To use the Chi-square test, we must we must construct data with both expected and actual claim counts. From the input parameters, there is only one Line, with a claim for this line being Type 1 (Type 2) 75% (25%) of the time. Also, the probability of the claim being Status CWP (CNP) is 60% (40%). Generally we wish to compare the actual and expected claims by unique “cells.” Here the

² The simulator converts annual frequency to monthly frequency when performing the simulations.

³ [2], pp. 103-104

Modeling Loss Emergence and Settlement Processes

cells are defined by unique combination of all the predictors. Here are some S-PLUS⁴ statements that produce all possible combinations:

```
temp1a <- unique(dataacc$Simulation.No)
temp1b <- unique(dataacc$Accident.Year)
temp1b2 <- unique(dataacc$Accident.Month)
temp1c <- unique(dataacc$Line)
temp1d <- unique(dataacc$Type)
temp1e <- unique(dataacc$Status)
temp1 <- list(
  Simulation.No=temp1a, Accident.Year =temp1b, Accident.Month=temp1b2,
  Line=temp1c, Type=temp1d, Status=temp1e)
temp2 <- expand.grid(temp1)
  for (j in 1:length(temp2)) {
    if(is.factor(temp2[,j]))
      temp2[,j] <- as.character( temp2[,j])
  }
temp7 <- order(temp2$Simulation.No,temp2$Accident.Year,temp2$Accident.Month,
  temp2$Line,temp2$Type,temp2$Status)
temp3 <- temp2[temp7,]
temp3$timeyear <-
  timeyear(base.year,temp3$Accident.Year,temp3$Accident.Month)
```

The “expand.grid” statement creates a data frame⁵ with all possible combinations of the individual variables. It is important to note that the expand.grid statement produces more unique combinations of the individual variables than are found in the actual data. For this run, there are 144,000 combinations (1000 x 3 x 12 x 1 x 2 x 2). The last statement calls a user-defined function “timeyear”, which converts accident year and month to a numeric variable that corresponds to the length of the trend period over which the simulator’s “trend” parameter operates.

To apply Chi-square testing, we want to calculate the expected number of claims for each record in the data frame and to tabulate the actual number of claims. The total claims for each simulation and accident year/month equals $10 * (1.02)^{\text{timeyear}}$. Then the expected count for each type and status combination is just this number multiplied by the probabilities mentioned earlier. After this “expected count” dataframe is built, it can be merged with the actual counts to produce a dataframe, which we will call “ccstacked”, with both actual and expected counts. Here is the image of the first few records:

Data frame “ccstacked”

⁴ SPlus and R are implementations of the statistical language S.

⁵ Think of a data frame in R or SPlus as a matrix, with each column representing either a predictor or a value and each row representing a combination of predictors or an observation.

Modeling Loss Emergence and Settlement Processes

Simulation. No	Accident Year	Accident Month	Line	Type	Status	timeyear	expec.count	count
1	2000	1	1	1	CNP	0.08333 3	3.004955	2
1	2000	1	1	1	CWP	0.08333 3	4.507432	6
1	2000	1	1	2	CNP	0.08333 3	1.001652	0
1	2000	1	1	2	CWP	0.08333 3	1.502477	3
1	2000	2	1	1	CNP	0.16666 7	3.009918	2
1	2000	2	1	1	CWP	0.16666 7	4.514877	4
1	2000	2	1	2	CNP	0.16666 7	1.003306	1

Proper application of the Chi-Square test requires that the expected counts from the theoretical distribution be modified so that the total number of claims and the marginal probabilities match the actual data.

	Actual Counts				Expected Counts		
	Type 1	Type 2	Margin		Type 1	Type 2	Margin
CNP	111,066	37,007	0.39890 6	CNP	111,029.0	37,044.0	0.398906
CWP	167,268	55,857	0.60109 4	CWP	167,305.0	55,820.0	0.601094
Margin	0.74982 6	0.25017 4	371,198		0.749826	0.25017 4	371,198

$$\chi^2 = \sum_i \sum_j \frac{(Actual_{ij} - Expected_{ij})^2}{Expected} = 0.0819.$$

$$\Pr [\chi^2 > 0.0819] = 0.775.$$

The independence of Type and Status on counts is supported.

If we instead use the crosstabs command in S-PLUS, we obtain the output:

```
> temp1 <- crosstabs( count ~ Status + Type ,
+ data=datacc.stack,
```


Modeling Loss Emergence and Settlement Processes

```

+           na.action=na.fail,                ##### No missing data should
exist
+           drop.unused.levels=F
+       )
> temp1
Call:
crosstabs(formula = count ~ Status + Type, data = dataacc.stack, na.action =
           na.fail, drop.unused.levels = FALSE)
371198 cases in table
+-----+
|N      |
|N/RowTotal|
|N/ColTotal|
|N/Total |
+-----+
Status  |Type          |2          |RowTotal |
+-----+-----+-----+
CNP     |1.11066e5    |3.7007e4   |148073   |
        |0.75         |0.25       |0.4      |
        |0.4          |0.4        |         |
        |0.3          |0.1        |         |
+-----+-----+-----+
CWP     |1.67268e5    |5.5857e4   |223125   |
        |0.75         |0.25       |0.6      |
        |0.6          |0.6        |         |
        |0.45         |0.15       |         |
+-----+-----+-----+
ColTotal|278334       |92864      |371198   |
        |0.75         |0.25       |         |
+-----+-----+-----+
Test for independence of all factors
Chi^2 = 0.081902898185 d.f.= 1 (p=0.77473503489)
Yates' correction was not used

```

Thus, we can see that the “crosstabs” command produces the Chi-square analysis.

Use of Poisson Generalized Linear Model (GLM) for testing.

The hypotheses underlying the model parameterization imply that covariates such as Type and Status have multiplicative effects on the number of claims. Such situations are conveniently modeled using a log-linear model such as a Poisson GLM with categorical and/or numeric predictor variables. This model is more flexible than just relying on cross-tabulations. Since this is a powerful predictive modeling technique, I will describe its use here in some detail. In particular, we will discuss how the Type and Status variables affect the claim counts and whether there is an interactive effect.

In the Poisson GLM, the dependence of $\mu_i = E(Y_i)$, the variable of interest, is related to the covariates through the relationship: $\log \mu_i = \sum_{j=1}^p \beta_j x_{ij}$ for $i=1, 2, \dots, n$. Here p is the number of

Modeling Loss Emergence and Settlement Processes

covariates, n is the number of observations for which the fitting takes place, and β is the vector of coefficients determined by the model.⁶

To model the effects of Type and Status, we can use the data frame “ccstacked” described earlier. We run a “full” model, using Type, Status, and their interaction as predictors, and a “reduced model” without the interactive effect. If our model parameters work the way we think they should, then the interactive term should add very little predictive power and the reduced model should fit very well. The details of the actual runs are contained in Appendix B. We outline the results here.

The S-PLUS statement for the reduced model is

```
model15x<- glm(count ~ + Type + Status,  
               data = temp.datacc.stack,  
               family = poisson,  
               x=T)  
summary(model15x,correlation=F)
```

The predicted value for this fitted model is

$$\log \hat{\mu} = 1.126272675 - 1.097685774 * Type + 0.410026757 * Status \quad (6.2.1)$$

How do we interpret these equations? Type and Status in the “ccstacked” file each take two possible values. When converted to the indicator variables, the translations are that $Type = 0$ (1) when the original variable $Type = 1$ (2), and that $Status = 0$ (1) when the original variable $Status =$ CNP (CWP). There are 36,000 combinations of simulation number, accident year, and accident month in the data.

Equation (6.2.1) gives the log of the predicted number of claims by Type and Status for each of these combinations. The sum of the predicted values for the four combinations of Type and Status equals 10.31. Multiplying this by 36,000 produces the total number of claims 371,198. This is almost exactly equal to the a priori expected number of claims.

The coefficient of -1.097685774 for Type implies that

$$\frac{\Pr[Type = 2]}{\Pr[Type = 1]} = \exp(-1.097685774) \approx 0.3336.$$
 Indeed, this ratio is 1/3 from the input parameters.

The coefficient of 0.410026757 for Status implies that

⁶ [22] is an excellent reference on GLMs and Poisson GLMs.

$$\frac{\Pr[\text{Status}=CWP]}{\Pr[\text{Status}=CNP]} = \exp(0.410026757) \approx 1.507. \text{ The ratio of the input values is 1.500.}$$

Thus, we have good idea that the model fits the data well.

Can the model be improved by including the interactive variable *Type*Status* ?

To check this, we run the model statement

```
model6x<- glm(count ~ Type + Status + Type*Status
              data = temp.dataacc.stack,
              family = poisson,
              x=T)
```

To test the additional value of the interactive variable we compare the deviations:

```
> anova(model5x,model6x,test="Chi")
Analysis of Deviance Table
```

Response: count

	Terms	Resid. Df	Resid. Dev	Test Df
1	+ Type + Status	143997	160969.366	
2	Type + Status + Type * Status	143996	160969.284	+Type:Status 1

	Deviance	Pr(Chi)
1		
2	0.0819088429	0.774727081

The “Resid.dev” is the residual deviance as defined for GLMs and will not be discussed further here. What is interesting is the reduction in deviance from introducing the interactive variable and the corresponding probability exactly match those from the earlier “crosstabs” calculation. Thus we cannot reject the null hypothesis that there is no interactive term.

In the actual frequency testing, we ran GLMs that include the trend period as a numeric predictor, since we are assuming 2% annual frequency trend. These models provided better predictions, as we would expect. Even in the presence of a numeric time variable, the *type*status* interaction is not indicated.

6.2.2 Testing Size of Loss Distributions

The simulator was run May 12, 2010 to produce output on which one can test the distribution of claim severities (i.e., size of loss). The simulator produced claims from three different lines, with the severity distributions by line set as lognormal, Pareto, and Weibull. The frequency parameters assure that the number and size of claims by lines are mutually independent.

Description of Run and Parameters:

Set up a test for univariate severity distributions and some time distributions.

Project name: Test severities 20100512

Purpose: Test univariate ultimate severities. Set up three lines of business with no correlation in frequency among the three lines. For each simulation, the number of occurrences by line by accident year is as described in the frequency parameters below.

- 3 Lines. For each line, every occurrence generates one claim of one type.
- Zero trend for frequency, zero trend for severity.

Accident Years: 2000-2001

Random Seed: 16807

Frequency correlation -- uncorrelated

Each line has same frequency parameters

Annual Frequency: Poisson(**600**)

Monthly exposure: (1), zero trend, no seasonality

Set claim/acc distribution matrix as follows:

each occurrence generates 1 claim of one type

P(0) = 0.0, EstP(0) = 0.0

Lags: Irrelevant for this run except for report lag. We are not testing the reserve change process.

Report lag: Exponential with Rate = 1/365, mean=365 days. Max=3650.

Payment lag: Maximum one day

Inter-valuation lag: Maximum one day

Correlation of Amount with lag: normal Correlation=c() Dim=2

Reserve adequacy: Irrelevant, leave at default values.

Run: 100 simulations.

Severity parameters – these vary by Line.

- Line 1 Type 1: Lognormal, mean=100,000, stnd. dev. 100,000, max 10,000,000. This means that the input lognormal parameters are $\mu = \text{meanlog} = 11.16636357$, $\sigma = \text{sdmean} = 0.832549779$.
- Line 2, Type 2: Pareto with α (shape) =6, θ (scale) = 500,000. This results in a mean of 100,000 and standard deviation 122,474.5.

Modeling Loss Emergence and Settlement Processes

- Line 3, Type 3: Weibull with θ (scale) = 95,000, τ (shape) = 0.9. This results in mean of 99,957 and standard deviation of 111,256.

Expected total # claims = 600 (freq) x 100 (# sims) x 3 (lines) x 2 (years) = 360,000.

Actual # claims: 359,819.

Results of simulation:

The output of the simulation was summarized to the “Ultimate Loss” file format described above in section 6.1.⁷ This file contains one record for each claim, with the ultimate value of the claim in the “payment” field and the case reserve equal to zero, since all claims are at ultimate value.

The work in the test was divided up between two people. Joe Marker ran the simulator, created the summary file, and wrote the main body of this section. Yuting Yang did the testing work described below, wrote the R code, and produced the graphs.

The R project included four libraries that can be added using the statements.

```
library(stats4)
library(MASS)
library(actuar)
library(graphics)
```

Three separate vectors, named `ultloss1`, `ultloss2`, and `ultloss3`, were created using R. For $k = 1$ to 3, `ultloss k` is the vector of loss sizes for the claims from line k . To test whether the simulator actually generates claims sizes according to the parameters input, we fit each of the vectors to the lognormal, Pareto, and Weibull distributions, and then conducted goodness of fit tests.

First, exploratory data analysis was performed on the three data sets. This included:

- Calculating the standard statistics
- Generating histograms, empirical densities, empirical log densities and empirical cdfs.

Appendix B contains all the graphs.

The following illustrates how the R commands produced the graphs for Line 1.

```
hist(ultloss1,main="Histogram of observed data of Line 1",
     freq=FALSE,breaks=10000,xlim=c(0,1050000))
plot(density(ultloss1),main="Density estimate of Line 1",xlim=c(-1000,600000))
plot(density(log(ultloss1)),main="LogDensity estimate of Line 1",xlim=c(0,20))
plot(ecdf(ultloss1),main="Empirical cdf of Line 1",xlim=c(0,1e+06))
```

⁷ The name of this file is “ultloss20100520.csv”.

Modeling Loss Emergence and Settlement Processes

This analysis of the graphs suggests that Line 1 has a lognormal density. Lines 2 and 3 have similar shapes, but Line 2 has more data in the tail, which suggests a heavier-tailed distribution.

The initial guess for the severity distributions is:

Line 1 – Lognormal, Line 2 – Pareto, Line 3 – Weibull.

The remaining steps in fitting the data consisted of:

1. Calculating the maximum likelihood estimates for the three distributions.
2. Informally looking at the fits by constructing and examining the Q-Q Plots for the candidate distributions.
3. Producing a more formal goodness of fit test by binning both the expected and empirical distributions and then performing chi-square tests.

Calculate the maximum likelihood estimates of the parameters.

The next step was to use maximum likelihood estimation (“m.l.e.”) to determine the optimal parameters for each line. We used the R command “fitdistr” for the calculations. However, prior to using this command, we calculated the negative loglikelihoods (“n.l.l.”) using the density functions. For each observed value x , let’s review the both the p.d.f. $f(x)$ and the corresponding negative loglikelihood $-\ln f(x)$.⁸

a) If X has the lognormal distribution with parameters μ and σ , then $Y = \ln X$ has the normal distribution with parameters μ and σ . For the m.l.e., it is more convenient to fit observations $y = \ln x$ than it is to fit x . We have

$$f_Y(y) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{z^2}{2}}, \text{ where } z = \frac{y - \mu}{\sigma}, \text{ and n.l.l.} = \ln(\sigma\sqrt{2\pi}) + \frac{z^2}{2}.$$

b) For X distributed Pareto⁹ with shape α and scale θ , we have

$$f(x) = \alpha \theta^\alpha (x + \theta)^{-(\alpha+1)}, \text{ and n.l.l.} = -(\ln \alpha + \alpha \ln \theta) + (\alpha+1) \ln (x+\theta)$$

c) For X distributed Weibull with shape τ and scale θ , we have

⁸ The parameterizations here follow those from Appendix A of [2].

⁹ Pareto refers to Pareto Type II, Lomax

Modeling Loss Emergence and Settlement Processes

$$f(x) = \frac{\tau x^{\tau-1}}{\theta^\tau} \exp\left[-(x/\theta)^\tau\right], \text{ and n.l.l.} = -\ln \tau + \tau \ln \theta - (\tau-1) \ln x + (x/\theta)^\tau$$

The parameter values are then determined by finding the parameter values that maximize the total n.l.l.

For Line 1, the R code below illustrates the use of the “fitdistr” command to fit $\ln X$ to a normal distribution::

```
fit1.ln <- fitdistr(log(ultloss1),"normal")
fit1.ln$estimate # mean sd
#11.1659376 0.8361509
-fit1.ln$loglik #148761.9
```

The output of any R command is an “object” with “properties”. The first line above gives this object the name “fit1.ln”. The vector fit1.ln\$estimate of length 2 contains the optimal parameter values $\mu = 11.1659376$ and $\sigma = 0.8361509$. The number -fit1.ln\$loglik is the minimized negative loglikelihood.

The same command was applied to the Line 2 data for both the Pareto and Weibull distributions.

```
## 2.1-Pareto ##
fit.p2<-fitdistr(ultloss2.0,dpareto,list(shape=6,scale=500000))
## list() provides initial values for optimization
fit.p2$estimate # shape scale
#5.97635e+00 5.00000e+05
-fit.p2$loglik #1500363
## 2.2-weibull (second method slightly better) ##
fit2.w <- fitdistr(ultloss2.0,"weibull")
fit2.w$estimate # shape scale
# 9.056193e-01 9.750673e+04
fit2.w$loglik # -1500950
fit.w2<- fitdistr(ultloss2.0,dweibull,list(shape=.9097626,scale=95000))
fit.w2$estimate # shape scale
# 9.009281e-01 9.500000e+04
-fit.w2$loglik # 1500926
```

The Line 3 results are shown:

```
## 3.1-Pareto ##
fit.p3<-fitdistr(ultloss3.0,dpareto,list(shape=7,scale=6.026793e+05))
```

Modeling Loss Emergence and Settlement Processes

```
fit.p3$estimate #   shape      scale
                #6.966806e+00 6.026793e+05
-fit.p3$loglik  #1499343
## 3.2-weibull (first method slightly better) ##
fit.w3 <- fitdistr(ultloss3.0,"weibull")
fit.w3$estimate #   shape      scale
                # 9.052532e-01 9.907429e+04
-fit.w3$loglik  # 1498920
```

For a set of observations, it is tempting to compare the n.l.l. from different models to label one of them as best. There is some logic in doing this for the Pareto and Weibull distributions, since they each have two parameters. However, there are better methods to differentiate them, which we discuss below.

It would be wrong to compare an n.l.l. from fitting the lognormal to either the n.l.l. for the Weibull or Pareto. This is because the observed variable $\ln X$ for fitting the lognormal is on a totally different scale than the variable X used for the Weibull and Pareto.

QQ Plots.

The QQ plot provides one of the best ways to visually compare two distributions. If one of the distributions is the empirical distribution of the observed values and the other is a c.d.f. of a random sample of the same size from the fitted distribution, then the closer the plot is to a 45 degree line, the better the fit.

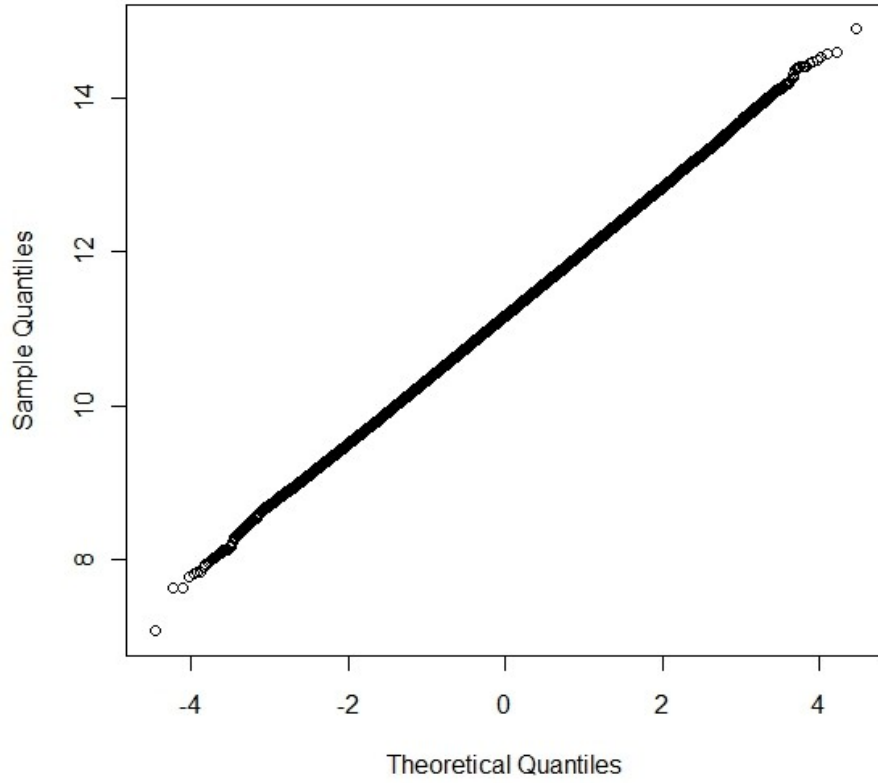
The following R statements create a QQ-plot for the Line 2 data and the fitted Weibull distribution.¹⁰

```
thqua.w2 <- rweibull(n2,shape=fit.w2$estimate[1],scale=fit.w2$estimate[2])
qqplot(ultloss2,thqua.w2,xlab="Sample Quantiles", ylab="Theoretical Quantiles",
main="Line 2, weibull")
abline(0,1,col="red")
```

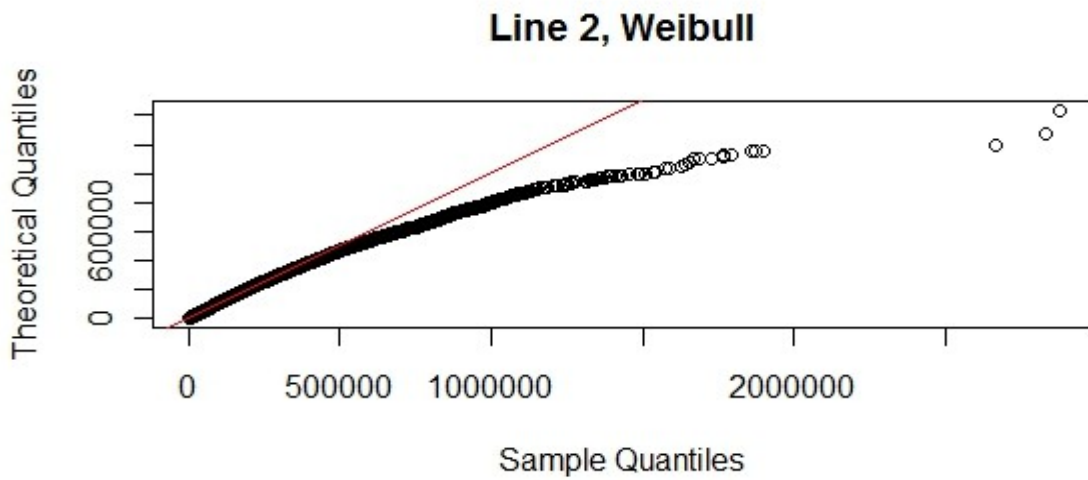
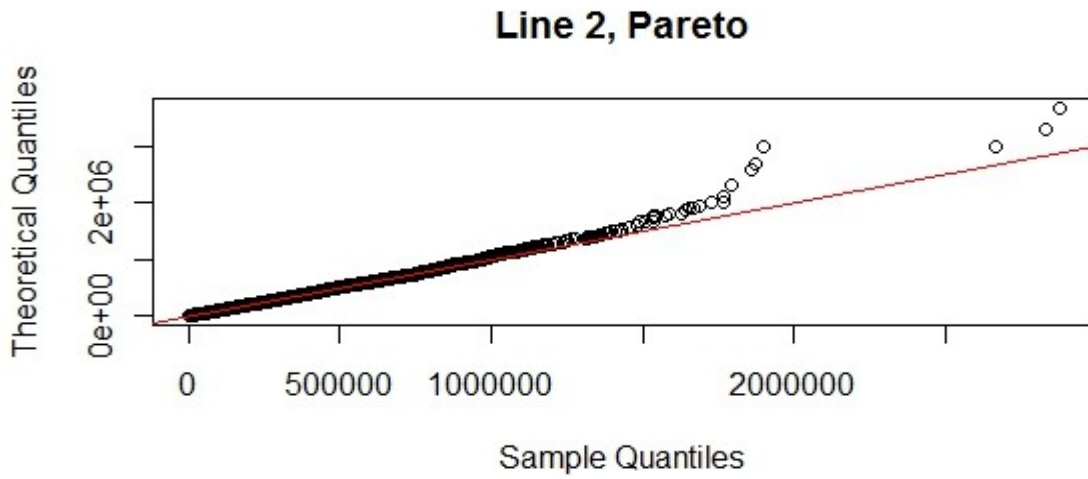
The QQ-plots for the fitted distributions for each line follow.

¹⁰ These plots show quantiles of the observed values versus quantiles of a random sample from the fitted distribution. More commonly, the corresponding quantiles of the fitted distribution are used rather than a sample. See Section 5.1 of [31], Venables, W.N. and Ripley, B.D., *Modern Applied Statistics with S, Fourth Edition*, Springer-Verlag (2002).

Line 1, Lognormal

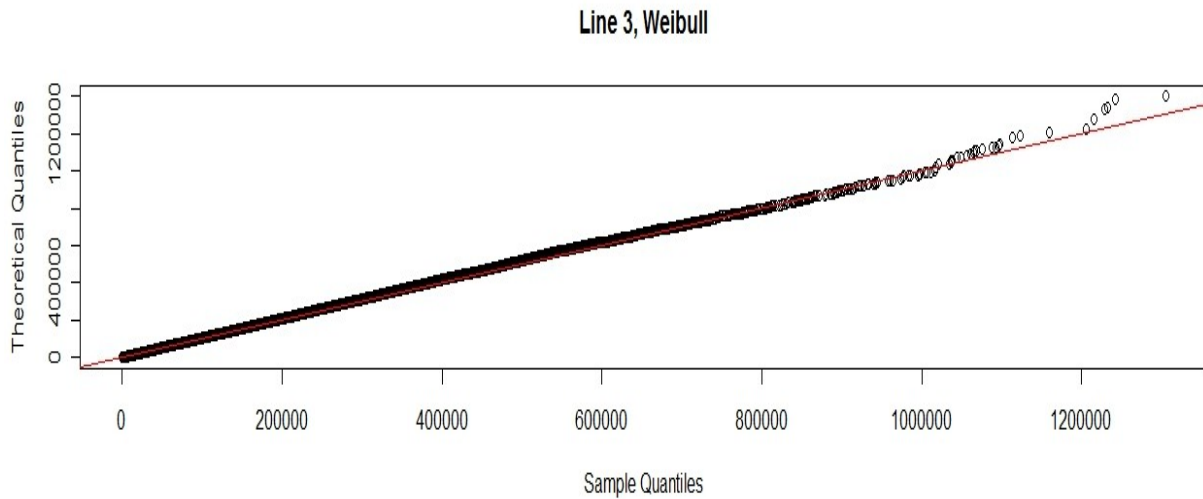
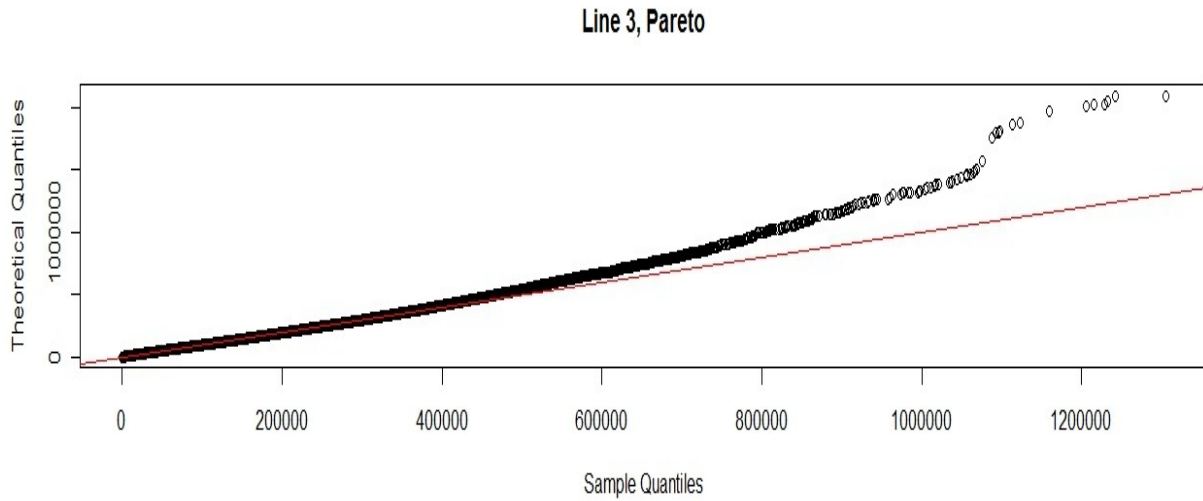


For Line 2, the Pareto fits the data much better over a longer range than the Weibull distribution.



Modeling Loss Emergence and Settlement Processes

For Line 3, the Weibull distribution is best.



The good fits for each of the three Lines indicates that the chosen models fit the ultimate loss size well. We can go further to test more formally the hypotheses that these models fit the data.

Goodness of Fit tests

We discuss the use of the Chi-square test for the Pareto fit to Line 2. First, define the break points and set up the bins according to these break points.

```
#2.3.1 Pareto Chi-Square #
m = mean(ultloss2)
s = sqrt(var(ultloss2))

ult2.cut <- cut(ultloss2.0,
               breaks = c(0,m-s/2,m,m+s/4,m+s/2,m+s,m+2*s,2*max(ultloss2))) ##binning data
table.ult2 <- table(ult2.cut) ## binned data table
ult2.os <- c(as.vector(table.ult2)) ## vectorization

b = length(ult2.os)

labs.2 <- levels(ult2.cut) ## extract the breakpoints
break.2 <- cbind(lower = as.numeric(sub("\\((.+),.*", "\\1", labs.2)),
                upper = as.numeric(sub("[^,]*,([^)]*)\\]", "\\1", labs.2)))
as.numeric(sub("[^,]*,([^)]*)\\]", "\\1", labs.2))
```

See the R code in Appendix B for an illustration of the previous two R statements.

To calculate the expected number of claims in each bin, use `ppareto` to generate the c.d.f. values at the break points. To calculate the expected number of claims in each bin, use `ppareto` to generate the c.d.f. values at the break points. Note that we need to define an “excess” interval to contain the expected number of claims larger than the last break point. Here we have chosen the last break point to be twice the maximum observed value.

```
ult2.p <- ppareto(break.2,shape=fit.p2$estimate[1],scale=fit.p2$estimate[2]) ##
Pareto cdf values at break points

ult2.prob <- (ult2.p[,2]-ult2.p[,1])[1:b-1] ## Probabilities of each interval
ult2.ex <- n2.0*c(ult2.prob,1-sum(ult2.prob))
## Expected frequency of each interval and the "excess" interval
```

The expected # claims by interval are in `ult2.ex` and the observed # claims are in `ult.os`.

```
E.2 = ult2.ex
O.2 = ult2.os

x.sq.2 = (E.2-O.2)^2/E.2

cbind(E.2,O.2,x.sq.2) ## expected, observed, and chi-square of each interval
after full adjustment

##### Test Statistic Calculation #####
#
#-----
# E.2 O.2 x.sq.2
#[1,] 43993.890 44087 0.19705959
#[2,] 35651.989 35680 0.02200752
#[3,] 10493.758 10323 2.77864169
#[4,] 7240.583 7269 0.11152721
#[5,] 9277.383 9164 1.38570182
#[6,] 8063.576 8176 1.56743997
#[7,] 5289.820 5312 0.09299630
```

Modeling Loss Emergence and Settlement Processes

Next we calculate chi-squared statistic and its critical value, which shows that the null hypothesis that the Pareto model fits the data cannot be rejected. **It is important to note that the degrees of freedom are not six as we might expect, but rather $6-2 = 4$.** This is because the Pareto distribution used for the expected values was the best fit Pareto distribution from the data, which is a two-parameter distribution.¹¹

```
##chi-square test statistic##
df=length(E.2)-1-2          ## df = 4
chi.sq.2 <- sum(x.sq.2)     ## test statistic
chi.sq.2                    ## 6.155374
qchisq(.95,df)              ## critical value ## 9.487729
1-pchisq(chi.sq.2,df)      ## p-value ## 0.1878414
```

Using the chi-squared test in R directly would produce a wrong p -value, to wit:

```
## chi-square goodness-of-fit test from R ##
chisq.test(0.2,p=E.2/n2.0)
#####
# Chi-squared test for given probabilities
# data: 0.2
# X-squared = 6.1554, df = 6, p-value = 0.406
#####
```

Summary.

The initial motivation for the procedures in this section was to test the output of the simulator to see whether it fit the severity distribution specified by the input parameters. We first used graphs to select candidates for fitted distributions. For example, for Line 2, this led to selecting Weibull and Pareto as candidates. Recall that the size of loss distribution for the Line 2 simulation was set as Pareto with with α (shape) =6, and θ (scale) = 500,000. The m.l.e. estimates for the Pareto parameters are $\hat{\alpha} = 5.97635$ and $\hat{\theta} = 500,000$.

The Q-Q plots of the observed values versus sample values from the fitted theoretical distribution gives a visual indication of which distribution best fits the observed values for each line. Appendix B contains all the graphs. Next we used the chi-square test to decide whether the selected

¹¹ See, for example, [32] Devore, J. L., *Probability and Statistics for Engineering and the Sciences*, Fifth Edition (2000), Duxbury Thomson Learning.

models fit the data well. We did not calculate other appropriate tests such as Kolmogorov-Smirnov or Anderson-Darling.¹²

One can also use the procedures in this section to fit a set of actual claim size data, once it is converted to a format similar to the “ultloss” file from this chapter. To help with this, we included in Appendix B all the R statements used to fit and test the various distributions we tested for each line.

6.2.3 Testing Correlated Frequencies

In order to make sure the Copula feature in the model is appropriate, we implemented statistical techniques to test the Gaussian Copula setting based on correlated set of frequency data simulated by the model. The simulation run results files “c_20100720_1900.csv” and “bymonth.csv” are used, which have 1 year frequency data for 3 lines. The number of simulation is 1,000. The correlation assumption is as below:

Correlation	Line 1	Line 2	Line 3
Line 1	1	0	0.99
Line 2	0	1	-0.01
Line 3	0.99	-0.01	1

Some key parameters for the simulator run that produced the dataset for this test include:

Annual frequency for each line is Poisson with mean 96 occurrences.

Each occurrence generates exactly one claim.

The simulator produces claims for accident year 2000 only.

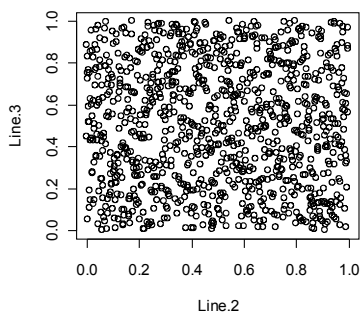
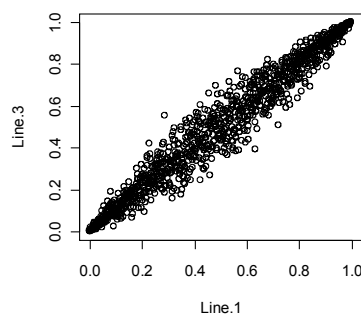
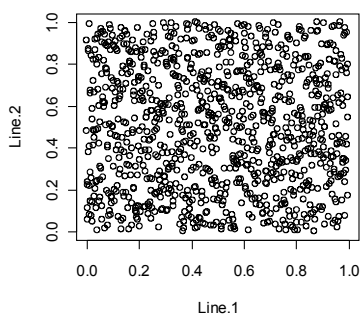
Descriptive statistics

The annual frequency number for each line and each simulation is extracted based on the run results file. In addition, the 1st month frequency number is also extracted for a parallel test run. The two sets do not show too much difference regarding testing results. The test results below are based on annual frequency only. See Appendix B for the layout of the data set used for this test.

¹² See [2], pp. 448-458, for a fuller discussion of these tests.

Modeling Loss Emergence and Settlement Processes

A scatter-plot of the empirical cumulative distribution function for each pair of lines is displayed below, showing a high positive correlation between line 1 and line 3. This is in line with our assumption for the correlation between Line 1 and Line 3: 0.99. Other pairs' correlation cannot be identified, which is also consistent with our parameter assumption.



Test of Correlation

- (1) Fit simulated data set (data pair) to normal copulas.

Two methods have been used for the copula fitting. We can see that Rho 1, Rho 2 and Rho 3 MLE are very close to our assumption (0, 0.99, -0.01).

Using maximum likelihood method

The estimation is based on the maximum likelihood

and a sample of size 1000.

	Estimate	Std. Error	z value	Pr(> z)
rho.1	-0.002112605	0.031977597	-0.06606516	0.9473259
rho.2	0.979258746	0.000921392	1062.80366235	0.0000000
rho.3	-0.010486832	0.031974114	-0.32797880	0.7429277

The maximized loglikelihood is 1591.565

The convergence code is 0

Using Inversion of Kendall's tau

The estimation is based on the inversion of Kendall's tau

and a sample of size 1000.

	Estimate	Std. Error	z value	Pr(> z)
rho.1	-0.01420116	0.033302051	-0.4264349	0.6697910
rho.2	0.97843954	0.001595654	613.1904797	0.0000000
rho.3	-0.01938295	0.033034985	-0.5867400	0.5573783

(2) Apply goodness-of-fit test.

- a) The empirical copula (from simulation model) and hypothesized copula are compared under the null hypothesis that they are from the same copula. Cramér-von-Mises (“CvM”) statistic S_n is used and the p -value is estimated using parametric bootstrapping method (simulated p value)

$$S_n = \int_{[0,1]^d} C_n(\mathbf{u})^2 dC_n(\mathbf{u}) = \sum_{i=1}^n \{C_n(\hat{\mathbf{U}}_i) - C_{\theta_n}(\hat{\mathbf{U}}_i)\}^2$$

$$C_n(\mathbf{u}) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}(U_{i1} \leq u_1, \dots, U_{id} \leq u_d),$$

$$\mathbf{u} = (u_1, \dots, u_d) \in [0, 1]^d.$$

This test can be realized by using function “gofCopula”. Please note that, due to the simulation number chosen to do parametric bootstrapping, the run takes a while to finish.

Modeling Loss Emergence and Settlement Processes

The goodness-of-fit test using CvM statistics was intended to be applied at dimension 3. However, the run speed is very slow and impractical. Only 100 simulations for calculating the p value take more than 20 minutes to finish. Therefore, I used 2 dimension tests 3 times to accomplish it. Only 100 simulations for parametric bootstrapping will take around 10 minutes for 2 dimensions. Ties will produce an incorrect p value calculation. We used pseudo observations to randomly break the ties as suggested in “Modeling Multivariate Distributions with Continuous Margins Using the copula R Package” by Ivan Kojadinovic and Jun Yan, *Journal of Statistical Software* May 2010, Vol 34, Issue 9.

The results below show that the simulated correlated frequencies can fit to normal copula very well.

Line 1&2

Parameter estimate(s): -0.002100962
Cramer-von Mises statistic: 0.0203318 with p -value 0.4009901

Line 1&3

Parameter estimate(s): 0.97926
Cramer-von Mises statistic: 0.007494245 with p -value 0.3811881

Line 2&3

Parameter estimate(s): -0.01049841
Cramer-von Mises statistic: 0.01614539 with p -value 0.5891089

- b) Kolmogorov-Smirnov (K-S) test. Only 2 dimensions are currently the maximum number allowable in R. Based on fitted copula, a random data set is generated and the empirical data and generated data are used to do 2-sample K-S test. The null hypothesis is that the two data sets are from the same continuous distribution. The statistic is

$D_{K-S} = \max_i |F_E(x_i) - F_H(x_i)|$. The p value will be calculated exactly if the sample size is less than 10,000.

Line 1&2 (Null hypothesis: correlation coefficient = 0)

Two-sample Kolmogorov-Smirnov test

data: x12 and y12
D = 0.016, p-value = 0.96
alternative hypothesis: two-sided

Modeling Loss Emergence and Settlement Processes

Line 1&3 (Null hypothesis: correlation coefficient = 0.99)

Two-sample Kolmogorov-Smirnov test

data: x13 and y13

D = 0.014, p-value = 0.9895

alternative hypothesis: two-sided

Line 2&3 (Null hypothesis: correlation coefficient = -0.01)

Two-sample Kolmogorov-Smirnov test

data: x23 and y23

D = 0.0135, p-value = 0.9933

alternative hypothesis: two-sided

The P value of the K-S test shows that we cannot reject the null hypothesis. This is also quite consistent with what we can spot from the scatter plots.

(3) Conclusion

Based on the scatter plots, normal copula fitting parameters, and also the goodness-of-fit test using CvM statistics, we can conclude that the correlation in simulated frequency data is consistent with our assumption.

7. POTENTIAL APPLICATIONS AND MODEL ENHANCEMENTS

The CAS Committee on Dynamic Risk Modeling plans to hold a 2011 Call Paper Program soliciting enhancements to the model, additional testing as well as papers applying the model to test alternative loss reserving methods and models. The following model enhancements would be welcome:

1. Include covariates as categorical or numeric variables in setting the parameters for the distributions. For example, in the modeling for claim size, if “state” is a categorical variable that affects the parameter of the distribution, we can include this information as input for the model, rather than setting up a separate “type” for each state. The covariate should be passed to the output detailed claim file.
2. Additional pairs (or groups) of variables whose sample values are correlated.
3. Input and output to and from standard database programs.

The following additional tests of the Loss Simulation Model would be welcome, and could be performed while testing alternative loss reserving methods and models:

Modeling Loss Emergence and Settlement Processes

1. Poisson frequencies have been tested, while negative binomial frequencies have not been tested.
2. Correlations of mean frequencies across lines of business have been tested, while correlations between size of loss and report lag within a Type have not been tested.
3. Ultimate claim values have been tested, while the sequence of loss reserve changes have not been tested.
4. The Single Payment Model has been tested, while the Periodic Payments Model and the Multiple Random Payments Model have not been tested.
5. The mechanics of certain parameters such as alpha, severity trend, inertia, etc. have been verified by examining the code. However, output from runs using these parameters has not been verified.
6. Apply the testing approach described in Section 3.2 above to assess the ability of simulated data from the Loss Simulation Model to represent “real” data.

Other tests, enhancements or topics not on these lists are also welcome as long as they address something specific to the Loss Simulation Model or its intended uses.

8. CONCLUSIONS

The LSMWP has developed a model that we hope will become a valuable tool in researching reserving methods and models. We hope that actuaries will use this model to:

- Better understand the underlying loss development process.
- Determine which methods and models work best in different reserving situations.
- Reflect this knowledge in evolving loss reserving practices.

Modeling Loss Emergence and Settlement Processes

BIBLIOGRAPHY

Categories:

1. General Interest
2. Building Models
3. Testing Simulated Data
4. Testing Reserving Methods

Abbreviations:

- CAS – Casualty Actuarial Society
 IAA – International Actuarial Association
 SOA – Society of Actuaries

REF #	TITLE:	CATEG ORIES:	PUBLICATION :	AUTHOR[S]:	PUBLISHER :	YEAR:	TOPICS:	TYPE OF MATERIAL:
1	Modeling Multivariate Distributions with Continuous Margins Using the copula R Package	1,2,3	Journal of Statistical Software	Kojadinovic, Ivan and Yan, Jun.	American Statistical Association	May 2010, Vol 34, Issue 9	Copulas, incl. discrete variables	Refereed Paper/Article
2	Loss Models From Data to Decisions	1,2,3	CAS Syllabus	Stuart Klugman, Harry Panjer, and Gordon Willmot	Wiley	Third ed., 2008	Probability Distributions	Book
3	Thinking Outside the Triangle	2,3,4	ASTIN Colloquium Papers 2007	Glenn Meyers	CAS: Arlington, Virginia	2007	Reserves	Refereed Paper/Article
4	Estimating Predictive Distributions for Loss Reserve Models	2,3,4	Casualty Actuarial Society Forum	Glenn Meyers	CAS: Arlington, Virginia	2006	Reserves	Refereed Paper/Article
5	Probability and Statistical Inference, 7th ed.	1,2		Hogg, Robert V. and Tanis, Elliott A.	Pearson Prentice Hall	2006	Probability and Statistics	Book
6	Extending the Linear Model with R	1,2		Faraway, Julian J.	Chapman and Hall/CRC	2006	Non-Linear models, GLM.	Book
7	Introduction to Mathematical Statistics 6th ed.	1,2		Hogg, R.V., McKean, J.W, and Craig, A.T.	Pearson Prentice Hall	2005	Probability and Statistics	Book
8	Survival Analysis Using S: An Analysis of Time-to-Event Data	1,2		Tableman, Mara and Kim, Jong Sung	Chapman and Hall/CRC	2004	Survival Models	Book
9	Toward a Unified Approach to Fitting Loss	2,3	Proceedings of the CAS	Jacques Rioux and Stuart Klugman	CAS: Arlington, Virginia	Working Paper, 2003	Reserves	Non-Refereed Paper/Article
10	On the Safety Loading for Chain Ladder Estimates: A Monte Carlo	4	CAS Forum	Schiegl, M	IAA: Brussels, Belgium	2002 Vol: 32:1 Page(s)		Refereed Paper/Article

Modeling Loss Emergence and Settlement Processes

REF #	TITLE:	CATEGORIES:	PUBLICATION :	AUTHOR[S]:	PUBLISHER :	YEAR:	TOPICS:	TYPE OF MATERIAL:
	Simulation Study); 107-128		
11	Modeling Size-of-Loss Distributions for Exact Data in WinBUGS	3	Journal of Actuarial Practice	David P.M. Scollnik	CAS: Arlington, Virginia	2002 Pages 10, 193-218	Reserves;Reinsurance Research	Non-Refereed Paper/Article
12	Understanding Relationships Using Copulas	1,2,3	North American Actuarial Journal	Frees, Edward W and Valdez, Emiliano A.	SOA : Schaumburg, Illinois	Volume 2, Number 1 (1998)	Copulas and related topics	Refereed Paper/Article
13	Some Extensions of J.N. Stanard's Simulation Model for Loss Reserving	2,4	CAS Forum	Vaughan, Richard L.	CAS: Arlington, Virginia	1998 Vol: Fall Page(s) : 415-498	Reserves	Non-Refereed Paper/Article
14	Testing the Assumptions of Age-To-Age Factors	2,4	ACM Portal: The Guide to Computing Literature	Venter, Gary G.	CAS: Arlington, Virginia	1998 Vol: LXXXV Page(s): 807-847		Refereed Paper/Article
15	Performance Testing Aggregate and Structural Reserving Methods: A Simulation Approach	2	CAS Forum	Rollins, John W.	CAS: Arlington, Virginia	1997 Vol: Summer, Vol 1 Page(s): 137-174	Reserves	Non-Refereed Paper/Article
16	A Comparative Study of the Performance of Loss Reserving Methods Through Simulation	2,4	CAS Forum	P. Narayan and T. Warthen	CAS: Arlington, Virginia	1997Summer, Volume 1, Pages 175-195	Reserves	Non-Refereed Paper/Article
17	Simulation Procedure for Comparing Different Claims Reserving Methods; A	2,4	CAS Forum	Pentikäinen, Teivo; Rantala, Jukka	CAS: Arlington, Virginia	1995 Vol: Fall Page(s) : 128-156	Reserves	Non-Refereed Paper/Article
18	Note on Simulation of Claim Activity for Use in Aggregate Loss Distributions; A	2	Casualty Loss Reserve Seminar Transcript	Lyons, Daniel K.	CAS: Arlington, Virginia	1994 Vol: Spring, Vol 1 Page(s): 357-392	Reserves;Reinsurance Research	Non-Refereed Paper/Article
19	Simulation Models for Reserve and Surplus Analysis	2		Kreps, Rodney E.; Englander, Jeffrey A.	CAS: Arlington, Virginia	1993	Reserves	Casualty Loss Reserve Seminar Transcript
20	The Computation of Aggregate Loss Distributions	2	Proceedings of the CAS	John P. Robertson	CAS: Arlington, Virginia	LXXIX 57-133, 1992	Reserves	Non-Refereed Paper/Article
21	Validation and verification of simulation models	3	Proceedings of the CAS	Robert G. Sargent	ACM Press	1992	Model Validation and Verification	
22	Generalized Linear Models	1,2,3		McCullagh, P. and Nelder, J.A.	Chapman & Hall/CRC	2nd ed., 1989	Predictive Modeling, GLMs	Book

Modeling Loss Emergence and Settlement Processes

REF #	TITLE:	CATEGORIES:	PUBLICATION :	AUTHOR[S]:	PUBLISHER :	YEAR:	TOPICS:	TYPE OF MATERIAL:
23	Guide to Simulation; A	1,2,3	Proceedings of the CAS	Bratley, P.; Fox, B. L.; Schrage, L. E.	Springer-Verlag : New York	1987 Vol: 2	Risk Theory	Book
24	Simulation Test of Prediction Errors of Loss Reserve Estimation Techniques; A	2,4	Proceedings of the CAS	Stanard, James N.	CAS: Arlington, Virginia	1985 Vol: LXXII Page(s): 124-148	Reserves; Risk Theory	Refereed Paper/Article
25	Transformed Beta and Gamma Distributions and Aggregate Losses	2	CAS Forum	Venter, Gary G.	CAS: Arlington, Virginia	LXX 156-193, 1983		Refereed Paper/Article
26	The Calculation of Aggregate Loss Distributions from Claim Count Distributions	2,3	Proceedings of the CAS	Phil Heckman and Glenn Meyers	CAS : Arlington, Virginia	LXX, 22-61, 1983	Reserves	Non-Refereed Paper/Article
27	Computer Simulation and the Actuary: A Study in Realizable Potential	1,2	Proceedings of the CAS	Arata, David A.	CAS : Arlington, Virginia	1981 Vol: LXVIII Page(s): 24-64		Refereed Paper/Article
28	Estimating Casualty Insurance Loss Amount Distributions	2,3	Transactions of the SOA	Venter, Gary G.	CAS : Arlington, Virginia	LXVII, 57-109, 1980	Reserves	Casualty Loss Reserve Seminar Transcript
29	The Aggregate Claims Distribution and stop-Loss Reinsurance	2	Wiley	Panjer, Harry	SOA : Schaumburg, Illinois	XXXII, 523-535, 1980	Risk Theory	
30	Automobile Collision Deductibles and Repair Cost Groups: The Lognormal Model	1,2		Bickerstaff, David R.	CAS : Arlington, Virginia	1972	Collision losses, lognormal distribution	Proceedings Article
31	Modern Applied Statistics with S, Fourth Edition	2		Venables, W.N. and Ripley, B.D.	Springer-Verlag	2002	Statistics and "R" language	Book
32	Probability and Statistics for Engineering and the Sciences, Fifth Edition	2		Devore, J.L.	Duxbury Thomson Learning	2000	Probability and Statistics	Book

Modeling Loss Emergence and Settlement Processes

REF	LINK	COMMENTS
1	http://www.jstatsoft.org/v34/i09/paper	Comments: Thorough discussion of using copulas in R. Describes how to adapt copulas to discrete random variables. Goodness of fit tests. Illustrations using insurance Loss Data
2	-	Comments: This book is devoted to the problem of fitting parametric probability distributions to data. This treatment unifies loss modeling in one book. Emphasis is made on the distribution of single losses related to claims experienced against various types of insurance policies. The book includes five sets of insurance data as examples.
3	http://www.actuaries.org/ASTIN/Colloquia/Orlando/Papers/Meyers.pdf	Comments: Overfitting is the result of having a model that is too complex for the amount of data that is available. This is typically the case when for a loss reserve model which has a large number of parameters on a smallish triangle of data. This paper uses simulation methods to diagnose overfitting in estimating the predictive distribution of loss reserves by the method of maximum likelihood. This paper then shows how to use a Bayesian fitting methodology to overcome overfitting using prior information that is "outside the triangle."
4	http://www.casact.org/pubs/forum/06forum/163.pdf	Comments: This paper demonstrates a Bayesian method for estimating the distribution of future payments of individual insurers. The main features of this method are: (1) the stochastic loss reserving model is based on the collective risk model; (2) predicted loss payments derived from a Bayesian methodology utilizing large/stable insurers as its prior information; (3) applying tests on large numbers of insurers and demonstrating its predictions are within the statistical bounds expected for a sample of its size. It concludes with an analysis of reported reserves and their subsequent developments in terms of the predictive distribution calculated by this Bayesian methodology.
5	-	Comments: A book on probability and statistics. This book is easier to understand than the Hogg, McKean, and Craig book.
6	-	Comments: This is an advanced textbook covering topics beyond linear models, such as Generalized Linear Models, mixed-effects models and non-parametric models. The book uses the programming language "R". It assumes knowledge of linear models. The author has written a predecessor book on Linear Models.
7	-	Comments: A book on probability and statistics frequently used in university statistics classes. This book contains chapters on Statistical Inference, Maximum Likelihood Methods, and Inferences about Normal Models.
8	-	Comments: This book covers parametric and non-parametric survival models and provides modeling examples in the statistical language S (note: the language R is closely related to S). Survival models are useful in modeling time-to-event variables such as report lag and settlement lag and also claim size distribution. They handle censored data well.
9	-	Comments: There are two components to fitting models - selecting a set of candidate distributions and determining which member fits best. It is important to have the candidate set be small to avoid overfitting. Finite mixture models using a small number of base distributions provide an ideal set. Because actuaries fit models for a variety of situations, particularly with regard to data modifications, it is useful to have a single approach. Though not optimal or exact for a particular model or data structure, the method proposed in this paper should be reasonable for most all settings. A computer program implementing these models and techniques is provided.
10	http://www.casact.org/library/astin/vol32no1/107.pdf	Comments: A method of analyzing the risk of taking a too low reserve level by the use of Chain Ladder method is developed. We give an answer to the question of how much safety loading in terms of the Chain Ladder standard error has been added to the Chain Ladder reserve in order to reach a specified security level in loss reserving. This is an important question in the framework of integrated risk management of an insurance company. Furthermore we investigate the relative bias of Chain Ladder estimators. We use Monte Carlo simulation technique as well as the collective model of risk theory in each cell of run-off table. We analyse deviation between Chain Ladder reserves and Monte Carlo simulated reserves statistically. Our results document dependency on claim number and claim size distribution types and parameters.

Modeling Loss Emergence and Settlement Processes

REF	LINK	COMMENTS
11	-	Comments: This paper discusses how the statistical software WinBUGS can be used to implement a Bayesian analysis of several popular severity models applied to exact size-of-loss data. The particular models targeted are the gamma, inverse gamma, loggamma, lognormal, (two-parameter) Pareto, inverse (two-parameter) Pareto, Weibull, and inverse Weibull distributions. It is possible to implement additional size-of-loss models (including those for truncated data) using analogous methods.
12	www.soa.org/library/journals/north-american-actuarial-journal/1998/january/naaj9801_1.pdf - 2009-05-07	Comments: Excellent introduction to copulas. Includes fitting models to data, simulation, and several insurance company data applications.
13	http://www.casact.org/pubs/forum/98forum/vaughn.pdf	Comments: The loss process model and simulation procedures proposed by James M. Stanard in 1985 are extended in numerous ways, including provision for serial autocorrelation of parameters, mixtures of claim types, conditional selection of sample points, and a much greater variety of reserving methods. The extended model is used to explore many questions arising in practical loss reserving and to assist the loss reserver in choosing the best estimator for particular data conditions.
14	http://www.casact.org/pubs/proceed/proceed98/980807.pdf	Comments: The use of age-to-age factors applied to cumulative losses has been shown to produce least-squares optimal reserve estimates when certain assumptions are met. Tests of these assumptions are introduced, most of which derive from regression diagnostic methods. Failures of various tests lead to specific alternative methods of loss development.
15	http://www.casact.org/pubs/forum/97sforum/97sf1137.pdf	Comments: Aggregate financial data histories are extensively used by actuaries in projecting ultimate liabilities, but the claim occurrence, reporting, and settlement process which generates these data is not perfectly understood, rarely modeled directly, and not incorporated into the structure of most popular reserving methods. This paper utilizes today's computer applications to create an automated simulation tool. This module allows the user to choose statistical assumptions for each element of the claims process and specify the structure of the simulation experiment. It then generates random paid loss and claim count histories based on these inputs. The module is used to perform an experimental test of the performance of aggregate versus structural reserving methods. The methods chosen are the paid loss development method and a new "closed claim cost" structural method. In each trial, a database of ten accident years at ten annual evaluations is simulated. Then both methods are run at five successive calendar year evaluations of the simulated data. Several error functions are tabulated at each valuation date and the speed of the approach of each method's indication to the true value of the aggregate costs is examined.
16	http://www.casact.org/pubs/forum/97sforum/97sf1175.pdf	Comments: This paper uses four methods to simulate loss development triangles and uses the simulated triangles to compare two traditional actuarial reserving methods against three regression methods of reserving. Section IV on Comparison of Procedures provides good insights into using simulation to test reserving methods.
17	http://www.casact.org/pubs/forum/95fforum/95ff128.pdf	Comments: The estimation of outstanding claims is one of the important aspects in the management of the insurance business. Exploration of the inaccuracies involved is traditionally based on a post-facto comparison of the estimates against the actual outcomes of the settled claims. However, until recent years it has not been usual to consider the inaccuracies inherent in claims reserving in the context of more comprehensive (risk theoretical) models. Important parts of the technique which will be outlined in this paper can be incorporated into over-all risk theory models to introduce the uncertainty involved with technical reserves as one of the components in solvency and other analyses. The idea in this paper is to describe a procedure by which one can explore how various reserving methods react to fictitious

Modeling Loss Emergence and Settlement Processes

REF	LINK	COMMENTS
		variations, fluctuations, trends, etc. which might influence the claims process and how they reflect on the variables indicating the financial position of the insurer.
18	http://www.casact.org/pubs/forum/94spforum/94spforum/357.pdf	Comments: Aggregate loss distributions have been used in a number of different applications over the last few years. These applications have usually focused on the distribution of losses at ultimate or final values and have not studied how losses move to ultimate values over time. The approach outlined in this note models claim activity through the use of transition matrices. Individual claim activity is then incorporated into an aggregate loss simulation model to determine a number of distributions of interest. Keywords: Confidence Estimates, Loss Development, IBNR, Reinsurance Research - Loss Distributions, Size of
19		Comments: Simulation models have a long history of use in the actuarial profession for estimating confidence intervals for loss reserves and pricing estimates. More recently, models have been developed which can be used to derive probability levels for insurance company surplus. In this session we will describe procedures for modeling loss reserve uncertainty using simulation models. Models utilizing claim count and claim severity distribution will be discussed. The issue of parameter variance and its estimation will be presented. Techniques for estimating payout patterns and its impact on discounted reserve variability will be presented. A more comprehensive model of insurance company operations which can be used to assess insurance company solvency will then be presented. The items modeled will include premiums, claims outstanding, investment returns, business cycles, and effect of inflation.
20	http://www.casact.org/pubs/proceed/proceed92/	Comments: This paper provides an application of the Fourier transform as a tool for computation of aggregate loss distributions from arbitrary frequency and severity distributions. The paper offers a complete algorithm and provides examples to allow its implementation in various computer languages. The final section contains a discussion of excess loss distributions where computation is not limited to the fast Fourier transform based algorithm.
21	http://www.informs-sim.org/wsc99papers/005.PDF	Comments: Referenced in S.M. Ross Simulation Text (2002 Academic Press, 6277 Sea Harbor Drive, Orlando, FL 32887), this article defines model validation and verification. The section on operational validity describes approaches helpful in assessing testing model results against known data points.
22		Comments: This is an excellent reference on GLMs and Poisson GLMs. This is often considered the definitive early work on Generalized Linear Models.
23		
24	http://www.casact.org/pubs/proceed/proceed85/85124.pdf	Comments: Contains fully described simulation models. Abstract: This paper uses a computer simulation model to measure the expected value and variance of prediction errors of four simple methods of estimating loss reserves. Two of these methods are new to the Proceedings. The simulated data triangles that are tested are meant to represent sample sizes typically found in individual risk rating situations. The results indicate that the commonly used age-to-age factor approach gives biased estimates and is inferior to the three other methods tested. Theoretical arguments for the source of this bias and a comparison of two of the methods are presented in the Appendices.
25	http://www.casact.org/library/astin/vol32no1/107.pdf	Comments: Distribution functions are introduced based on power transformations of beta and gamma distributions, and properties of these distributions are discussed. The gamma, beta, F, Pareto, Burr, Weibull and loglogistic distributions are considered. The transformed gamma is used to model aggregate distributions by matching moments. The transformed beta is used to account for parameter uncertainty in this model. Calculation procedures are discussed and APL program listings are included. The transformed gamma distribution is compared to exact methods of computing the aggregate distribution function based on the entire frequency and severity distributions.

Modeling Loss Emergence and Settlement Processes

REF	LINK	COMMENTS
26	http://www.casact.org/pubs/proceed/proceed83/83022.pdf	Comments: This paper discusses aggregate loss distributions from the perspective of collective risk theory. An accurate, efficient and practical algorithm is given for calculating cumulative probabilities and excess pure premiums. The input required is the claim severity and claim count distributions. One of the main drawbacks of the collective risk model is the uncertainty of the parameters of the claim severity and claim count distributions. Modifications of the collective risk model are proposed to deal with these problems. These modifications are incorporated into the algorithm. Examples are given illustrating the use of this algorithm. They include calculating the pure premium for a policy with an aggregate limit; calculating the pure premium of an aggregate stop-loss policy for group life insurance; and calculating the insurance charge for a multi-line retrospective rating plan, including a line which is itself subject to an aggregate limit.
27	http://www.casact.org/pubs/proceed/proceed81/81024.pdf	Comments: This paper argues that computer simulation is an underappreciated and, therefore, underutilized casualty actuarial resource. In so contending, "Computer Simulation and the Actuary" discusses five applications of Monte Carlo computer simulation to everyday actuarial problems: establishing full credibility standards; testing the solidity of new, limited purpose insurance companies; pricing difficult or catastrophic exposures; customizing casualty insurance charges and excess loss premium factors; and developing loss reserve confidence intervals. Illustrations of appropriate simulation solutions to each of these problems are provided.
28	-	Comments: This paper concentrates upon probability model-building and statistical techniques for estimating and testing the model parameters. A general procedure for selecting a "best" parameterized model based upon loss amount data. This solves only part of a broader problem, which is to estimate loss amount distributions for future coverage periods or future final-valued loss amount distributions for past coverage periods where the losses are not all settled or even known. To solve this broader problem it is necessary to specify models of the overall insurance loss processes, defining how the future relates to the past and how the individual insured relates to the whole insurance portfolio.
29		Comments: This paper describes a method for approximating the distribution of aggregate claims for a group life insurance contract. It is found that a compound Poisson process appropriately modeled the aggregate claims distribution, based on the collective risk assumption, which states that each life which leaves the group by death claim is immediately replaced by a life with identical mortality characteristics. The resulting method produces accurate and useful results, and is relatively easy to use as long as the number of distinct amounts of insurance in the group is not too great. In practice, this condition is achieved by rounding all face amounts to be integral multiples of some convenient unit.
30	-	Comments: This paper applies the theory of lognormal distributions to the study of Collision losses, which are left-truncated and shifted. The book discusses many properties of the lognormal, such as the distribution of moments.
31		Comments: This book is a practical book covering Probability and Statistics. It has a large number of well-explained examples.
32		Comments: This book is the most well-known book covering R (and the related language S+) and its application to data analysis. The first part of the book discusses the structure and basic elements of the language.

Acknowledgment

The authors acknowledge Mark Shapland for his vision in creating the LSMWP and for his encouragement and very helpful advice as work proceeded. The authors also acknowledge Richard Vaughan and Huan Zhu for the very major contributions they made in developing and testing a prototype model during the LSMWP's planning phase.

Supplementary Material

See Section 5 above on Documentation of Open Source Model. Also, please refer to the results of our search of the actuarial literature summarized in the Bibliography.

Abbreviations and notations

AL-BI, Auto Liability Bodily Injury

AL-PD, Auto Liability Property Damage

APD, Auto Physical Damage

AY, accident year

DRM Committee, Committee on Dynamic Risk Modeling

GLM, Generalized Linear Model

K-S test, Kolmogorov-Smirnov test

LSMWP, Loss Simulation Model Working Party

MLE, Maximum Likelihood Estimation

NLL, Negative Log-Likelihood

P-P plot, Probability-Probability or Percent-Percent plot

Q-Q plot, a Quantile-Quantile plot

R, an implementation of the S programming language

S, a statistical programming language

S-PLUS, a commercial implementation of the S programming language sold by TIBCO Software Inc..

Working Party Oversight, Model Building and Testing Team Members

Robert Bear is principal and founder of RAB Actuarial Solutions LLC, which offers the following consulting services: (1) loss reserve studies and research on loss reserving methods and models (2) development of dynamic risk models to facilitate evaluation of profitability as well as risk load and capital needs (3) insurance and reinsurance pricing, including reinsurance commutation, excess pricing and price monitoring studies (4) resolution of loss reserve and coverage disputes subject to insurance arbitration, reinsurance arbitration or mediation (5) actuarial and reinsurance expert witness and litigation support.

Bear currently serves as Chairperson of the CAS Dynamic Risk Modeling Committee. He previously served as Chairperson of the Reinsurance Association of America Actuarial Committee and as President of Casualty Actuaries in Reinsurance. He has authored several CAS discussion papers and articles on reinsurance pricing, loss reserving, and risk modeling issues. Additional information is available on his web site, www.rabsolutions.net.

Modeling Loss Emergence and Settlement Processes

Joseph Marker, an FCAS and MAAA, teaches Actuarial Science and Financial Mathematics at the University of Michigan. Joe is principal and co-founder of Marker Actuarial Services, LLC, which he and his wife Connie started in 2001. Marker Actuarial provides Property-Casualty actuarial and management consulting services. Actuarial consulting concentrates on predictive modeling, pricing, loss modeling, financial analysis, actuarial management, and service to small and mid-size insurers. Additional information is available at www.markeractuarial.com.

Prior to consulting, Joe worked for four organizations in various actuarial capacities for 28 years, the last fifteen years as Chief Actuary at two regional insurance companies.

He is a past president of the Midwestern Actuarial Forum. He published the paper, *Studying Policy Retention using Markov Chains*, PCAS LXXXV, 1998. Joe also co-authored the paper *Rating Claims-Made Insurance Policies*, Call Paper Program 1980. This paper was part of the CAS exam syllabus for many years.

Hai You is the software developer of the Loss Simulation Model that is the subject of this paper and has co-authored sections of this paper that document the model as well as the associated documentation within the program and on the CAS LSMWP web page. He also contributed software that assisted with the testing of the Loss Simulation Model.

As VP, technology at Gooouon, Hai is responsible for actuarial modeling, product development and innovative solutions. He is the author of ReserveMaster, the cutting-edge loss reserving platform for P&C insurers. Hai developed Illustration Corpula and contributed to CAS. His 3-day-project of Two Stage Bootstrap in R language has more than 100 downloads since created. Before co-founding Gooouon, Hai has served as an IT consultant in various insurance companies for 10 years, where he got insight knowledge of the core systems such as underwriting, claim, and rate making. Hai masters various computer languages, securities, databases, 3D, email, UI, data communication. He also believes that Object-Oriented-R is 10 times stronger than R. Just for fun, he developed his own games like ChessBot, Maze (using AStar), Mine Digger, etc.

Gooouon is an actuarial engineering company, where people integrate actuarial knowledge and analytical solutions into services. The company offers consulting services on pricing, reserving, statistical modeling, data mining, and customized actuarial solutions. Besides ReserveMaster, the company also provides tailor-cut Loss Data Processing, Health Insurance Solutions, iERM - Insurance Enterprise Risk Management and Depict - a drag and drop Enterprise Workflow Management platform. Please visit <http://www.gooouon.com> for additional information.

Glenn Meyers is Vice President of Research for ISO Innovative Analytics. He earned a Ph.D. in mathematics from the State University of New York at Albany and is a Fellow of the Casualty Actuarial Society and a member of the American Academy of Actuaries. Before joining ISO in 1988, Glenn worked at CNA Insurance Companies and the University of Iowa.

Glenn's current responsibilities at ISO include the development of insurance scoring products using predictive modeling techniques. Prior responsibilities have included working on ISO Capital Management products, increased limits, catastrophe ratemaking and reinsurance products.

Modeling Loss Emergence and Settlement Processes

Glenn's work has been published in Proceedings of the Casualty Actuarial Society (CAS) and the new CAS peer-reviewed journal, *Variance*. He also writes a regular column in the *Actuarial Review*. His papers have won numerous awards, and he is a frequent speaker at CAS meetings and seminars.

His service to the CAS includes membership on various education and research committees. He has served on the CAS Board of Directors and currently serves as the IAA's delegate to ASTIN. His work on the LSMWP was to develop statistical tests to compare the outcomes of a simulation model with outcomes produced by real data.

Curtis Parker, FCAS, CPCU, VP Chief Actuary Grange Insurance Companies. Curt is an actuary with 36 years of varying experiences in the Property/Casualty Insurance Industry including personal and commercial lines underwriting, pricing, product research and development, market planning, reserving, and predictive modeling. He has volunteered over the years on examination, ratemaking, and reserving committees/panels. He served as chairperson for the subcommittee of the LSMWP committee which completed research and documented literature providing general and more technical readings related to loss simulation considerations and techniques. His subcommittee was also charged with development of initial methods for testing output from the model with respect to the ability to distinguish it from actual historic data. The subcommittee is indebted to Glenn Myers for his significant contributions with this later portion of the subcommittee's charge.

Kailan Shang works in the area of financial risk management in AIA. Prior to this, he worked as a pricing actuary in a life insurance company for 2.5 years. Years of actuarial and risk management experience has allowed him to get a broad exposure, including Economic Capital, MCEV, financial engineering, dynamic management options, dynamic policyholder behavior modeling, product development and management, US GAAP reporting, dynamic solvency testing, etc.

As an FSA, CFA, PRM and SCJP, he is also an enthusiast of actuarial research through both volunteer works and funded research program. He participated in IAA Comprehensive Actuarial Risk Evaluation project and he is now working in a SOA research project "Valuation of embedded option in Pension plan" as the lead researcher.

Yuting Yang is a graduate student in Actuarial Science and Financial Mathematics at the University of Michigan. Prior to attending Michigan, she earned a master's degree in mathematics from the University of California at Davis. At UC Davis, she did research in probability, random matrices, and convex geometry. Yuting did most of the work in testing the severity distributions in section 6.2.2.

Working Party Members

<u>Working Party Members</u>	<u>Position</u>
Robert A. Bear	Co-Chairperson
Mark R. Shapland	Co-Chairperson
Ramzi AbuJamra	Member
Shobhit Awasthi	Member
Hassan A. Ayoub	Member
Glen Barnett	Member
Nebojsa Bojer	Member
Bhaskar Chattaraj	Member
Denise L. Cheung	Member

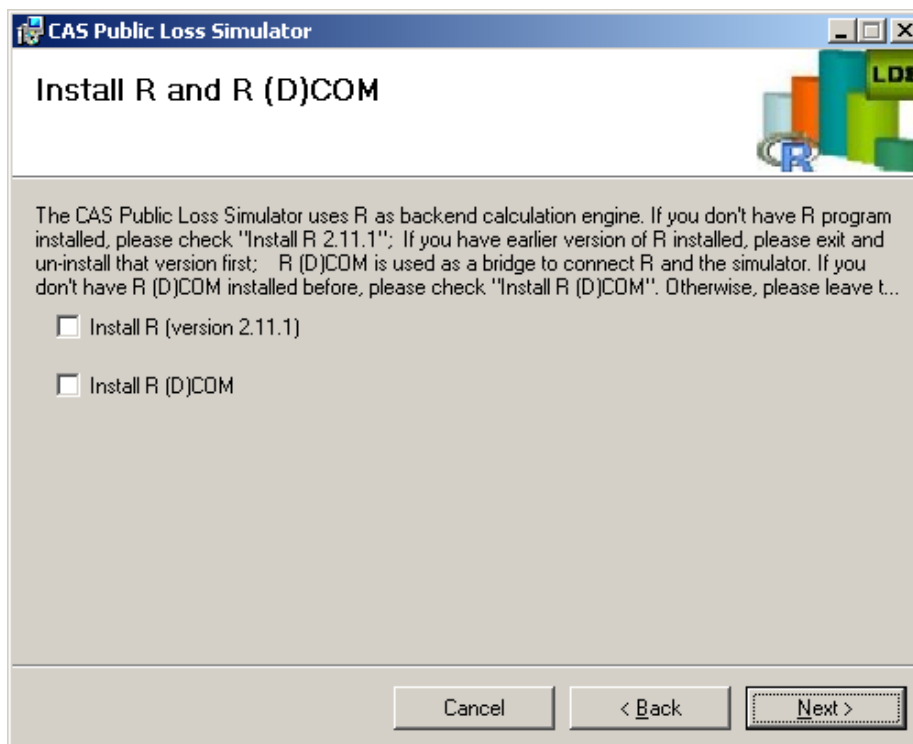
Modeling Loss Emergence and Settlement Processes

Wei Chuang	Member
Kevin M. Cleary	Member
Catherine Cresswell	Member
Salvatore Forte	Member
Bradford S. Gile	Member
Spencer M. Gluck	Member
Songling Guo	Member
Thomas Hartl	Member
Ping-Hung Hsieh	Member
Nicole Huang	Member
Turab Hussain	Member
Li Hwan Hwang	Member
Jan I. Iwanik	Member
Julia Jacobi	Member
Shiwen Jiang	Member
Nancy A. Kelley	Member
Stephen Jacob Koca	Member
Andrew M. Koren	Member
Scott C. Kurban	Member
Kin Hoe Lee	Member
Kin Yee Lee	Member
Stephen L. Lienhard	Member
Joseph O. Marker	Member
Glenn G. Meyers	Member
Jonathan E. Miller	Member
Jie Min	Member
F. James Mohl	Member
Curtis M. Parker	Member
Marco Pirra	Member
Arlie J. Proctor	Member
Ralph Stephen Pulis	Member
Keith A. Rogers	Member
Manalur S. Sandilya	Member
Kailan Shang	Member
Catherine E. Staats	Member
Christopher M. Steinbach	Member
Lin Yee Tan	Member
Varsha A. Tantri	Member
Jack T. Tower	Member
Daniel M. Van der Zee	Member
Justin M. VanOpdorp	Member
Richard L. Vaughan	Member
Gary G. Venter	Member
Yuting Yang	Member
Yuanhe (Edward) Yao	Member
Bo Zhou	Member
Hongbo Zhou	Member
Huan Zhu	Member
Jane E. Fulton	Staff Liaison

Appendix A

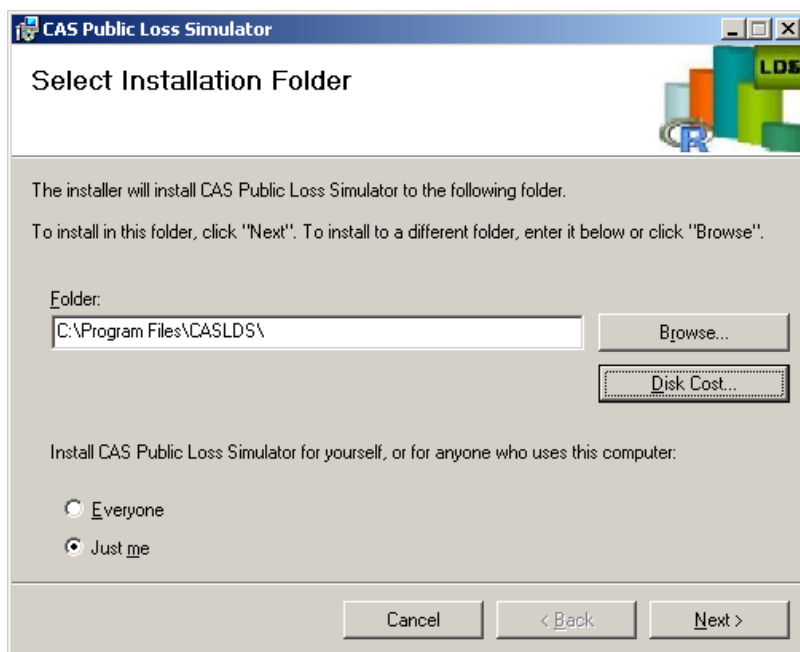
1) How to install the Public Loss Simulation Model

- a. To install the Loss Simulator Model to your computer, please download and run the windows installation package LossSimulatorSetup.msi from the following link
http://www.gooouon.com/loss_simulator/bin/LossSimulatorSetup.msi
- b. Following the initial screen instructions, the installation package will ask you two questions, as shown in Picture (1). The model runs R as the background calculation engine, and requires R (D) COM as a bridge component between R and the front end application. If you have not installed these two services before, please check them and continue. If you have earlier version of R installed, please exit and uninstall that version first, then come back and check the “Install R (version 2.11.1)”. Ignore the “Install R” option if you have higher version of R installed already.



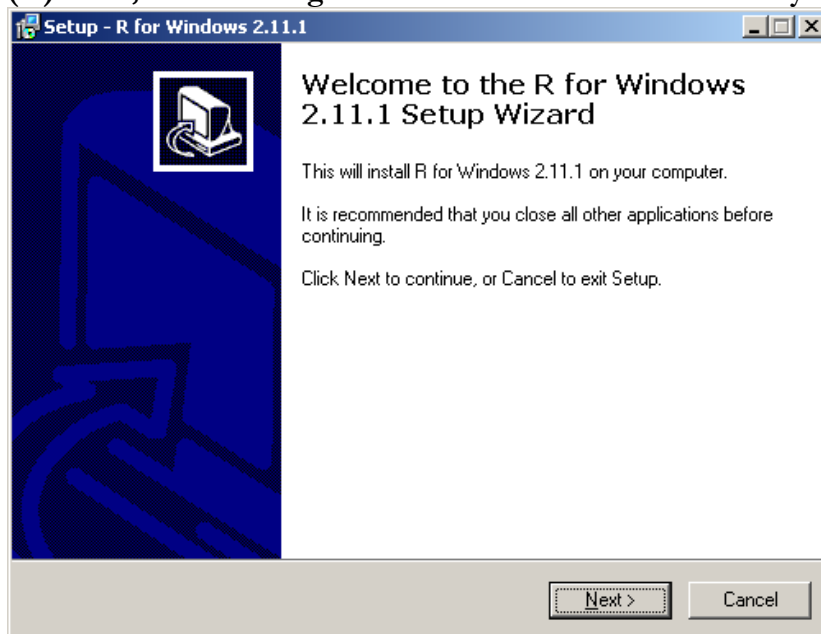
Picture (1) Installation Options

- c. Then provide an installation location, as shown in Picture (2)



Picture (2), Choose installation directory

- d. Click the next button to continue the installation. **If you choose to install R 2.11.1 and R (D)COM, the following two screens will start automatically.**

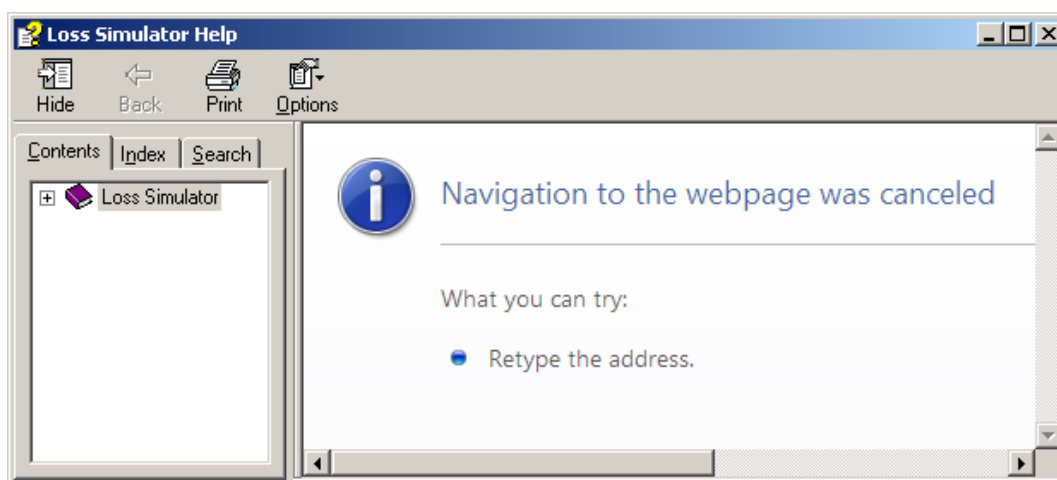


Picture (3), Install R 2.11.1



Picture (4), Install R (D) COM

- e. There is one more thing to be aware: besides the CAS online help at <http://www.casact.org/research/lsmwp/losshelp/index.cfm?fa=main>, Loss Simulator also contains a windows html help file as an attached help system. But due to Windows security restriction (<http://support.microsoft.com/kb/902225>), you may not be able to see the help content when launching help from the Simulator, and get the following error page instead.



Picture (5), Possible error when trying to open help file from Loss Simulator help menu.

In this case, you can just simply go to the **model installation folder** and right click the help file named **LossSimulator.chm**, and then select **Properties**, click **Unblock**. After this action, you will be able to see the help content.

2) How to run the Public Loss Simulation Model

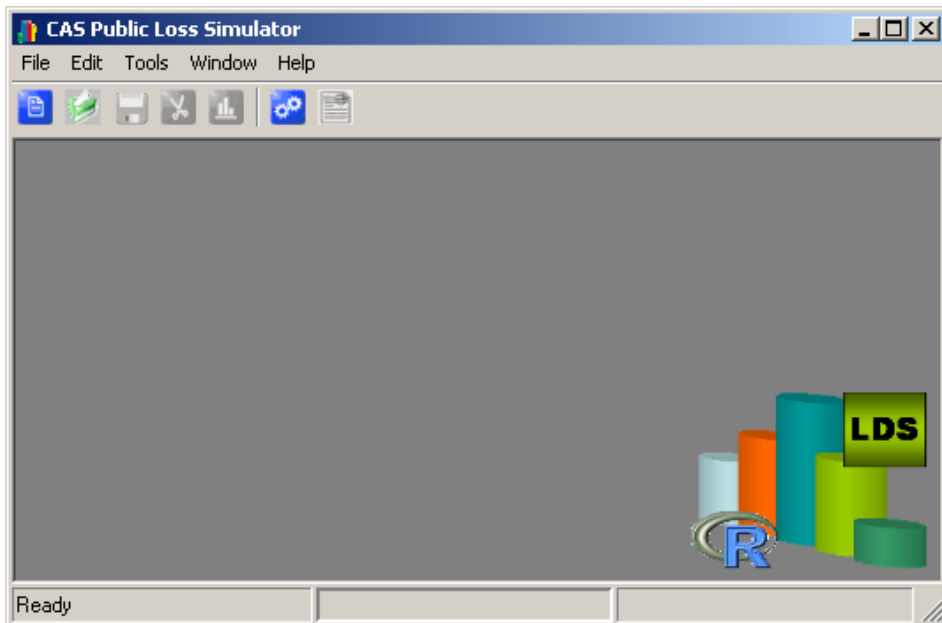


Public Loss Simulator.Ink

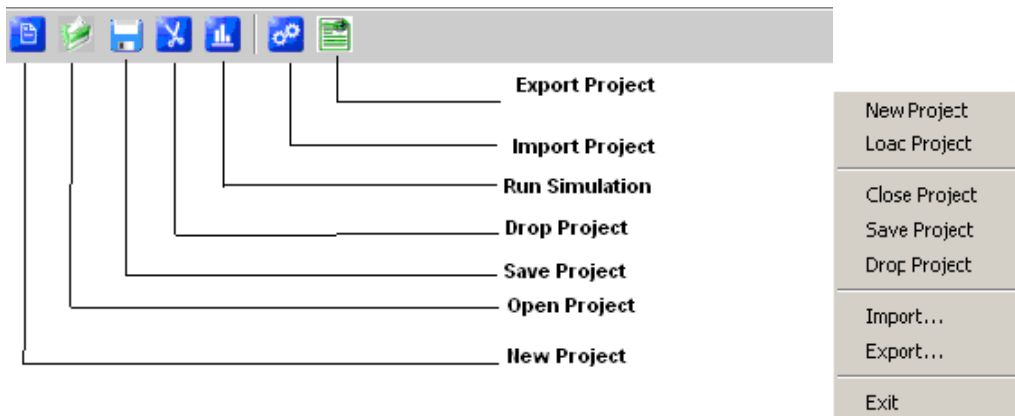
After installation, you will find the Public Simulator icon on both your desktop and Start Menu. You can just double click it to launch the application.

2.1 System Overview.

Initial screen will be like Picture (5). The application is developed within Windows UI standard, so that it contains menus and tool bar buttons associated with each menu items. Here is a brief explanation of each menu item and the later chapters will explain each of them in detail.



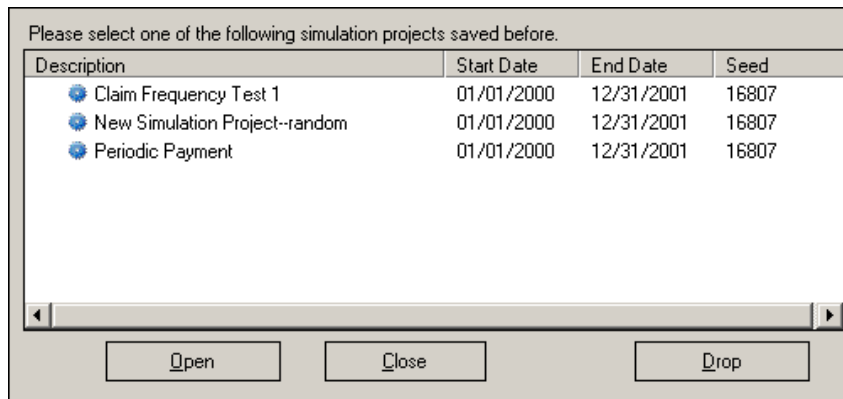
Picture (5) Public Loss Simulator main window



System Tool Bars and Menus

File, New Project : Initialize a brand new simulation project, which contains default settings for one line, with one Single Payment type. You can save it later.

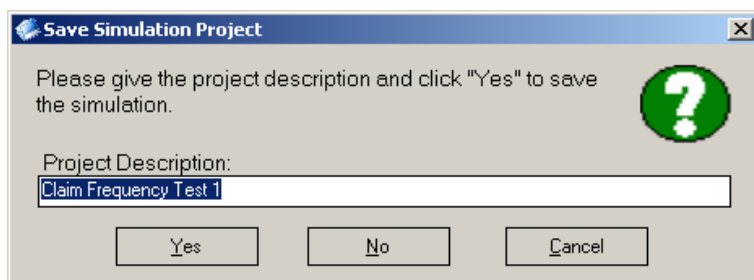
File, Load Project : Open a previously saved simulation project from the database.



Picture (6) Load previous saved project from database

File, Close Project : Close the current opened project. If system detects any change to the project properties, it will first ask for saving the project.

File, Save Project : Save the current opened project into database. When saving a simulation project, all the simulation properties will be saved into database. If it is initialized from a new project menu, system will ask you for a project description, as shown in Picture (7). System has a default MS ACCESS database attached; you can configure the simulator to connect to any kind of database server also.



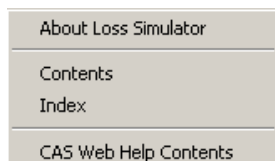
Picture (7) Project Save Window will pop up asking for project description

File, Drop Project : Permanently delete a project previously saved in the database. This menu item is enabled when a project is open.

File, Import... :

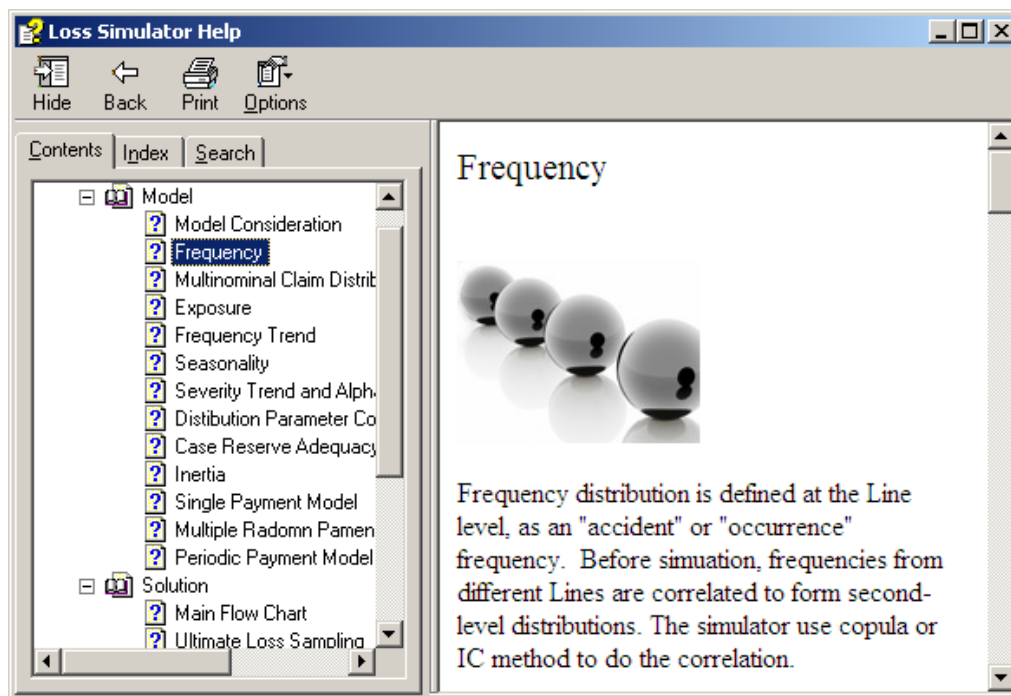
File, Export... : These two menus enable project sharing among users. You can export a project into a XML file and email it to another person. That person can import it and run simulation from his workstation, with all the project properties set by you. Or you can export the project and later import it again for different testing stage purpose. In that case, it is pretty similar to saving project to the database and loading project back.

This is help menu



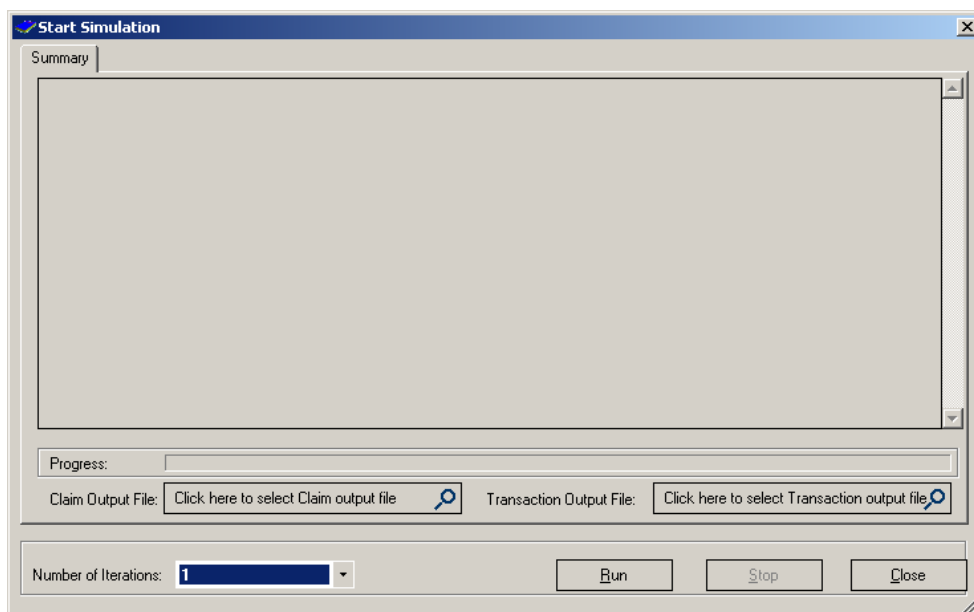
Help, CAS Web Help Contents : This will lead user to the help contents from CAS simulator working party web site.

Help, Contents : This will open the attached simulator help system, as shown in Picture (8). **If you cannot see any help contents from right side panel, please refer to section 7.1 (e) in this paper for proper configuration.**



Picture (8) Loss Simulator Help system

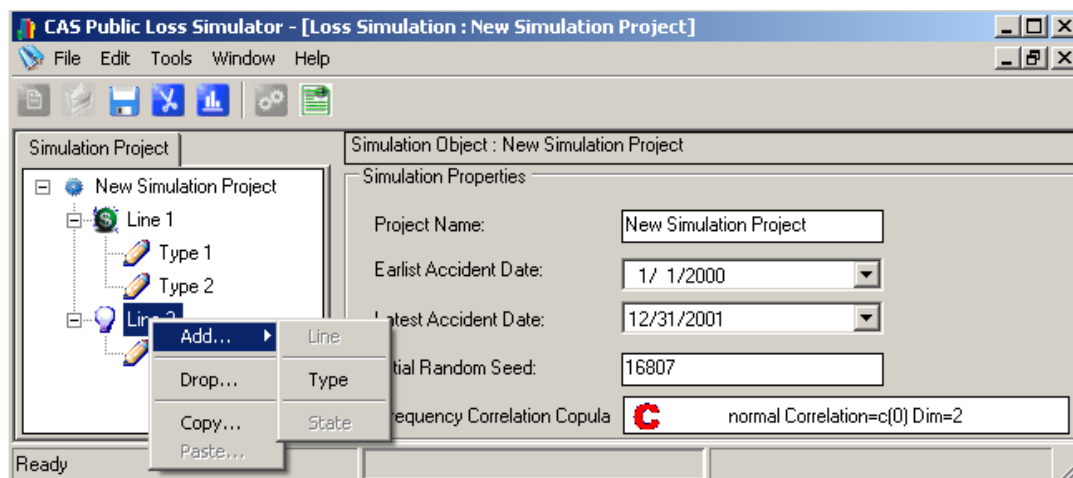
Tools, Run Simulation : This menu will let you start to run the simulation. As shown in Picture (9), you need to provide a claim file name, a transaction file name, and depending upon your needs, choose how many iterations you prefer. For a full scale simulation, you may try 1000 iterations. That could take hours to finish, generating hundred thousands of claims and millions of transactions. We will explain the simulation result in detail from later chapter.



Picture (9) Run Simulation window

2.2 Simulation Project.

A simulation project is one where modeler enters properties for line of businesses and types and initializes the simulation. A typical project can contain multiple lines of business (LOB), with each line of business consisting of one or more types of claims. Type can be treated as coverage in real case. A tree structure is the best way to describe the relationships among them, as shown in Picture (10).



Picture (10), a simulation project is constructed by tree structure

The left panel in the project window uses popup menu to allow user to add, drop, update, copy and paste nodes to build up a company business model; while the right panel of

the project displays properties of the selected node from left side. It could be a project properties, line level properties, or type level properties. This section of the paper won't give a detailed explanation of each property listed, most of which are explained in the help file or the CAS website online help for this model.

2.3 Overall Simulation Properties.

According to the tree relationship described in previous section, the simulation project will have properties that control line level correlations, as well as main simulation features such as the earliest accident date and the latest accident date.

As seen in Picture (11) below, when you double click the Simulation Project node, you will see the screen on right side panel. You can define a 10-year accident date range, starting from 01/01/2000 to 12/31/2009, as shown in the sample. This is a common time frame for most loss reserve analysis. In the later simulation stage, system will loop through each year and apply the line level properties, such as frequencies, etc, to generate synthetic claims and transactions.

The image shows a software dialog box titled "Simulation Properties". It contains several input fields and a button. The fields are: "Project Name" with the value "Test Simulation 1"; "Earliest Accident Date" with a dropdown menu showing "1/ 1/2000"; "Latest Accident Date" with a dropdown menu showing "12/31/2009"; "Initial Random Seed" with the value "16807"; and "Frequency Correlation Copula" with a green button labeled "C" and the text "normal Correlation=c(0) Dim=2".

Picture (11), Overall Simulation Properties

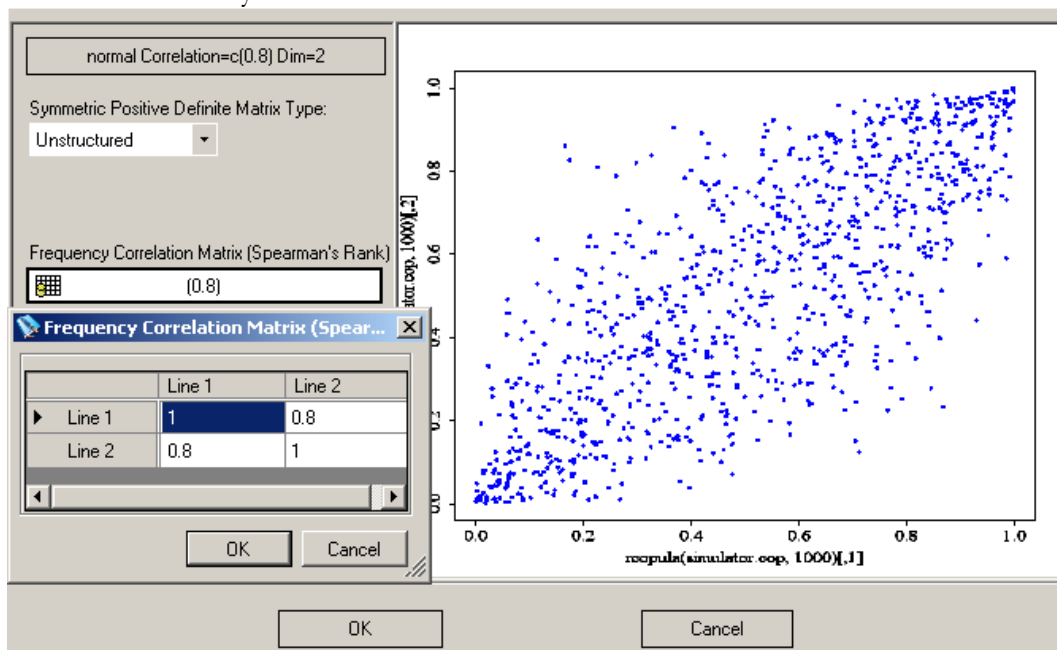
The **Initial Random Seed** (or **seed state**, or just **seed**) is a number used to initialize a pseudorandom number generator. The choice of a good random seed is crucial in the field of computer security. However, in our simulator it is mainly used for testing the simulated result. By fixing the seed, we ensure that users can generate the same result with the same seed number when the simulation is run under the same assumptions. Please note that the simulation results are still truly random generated when we permit the seed to be randomly selected.

The correlation among line level frequencies is defined by **Copula**. Please go to the help system to read more about how the copula is applied in the simulator on a multivariate distribution. Or you can read more Copula applications from R help.

You can **ignore** copula correlation with the default value set, as seen in Picture (11). The green button shows property of a 2-dimension Normal Copula, but correlation is $c(0)$. The button is enabled only for multiple lines of businesses. In our example, it has two dimensions because the sample created two lines of businesses under the simulation project,

and so there are two marginal frequency distributions to be correlated. The $c(0)$ correlation means there is no correlation between the variants. Thus when simulating, system will instead bypass the copula calculation and run the typical distribution individually defined by those variants.

When you are confident enough to apply Copula for correlations, you can click the green button above, and see the Picture (12) for a typical copula screen. You will notice a 2-dimensional copula scatter plot on the right side of the screen. In this sample, it is a **normal copula**. Please pay attention to the shape of the plot. The simulator will show copula plot up to 3 dimensions. The plot will be refreshed once you change any property. Correlation Matrix can be entered by clicking the button of “Frequency Correlation Matrix”. It is actually the Spearman’s Rank matrix for normal copula, as sampled. The matrix is symmetrical, so that you can just give upper left values, and the system will fill in the lower right corner values automatically.



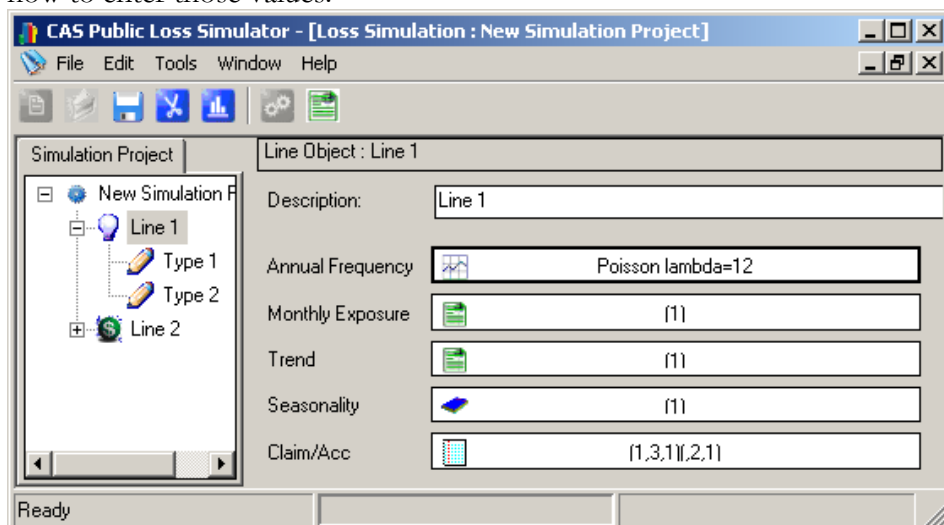
Picture (12), Define a Copula between LOB frequencies

The marginal distributions of this copula will be the line level Frequencies.

2.4 Line Level Properties


The Simulator defines most of frequency related properties at the line of business level, as shown in Picture (13). They are Annual Frequency, Monthly Exposure, Trend, Seasonality and Multinomial Claim Distribution among types. The help system has provided very detailed explanations on each property topic and how monthly frequency is generated from an Annual Frequency by these properties. The help system also illustrates the detail

workflow of: monthly frequency -> occurrences -> claims -> claims for each type. So, in this section of the paper, we won't go through them in detail again, but instead, will focus on how to enter those values.



Picture (13), Line level properties

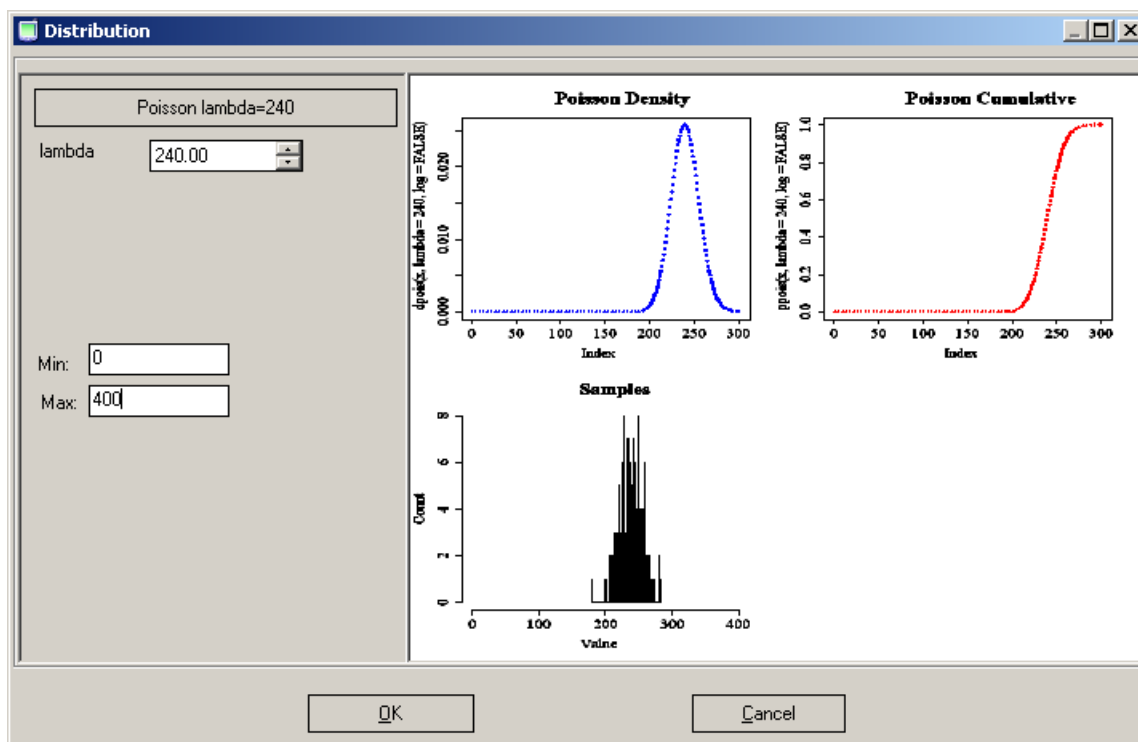
2.4.1) How to define distribution in simulator

Annual Frequency is a distribution value. The button on the right side will show the little  icon, and lists out the distribution property, such as, it is a Poisson distribution with $\lambda = 12$. This is a very small sample. You can change it by clicking the button, and you will see screen as Picture (14).

In fact, all the distributions, either continuous or discrete, are defined with screens similar to Picture (14). Any distribution defined in the simulator will have properties such as name, parameters, and min/max value. The upper left corner contains a button that shows the currently selected distribution. You can change the distribution by clicking it and picking from a dropdown list. Then you can change the parameter(s) for the selected distribution. Every time you make a change, you will notice that distribution graphs such as Density, CDF, and a Histogram are refreshed on the right side of the panel.

The **min/max values** are applied when sampling from the defined distribution, so that any generated random number will be limited to the $[\min, \max]$ range. With this said, when the associated distribution is the severity distribution, the max value also physically represents the **coverage limit**.

The Histogram shown in the picture is a visual representation, between the $[\min, \max]$ values, of the possible random numbers sampled in the future simulation run.



Picture (14), A Poisson distribution

All the distributions in the simulator are defined in the same manner as is done in the R language. For example, the user enters the meanlog (μ) and sdlog (sigma) parameters for all lognormal distributions in the model. Thus, if you want to define the size of loss to be lognormal with mean \$100,000 and standard deviation \$100,000, you need to perform a **conversion calculation** first.

To achieve the given size of loss distribution above, you can set the parameters

$$\text{sigma} = \text{sdlog} = \{\ln(1 + \text{CV}^2)\}^{.5}, \text{ where } \text{CV} = \text{mean}/(\text{standard deviation})$$

Since $\text{CV} = 100,000/100,000 = 1$ in our example, $\text{sigma} = \{\ln(2)\}^{.5} = .833$ and

$$\mu = \text{meanlog} = \ln(\text{mean}) - (\text{sigma}^2)/2 = \ln(100,000) - .693/2 = 11.166$$

Please recall that $\ln(x)$ represents the natural logarithm of x , while $\exp(x) = e^x$, the constant e raised to the power x .

This will produce a size of loss X with mean

$$EX = \exp(\mu + \frac{1}{2} \sigma^2) \approx 100,000 \text{ and second moment } E(X^2) = \exp(2\mu + 2\sigma^2) \approx 2 \times 10^{10}.$$

This results in $\text{Var}(X) = 2 \times 10^{10} - 10^{10}$ and Standard Deviation(X) = 100,000.

2.4.2) How to define simulation property by month

Most of the properties defined in simulator are **monthly** values. That means their value will

Modeling Loss Emergence and Settlement Processes

change by month and detail calculations are also carried out at monthly level. For instance, trend, exposure, seasonality, most of the severity distributions and lags, case reserves, etc, are all defined by month. Line level frequency is defined annually for simplicity, but it is still divided into monthly frequency in the initial simulation calculations.

This may cause some confusion, especially if you define severity distribution or lags by month. By default, the Simulator will set those values to all be the same value. When you change a value for one particular month, that value is automatically carried through to all subsequent months. This feature gives user a flexibility to treat all kinds of probabilities in their business, and provides a powerful ability to change the distribution for a certain variant. Please see Picture (16) for a Payment Lag Distribution example defined in next chapter.

At the line level, an example of frequency trend defined by each month in 2001 is illustrated by Picture (15). You can navigate to a different year using the blue arrows. Please be aware that every time you make a change to one cell, the values after that month will be changed automatically, and will remain changed across years. So, if you just want to change only one month's value, you should be sure to change the value for later months to the previous value.

	2001
Jan	0.9
Feb	0.9
Mar	0.9
Apr	0.8
May	0.8
Jun	0.8
Jul	0.8
Aug	0.8
Sep	0.8
Oct	0.7
Nov	0.9
Dec	0.9

Picture (15), an example of frequency trend

	2000
Jan	Exponential (rate=0.002739726)
Feb	Exponential (rate=0.002739726)
Mar	Exponential (rate=0.002739726)
Apr	Exponential (rate=0.002739726)
May	Exponential (rate=0.002739726)
Jun	Exponential (rate=0.002739726)
Jul	Exponential (rate=0.002739726)
Aug	Exponential (rate=0.002739726)
Sep	Exponential (rate=0.002739726)
Oct	Exponential (rate=0.002739726)
Nov	Exponential (rate=0.002739726)
Dec	Exponential (rate=0.002739726)

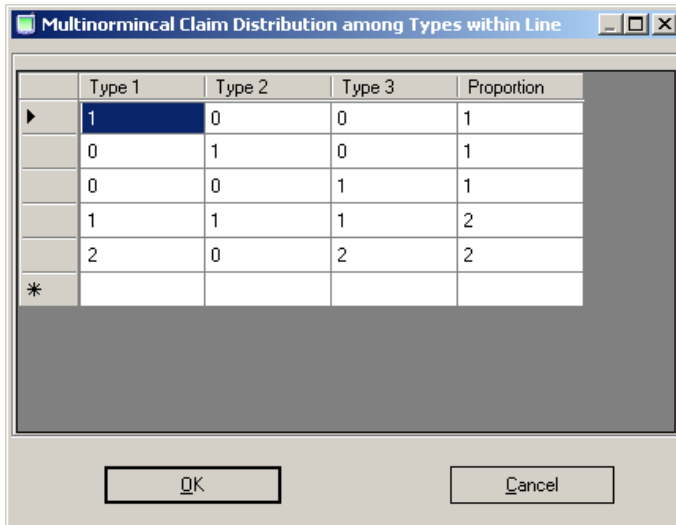
Picture (16), an example of Payment Lag monthly for each month

2.4.3) How to define claim/acc value

Modeling Loss Emergence and Settlement Processes

The number of claims of each Type arising from an occurrence is modeled via a multinomial distribution. This allows for the possibility of multiple claims from the same occurrence, either of the same Type or of multiple Types. It typically introduces correlations between frequencies across pairs of Types within the same Line.

Picture (17) demonstrates a screen example of the multinomial distribution parameters entered for a line with 3 types defined.



Picture (17), example of defining multinomial claim distribution among types

In this case, the multinomial matrix is marked with blue in the following table. For example, if you have 600 occurrences generated in some month for this line, then the Simulator will run this R command at backend to allocate the 600 occurrences to each possible combination of claim types.

```
> rmultinom(1, 600, c(1,1,1,2,2))
[1,]
[1,] 95
[2,] 85
[3,] 80
[4,] 178
[5,] 162
```

Type1	Type2	Type3	Proportion	Normalized Probability	Occurrence
-------	-------	-------	------------	------------------------	------------

Modeling Loss Emergence and Settlement Processes

1	0	0	1	0.1429	95
0	1	0	1	0.1429	85
0	0	1	1	0.1429	80
1	1	1	2	0.2857	178
2	0	2	2	0.2857	162

To interpret the above table, imagine that we have a die that is loaded so that the normalized probabilities in the above table represent the probability of landing on each of five sides of the die, while the probability of landing on the sixth side is 0. The Occurrence column summarizes the results of tossing the die 600 times. If we repeated this process 1,000 times and summed the resulting Occurrence columns from these simulations, we would expect that approximately 14.29% of the total occurrences generated would have one Type 1 claim and no Type 2 or Type 3 claims, 14.29% of the total occurrences generated would have one Type 2 claim and no Type 1 or Type 3 claims, 14.29% of the total occurrences generated would have one Type 3 claim and no Type 1 or Type 2 claims, 28.57% of the total occurrences generated would have one claim of each Type, and 28.57% of the total occurrences generated would have two Type 1 claims, no Type 2 claims and two Type 3 claims, **Claims generated from same occurrence will have the same accident date.**

The user can enter more combinations in the cells marked with * on the left, as shown in the figure. They can also just highlight a row and hit delete to remove a proportion. However, please note that if you add another row of data such as (1, 0, 0, 3), since the allocation of “1, 0, 0” already exist, the Simulator will combine them into one allocation as (1, 0, 0, 4).

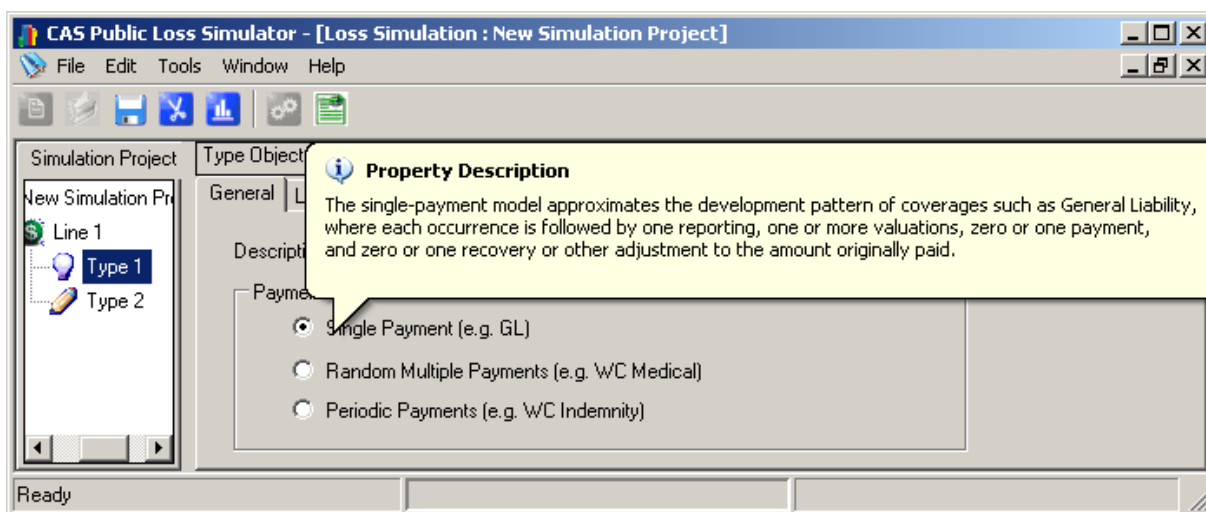
Please refer to the Loss Simulator Help section 5.4 - Multinomial Claim Distribution among Types for more technical model details on this topic.

2.5 Type Level Properties

The Type object defined in the Simulator can be used to define coverage for a line. It will contain most of the required coverage properties such as payment pattern, lags, severity, case reserve activities, and recovery adjustment properties.

2.5.1) Three Payment Patterns

This version of simulator contains three payment patterns defined in the first tab of each Type object, as sampled in Picture (18). **Loss Simulator contains another set of help system that, when you put your mouse over a property (most of them are property labels), you will see a popup tooltip help window for that property, explaining the physical meaning of that property.** Each payment pattern option has its own completely different simulation algorithm. Please refer to the help section online at: <http://www.casact.org/research/lsmwp/losshelp/index.cfm?fa=main> for a detailed discussion of each payment pattern option.



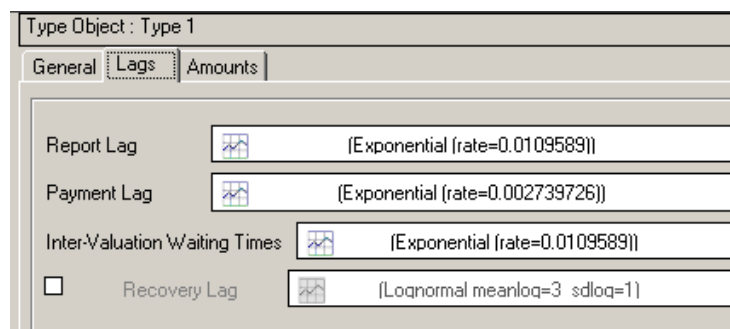
Picture (18), three payment patterns defined for each Type object

Once you click one of the payment pattern radio buttons, you will notice that the remaining two tabs of the Type object will also be changed associated with the selected payment pattern. We will explain them one by one, associating each property with each simulation algorithm.

2.5.2) Single Payment Pattern

A Type object of Single Payment Pattern will have the following two screens, as seen in Picture (19) and Picture (21), to define lags, severity properties, and recoveries.

2.5.2.1) Single Payment Pattern Lags



Picture (19), lags defined for single Payment Type object

Report Lag, a monthly defined continuous distribution, helps to generate Report Date of a claim from the accident date.

Payment Lag, a monthly defined continuous distribution, helps to generate Payment Date

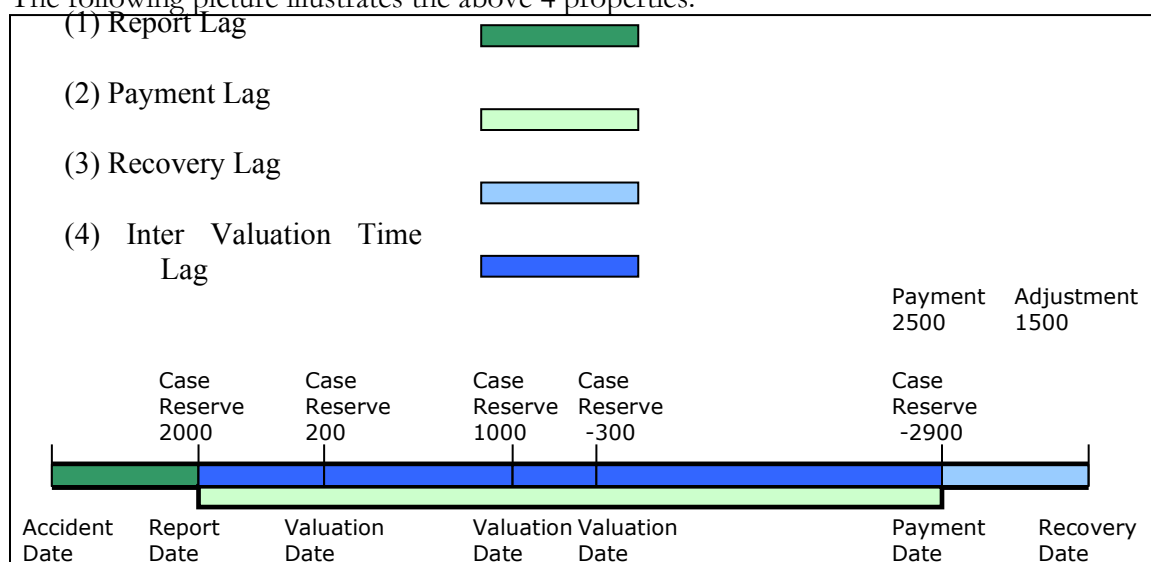
Modeling Loss Emergence and Settlement Processes

of the claim from the Report Date.

Inter-Valuation Waiting Times, a monthly defined continuous distribution, will define each Valuation Date (Case Reserve Date) between the Report Date and the final Payment Date.

Recovery Lag, a monthly defined continuous distribution, will define the Recovery Date after the Payment Date if the checkbox is selected.

The following picture illustrates the above 4 properties.



Picture (20), Illustration of Single Payment Pattern option in simulation

2.5.2.2) Single Payment Pattern Severity Properties

The following screen Picture (21) shows the properties to define severity, case reserve, and recovery (if selected). We will explain each property in detail.

The screenshot shows the 'Amounts' tab of a software interface. The 'Type Object : Type 1' is selected. The 'Amounts' tab is active, showing various input fields for defining severity and recovery properties. The fields are organized into sections:

- Size of Entire Loss:** (Lognormal meanloq=11.16636357 sdloq=0.832549779)
- Correlation with Payment Lag:** C (normal Correlation=c(0) Dim=2)
- Trend:** (1) Alpha: 0
- Deductible:** (0) P(0): (0.4)
- Case Reserve Adequacy:** (Lognormal meanloq=-0.105360516 sdloq=0.05)
- 40% Case Reserve Adequacy:** (Lognormal meanloq=0.5 sdloq=0.05)
- 70% Case Reserve Adequacy:** (Lognormal meanloq=0.25 sdloq=0.03)
- 90% Case Reserve Adequacy:** (Lognormal meanloq=0.05 sdloq=0.03)
- Est P(0):** (0.4) **Threshold:** (0)
- Minimum Change:** (100) **Min Rel Change:** (0.01)
- Inertia:** (0.2) **Fast Track:** (2000)
- Initial Payment Adequacy:** (Lognormal meanloq=0.5 sdloq=0.05)
- P(1):** (1)

Picture (21), Severity and recovery (if enabled) for single Payment Type object

Size of Entire Loss is a monthly defined continuous distribution that defines the severity distribution of the Type. It will determine the final claim payment. This variable can be correlated using a Copula object with the payment lag defined in the Type. However, NO correlation is selected in Picture (21).

P(0), a monthly defined value, defines the constant probability of closure without payment, for reasons other than failure to exceed the deductible.

The cumulative trend factors (“cum”) are calculated to the accident and payment dates, and then the trend multiplier is calculated as follows:

$$\text{trend} = (CUM_{acc_date}) \left(\frac{CUM_{pmt_date}}{CUM_{acc_date}} \right)^\alpha = (CUM_{acc_date})^{1-\alpha} (CUM_{pmt_date})^\alpha$$

We are applying the full trend factor from the average accident date underlying the assumed loss distribution to the accident date of the occurrence, and a portion “alpha” of the trend from the accident date to the payment date. This is an application of Butsic’s alpha parameter concept from his May 1981 CAS Discussion Paper Program entitled “The Effect of Inflation on Losses and Premiums for Property-Liability Insurers. “

Then the final payment of the claim is calculated by multiplying the trend factor and the simulated claim from the Size of Loss property and then applying the policy limit and deductible.

Case Reserve Adequacy

Case Reserve Adequacy at valuation time t is a lognormal random variable with

$\mu = (\text{meanlog at time } t)$ and $s = (\text{sdlog at time } t)$. Here t is the fraction:

$$t = \frac{\text{valuation time minus report date}}{\text{payment date minus report date}} = \frac{\text{Valuation lag}}{\text{Payment lag}}$$

The user enters the *meanlog* for times $t = 0\%$, 40% , 70% , and 90% (the time 0 value is labeled “case reserve adequacy” in the model). The *meanlog* for time 1.0 is set at 0.0. For other values of t , the Simulator applies linear interpolation to calculate values of *meanlog*. The modeler also inputted *sdlog* s for times 0% , 40% , 70% , and 90% , which provides flexibility in controlling the variance of the reserve adequacy factor.

On the Simulator screen, the user determines the case reserve adequacy parameters by selecting four lognormal distributions labeled as Case Reserve Adequacy, 40% Case Reserve Adequacy, 70% Case Reserve Adequacy, and 90% Case Reserve Adequacy, as displayed in Picture (21).

Threshold and EstP (0)

The Simulator introduces a non-negative threshold value that is associated with EstP(0) and is used when case reserves are set at each valuation time. If the claim’s ultimate size of loss value is strictly below this threshold, we would not apply the EstP(0) adjustment in setting the case reserve. If the user enters a zero value for threshold, then the model will apply the "EstP(0)", adjustment to all claims.

The algorithm for estimating the case reserve value of a claim is as follows:

- (1) *If threshold <= 0 OrElse payment > threshold Then*
value = (1 - EstPO at ValuationDate) payment * multiplier(apply policy limit and deductible)*
- (2) *If (threshold > 0 AndAlso payment <= threshold) Then*

Modeling Loss Emergence and Settlement Processes

$$value = payment * multiplier - deductible$$

(Multiplier is the value simulated from the interpolated case reserve adequacy lognormal distribution.)

Recovery. Recovery is optional.

By checking the Recovery Lag Checkbox, as shown in Picture (19), you enable the recovery calculation. **P(1)** is the probability that the claim closes with the initial payment amount, and is constant for all claims. The default value is 1, which means the initial payment covers the full payment amount and there is no recovery (even if you selected the checkbox).

InitialPaymentAdequacy is a monthly defined distribution that is used when P(1) is not 1. It represents the ratio of the ultimate payment after recoveries to the initial payment, and so defines the recovery amount as $Recovery = Initial\ Payment * (1 - ratio)$.

2.5.3) Multiple Random Payment Pattern

The multiple random payments model approximates the development pattern of coverage such as Medical Payments, where each occurrence is followed by a random number of reimbursable incurred expenses. Each expense is followed by one reporting and one payment. The final expense payment (which may or may not be the final expense incurred) is followed by zero or one recovery or other adjustment to the total of all previous payments.







The following pictures will show the screen of defining lags and severities properties of a multiple random payment type.

The screenshot shows a software interface with three tabs: 'General', 'Lags', and 'Amounts'. The 'Lags' tab is active. It contains a list of properties with their corresponding distributions:

Property	Distribution
Number of Separate Expenses	(Geometric prob=0.1)
Lag Between Expense Incurals	(Exponential (rate=0.0109589))
Expense Report Lag	(Exponential (rate=0.0109589))
Expense Payment Lag	(Exponential (rate=0.002739726))
<input type="checkbox"/> Recovery Lag	(Lognormal meanloq=3 sdloq=1)

Picture (22), lags defined for Multiple Random Payment Type object

Modeling Loss Emergence and Settlement Processes

General	Lags	Amounts
Size of Each Expense		(Lognormal meanlog=11.16636357 sdlog=0.832549779)
Trend		{1}
Decay		{1}
Deductible		{0}
Case Reserve Adequacy		(Lognormal meanlog=-0.105360516 sdlog=0.05)
Minimum Ratio		{0.01}

Picture (23), Size of Loss properties defined for Multiple Random Payment Type

The multiple-random-payments model approximates the development pattern of coverages such as Medical Payments, where each occurrence is followed by a random number of reimbursable incurred expenses. Each expense is followed by one reporting and one payment. The final expense payment (which may or may not be the final expense incurred) is followed by zero or one recovery or other adjustment to the total of all previous payments.

The number of reimbursable expenses per claim is assumed to follow a geometric distribution or a multinomial distribution specified by the user.

The severity parameters describe the distribution of severities for each individual expense, except that the deductible and maximum apply to all expenses in aggregate. There is a single trend factor (actually annual trend rate) applied to each expense through its incurral date, and there is decay factor allowing the user to specify a declining mean from one expense to the next expense arising from the same claim.

Case reserves are assumed to be revalued at each payment date. Their adequacy is measured relative to all expenses that have been or will be incurred but have not yet been paid, subject to a minimum that allows a reserve to be carried between the last payment and the recovery date. The "P(2 sig dig)" entry represents the probability that a case reserve will be estimated to a nearby "round" number -- in this case rounding to two significant digits -- rather than its exact value.

Recoveries are modeled as one-time adjustments to correct errors in the original amounts paid. For this purpose each amount paid is treated as an adequacy factor times the actual severity after application of the deductible and the maximum. Payment errors are reflected in the distribution of this factor less 1.00. In particular, if the adequacy factor is greater than 1.00, the initial payment will be too great and will produce a future recovery, represented as a negative payment. The simulator generates both the original overpayment and underpayment, spread uniformly across all payment dates, and the later recovery or adjustment.

2.5.4) Periodic Multiple Payment Pattern

The periodic payments model approximates the development patterns of coverage such as Group Long-Term Disability or the wage-replacement provisions of Workers' Compensation, where each occurrence is followed by a random number of regular periodic payments of equal amounts or of equal amounts subject to periodic inflation adjustments. The final payment is followed by zero or one recovery of any payments that were inadvertently made beyond the termination of disability.

Report Lag	(Exponential (rate=0.01095891))
Payment Frequency	<input type="radio"/> Weekly <input type="radio"/> Biweekly <input checked="" type="radio"/> Monthly
Payment Duration in Years	(Exponential (rate=0.02739726))
Initial Payment Lag	(Lognormal meanlog=3.5 sdlog=0.25)
<input type="checkbox"/> Recovery Lag	(Lognormal meanlog=3 sdlog=1)

Picture (24), lags defined for Multiple Periodic Payment Type

Size of Entire Loss	(Lognormal meanlog=13.16636357 sdlog=0.832549779)
Trend	(1.01)
COLA	(1.02)
Tabular Reserves	
CP Factor	(0.95)
Discount	(0.05)

Picture (25), severity and case reserve defined for Multiple Periodic Payment Type


Size of Loss, is typically related to salaries and may be approximated by a suitable distribution with a minimum and maximum. There is provision for trend which affects the payment size at time of occurrence, and for COLA factors which affects individual claims at annual intervals following commencement of payments.

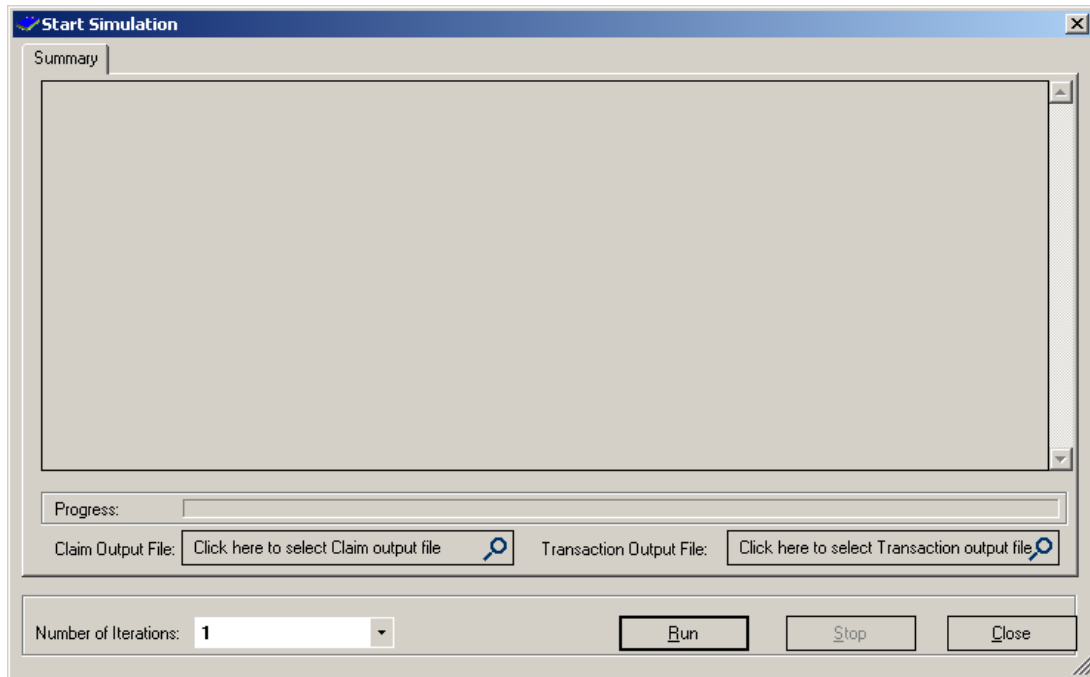
Trend, is the change in mean payment size for newly incurred claims. It is expressed as an annual factor but applied monthly.

COLA, is the change in payment size, at the end of the first and each subsequent year, for a given claim already in payment status.

CP factor, is the ratio of the periodic continuance probabilities assumed in the tabular reserves to the same probabilities implicit in the payment duration distribution.

2.6 Run Simulation

You can only start the simulation with an open project; otherwise, the simulation menu and the toolbar buttons are all disabled. Once you feel comfortable enough for the model configuration, you can go to Tools->Run Simulation, or simply click the  button from toolbar. The “Run Simulation” window pops up as shown in Picture (26)



Picture (26), Run simulation

2.6.1) Provide Claim and Transaction Output File

The Simulator will need a claim and a transaction output file. So from the above screen, you can click the two buttons individually and provide the **CSV** file name from the Windows File Save Dialog. The Simulator will then generate two **CSV** files after the simulation is done.

Here is an example of how the claim CSV file looks like (if opened with Excel),

Simulation 6/9/2010 11:51:33 AM

Simulation No	Occurrence No	Claim No	Accident Date	Report Date	Line	Type
1	1	1	20000117	2000110	3	1

Modeling Loss Emergence and Settlement Processes

				2000041		
1	2	1	20000112	8	1	1
				2000032		
1	3	1	20000108	0	1	1
				2000020		
1	4	1	20000129	1	1	1

And here is a sample transaction file,

Transactions 6/9/2010 11:51:33 AM

Simulation No	Occurrence No	Claim No	Date	Transaction	Case Reserve	Payment
			2000110			
1	1	1	3	REP	2000	0
			2001031			
1	1	1	3	RES	15439	0
			2001042			
1	1	1	5	RES	1330	0
			2001091			
1	1	1	5	RES	2938	0
			2002042			
1	1	1	9	RES	-3870	0
			2002050			
1	1	1	1	RES	-650	0
			2002052			
1	1	1	4	RES	-484	0
			2002060			
1	1	1	2	RES	-198	0
			2002070			
1	1	1	6	RES	-632	0
			2002082			
1	1	1	1	RES	-920	0
			2002090			
1	1	1	6	RES	-542	0
			2002100			
1	1	1	4	CLS	-14411	1233
			2000041			
1	2	1	8	REP	2000	0
			2000091			
1	2	1	8	RES	13113	0
1	2	1	2001040	RES	3830	0

Modeling Loss Emergence and Settlement Processes

			5				
			2001091				
1	2	1	3	RES	-1334	0	
			2002010				
1	2	1	1	RES	-2803	0	
			2002020				
1	2	1	5	RES	-950	0	
			2002031				
1	2	1	9	CLS	-13856	0	
			2000032				
1	3	1	0	REP	2000	0	
			2000040				
1	3	1	8	RES	29554	0	
			2000051				
1	3	1	1	RES	5778	0	
			2000062				
1	3	1	4	RES	7127	0	
			2000100				
1	3	1	4	RES	-984	0	
			2001010				
1	3	1	8	CLS	-43475	10866	

2.6.2) Number of Iterations

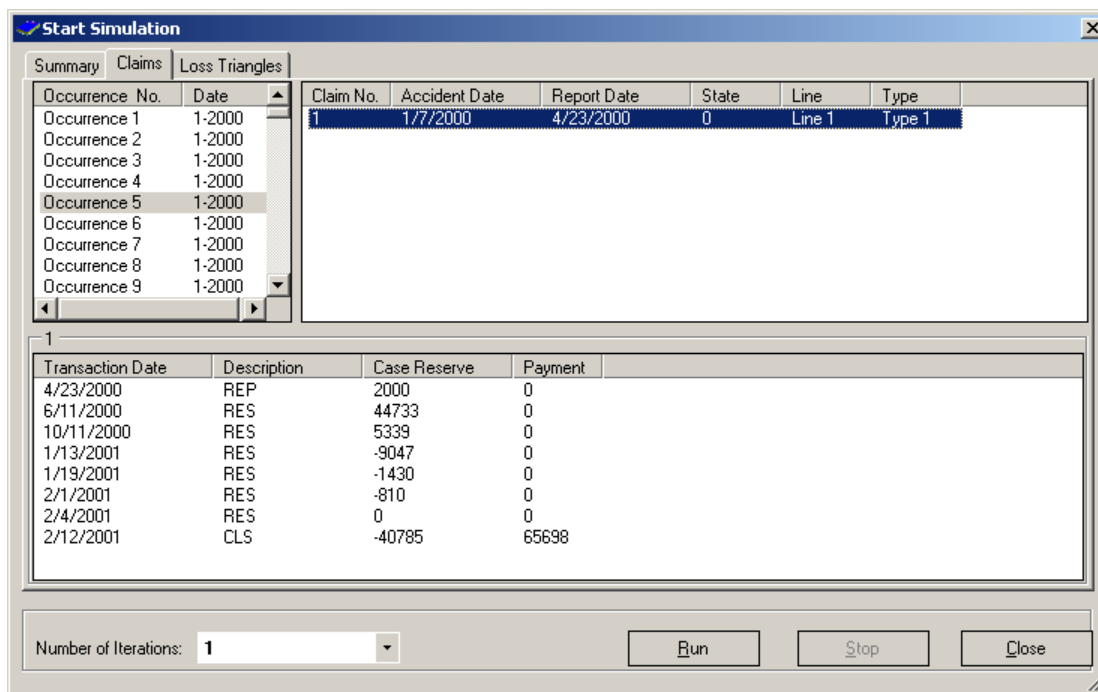
This controls how many iterations the Simulator will run. To get a feel for the simulator, you can simply choose 1 from the dropdown selection, and the simulation will finish after 1 iteration.

For real simulation, you can choose up to 1000 iterations. This could require many hours of CPU time, depending upon the complexity of your parameterization of the model and your computer speed. The final result may contain hundreds of thousands of claims and millions of transactions. In this case, you can let the program run overnight. However, you can cancel the process at any time by clicking the “Stop” button.

2.6.3) On Screen Simulation Output

The Run Simulation window will show the calculation progress, and at regular intervals, it also displays the simulation summary text. For a large simulation, the summary text will be refreshed after each 100 iterations of the simulation process, so it may not give you the full picture in this case.

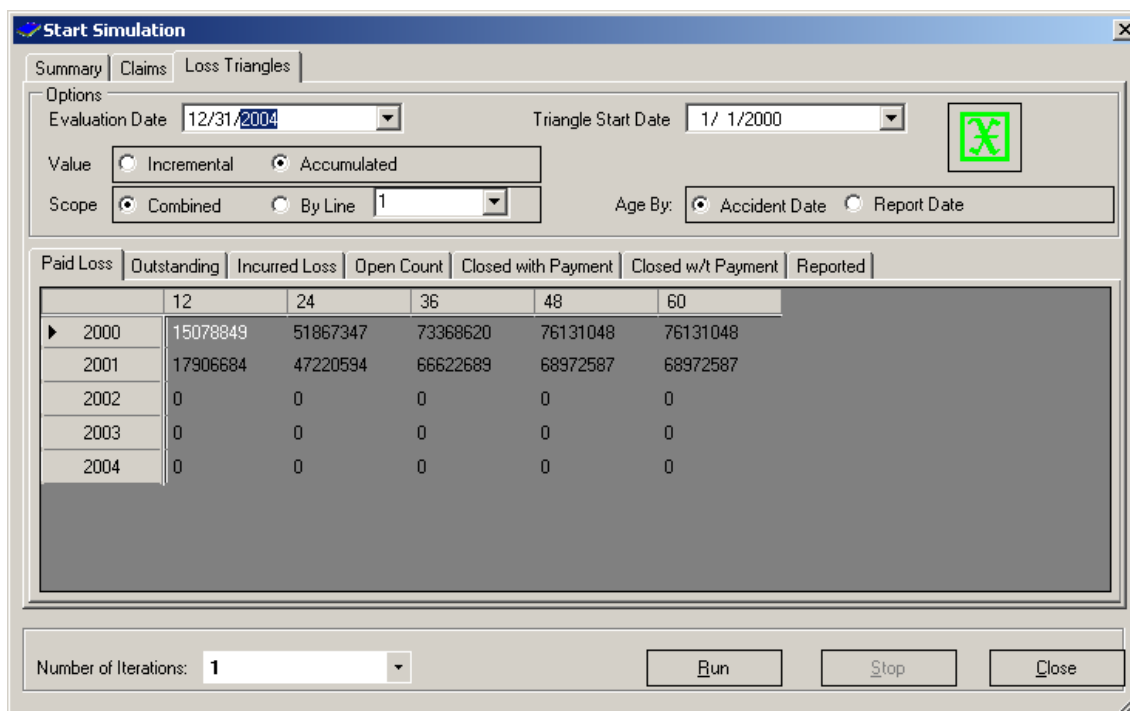
Modeling Loss Emergence and Settlement Processes



Picture (28), Occurrence, Claim and Transactions (if only 1 iteration is executed)

If you run the Simulator with one iteration, the program will also generate a variety of loss triangles as shown in Picture (29). The user may export the triangles to an Excel file. Please note that the “triangles” are presented in “rectangle” format, and that the user may also re-configure the triangles by choosing different properties on this screen.

Modeling Loss Emergence and Settlement Processes

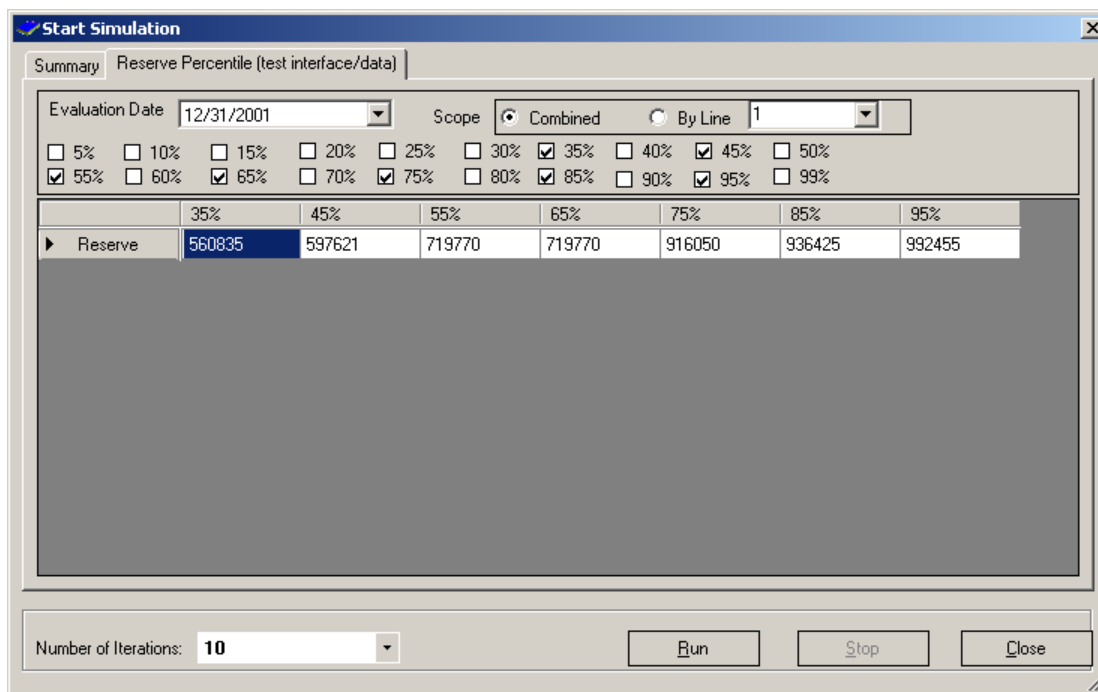


Picture (29), Loss Triangles from the output (if only 1 iteration is executed)

2.6.4) Reserve Percentile Output

If the simulation has more than 100 iterations, the Simulator will generate a Reserve Percentile table from the simulation results.

Modeling Loss Emergence and Settlement Processes



Picture (30), Reserve Percentiles from the output (if 100+ iteration is executed)

3) Simulation Example

The CAS Loss Simulator Working Party Testing Group, led by Professor Joe Marker, has conducted several testing scenarios to evaluate the simulation parameterization and simulation results. One scenario tests the severity distributions (i.e. size of loss). We select this test in order to illustrate how to set up the parameters to simulate a company's business.

The simulated results contains

File Name	File Description
c.csv	claims output by the simulator
t.csv	transactions output by the simulator
ultloss 20100512.csv	the file containing ultimate loss for each claim
Test severities 20100512.xml	The XML file containing the parameters used to run the model. This file can be imported into your Simulator directly to see the following parameter definitions, and will produce the same results as displayed in the above c.csv and t.csv files.

These files can be downloaded from the following links:

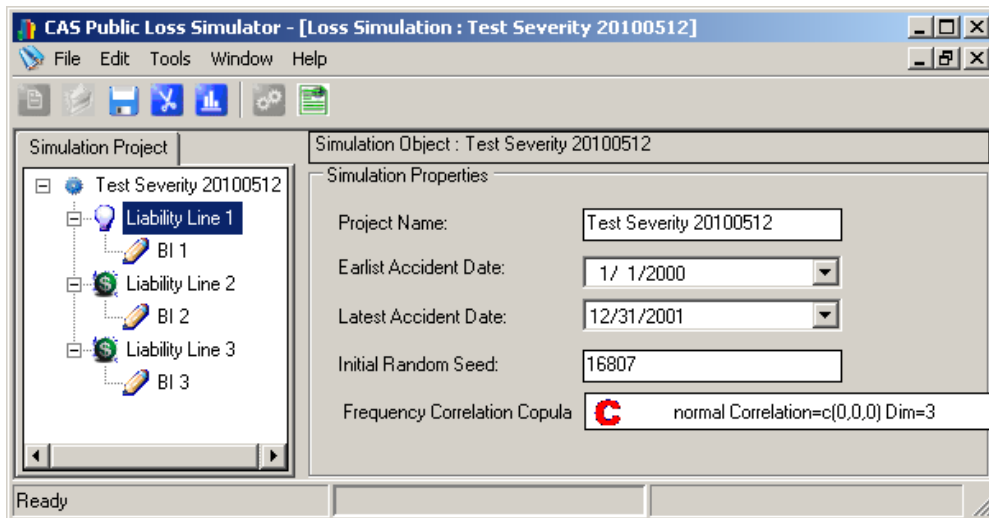
http://www.gououon.com/loss_simulator/file/simulation20100520

3.1 Test Objective: Test univariate ultimate severities

3.2 Simulation Project Level Parameters Setup.

Set up three lines of business with no correlation in frequency among the three lines. If you import the “Test severities 20100512.xml” file, you will see the exact same screens illustrated below, starting from Picture (31).

Project Name	Test Severity 20100512
Accident Years	2000-2001
Initial Random Seed	16807
Frequency correlation copula	(None) normal Correlation=c() Dim = 1
Line of Business (3)	Liability Line 1 Liability Line 2 Liability Line 3
Coverage	One BI Coverage (for each line), single payment model.

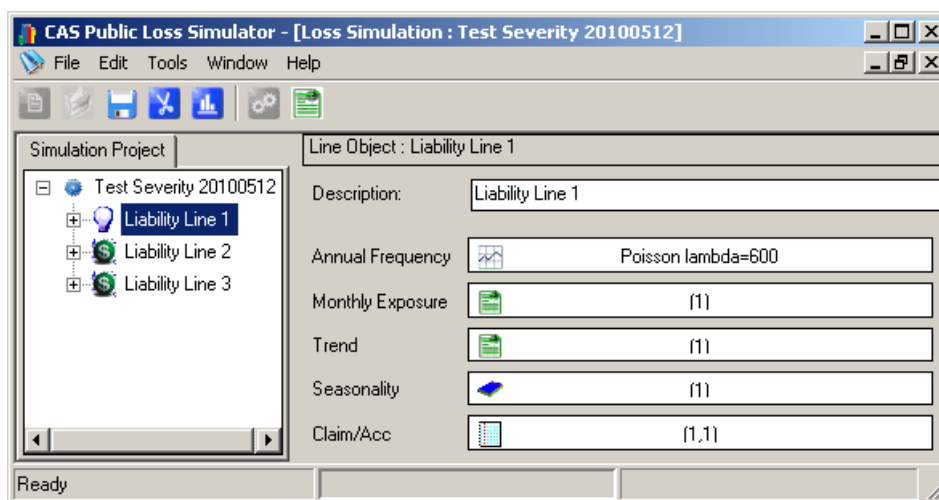


Picture (31), Severity Test with Project Level Parameters Setup

3.3 Line Level Parameters Setup.

For each line, an occurrence generates one claim of one unique type for each line. Zero trend for frequency, zero trend for severity. Most parameters are the same for all lines. We use Line 1 as representative of all three lines in the table below (please also refer to Picture (32)):

Annual Frequency	Poisson(600)									
Monthly Exposure	(1)—None									
Trend	(1)—None									
Seasonality	(1)—None									
Claim/acc matrix	each occurrence generates 1 claim of one type <table border="1" style="margin-left: 20px;"> <thead> <tr> <th></th> <th>Bl 1</th> <th>Proportion</th> </tr> </thead> <tbody> <tr> <td>▶</td> <td>1</td> <td>1</td> </tr> <tr> <td>*</td> <td></td> <td></td> </tr> </tbody> </table>		Bl 1	Proportion	▶	1	1	*		
	Bl 1	Proportion								
▶	1	1								
*										
Expected # claims	= 600 (freq) x 100 (# sims) x 3 (lines) x 2 (years) = 360,000.									
Actual # claims from the simulation result	359,819									



Picture (32) Line Level Parameters

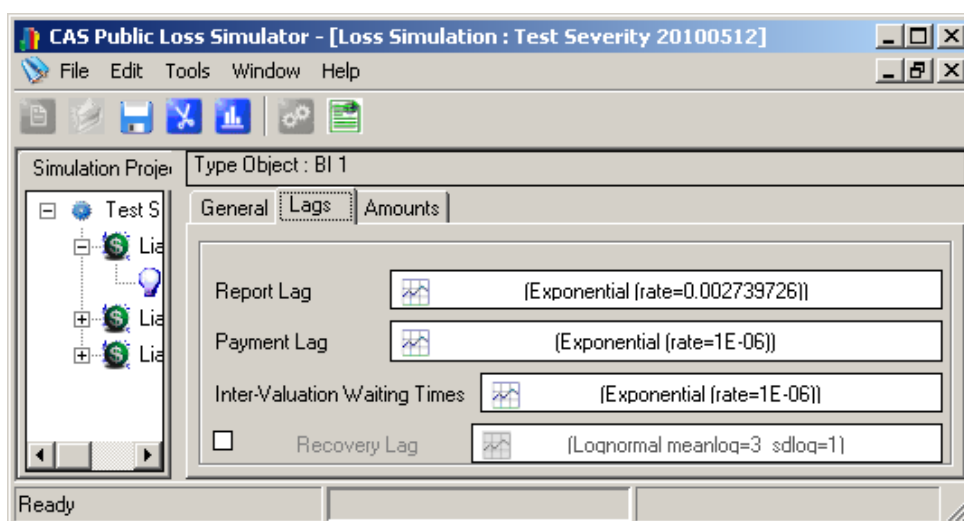
3.4 Coverage Level Parameters Setup.

Lags are irrelevant for this run except for report lag. We are not testing the reserve change process.

Report lag	Exponential with Rate = 1/365, mean=365 days. Max=3650.
Payment lag	Maximum one day (irrelevant)
Inter-valuation lag	Maximum one day (irrelevant)
Recovery lag box	Not checked (irrelevant)

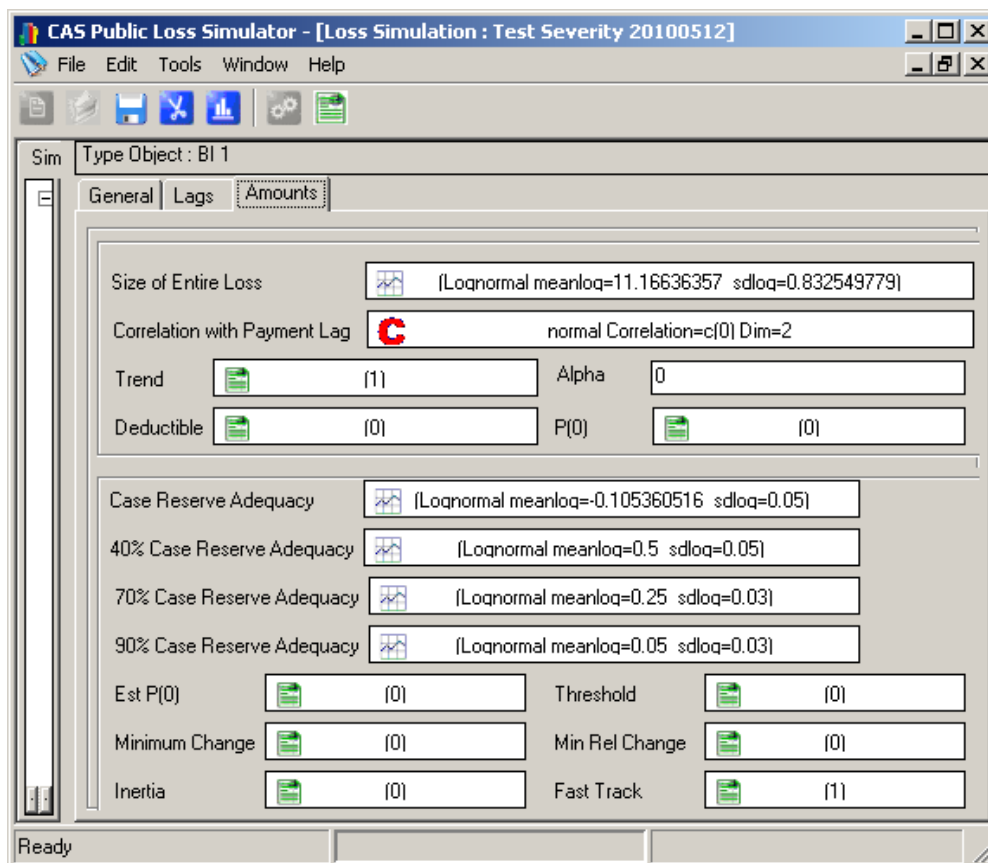
Modeling Loss Emergence and Settlement Processes

Correlation of Amount with lag	(None): normal Correlation=c() Dim=2
Alpha	0
Inertia	0
P(0)	0
EstP(0)	0
Case Reserve adequacy	Irrelevant, leave at default values.
Case Reserve Interpolation	Irrelevant, leave with default
Minimum change	0
Min Rel Chg	0
P(1)	Not set because recovery lag was not checked.
Initial payment adequacy	Not set because recovery lag was not checked.



Picture (33) Coverage Level Lag Parameters

Modeling Loss Emergence and Settlement Processes



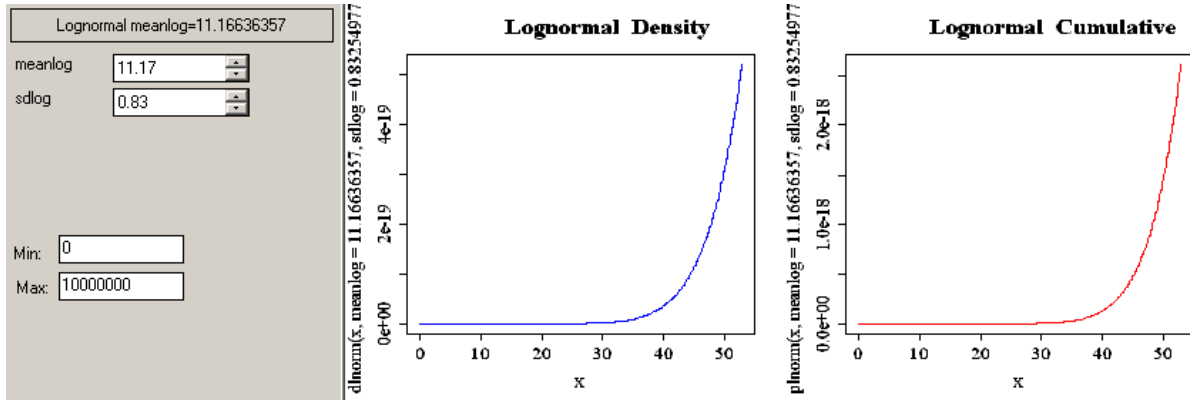
Picture (33) Coverage Level Amounts Parameters

3.5 Coverage Level Severity Distributions.

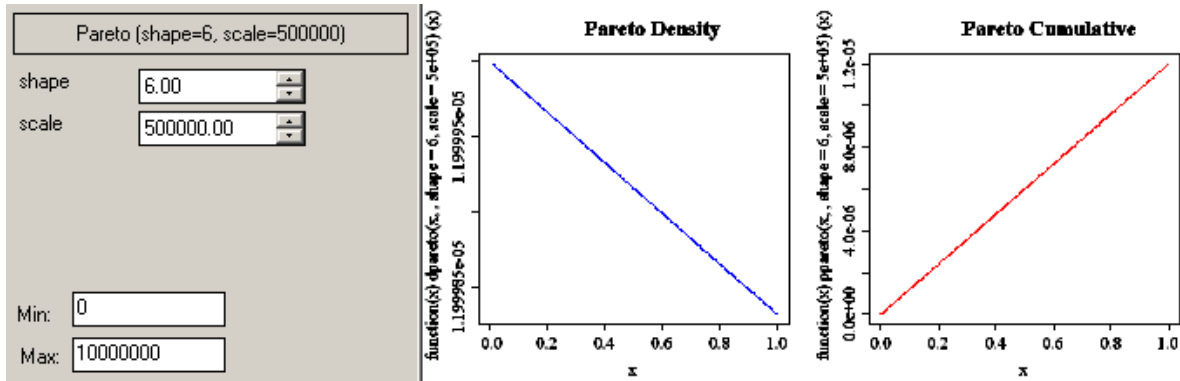
Size of Loss (severity) Parameters vary by Line

BI 1 for Liability Line 1 Picture (34)	Lognormal mean=100,000 standard. dev.=100,000 max =10,000,000	This means that the input lognormal parameters are $\mu = \text{meanlog} = 11.16636357$, $\sigma = \text{sdmean} = 0.832549779$.
BI 2 for Liability Line 2 Picture (35)	Pareto α (shape) =6 θ (scale) = 500,000 max =10,000,000	This results in a mean of 100,000 and stnd deviation of 122,474.5.
BI 3 for Liability Line 3 Picture (36)	Weibull θ (scale) = 95,000 τ (shape) = 0.9 max =10,000,000	This results in mean of 99,957 and stnd.dev. of 111,256.

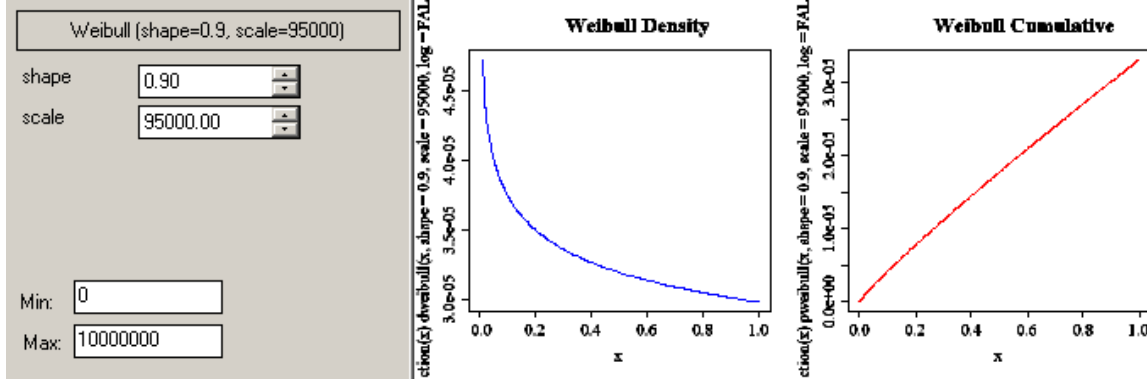
Modeling Loss Emergence and Settlement Processes



Picture (34) BI 1 for Line 1, Lognormal, mean=100,000 standard. dev.=100,000, max =10,000,000



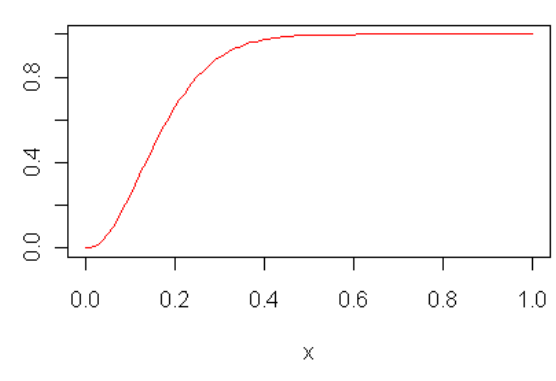
Picture (35) BI 2 for Line 2, Pareto, α (shape) =6, θ (scale) = 500,000, max =10,000,000



Picture (36), BI 3 for Line 3, Weibull, θ (scale) = 95,000, τ (shape) = 0.9, max =10,000,000

3.6 Number of Iterations: 100.

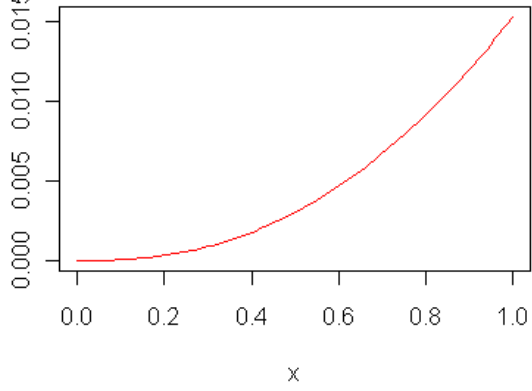
4) Distributions Used in the Public Loss Simulator and Their Parameterizations

Name	R Representation and Explanation	
Beta	Density	dbeta(x, shape1, shape2, ncp = 0, log = FALSE)
	Probability Function	pbeta(q, shape1, shape2, ncp = 0, lower.tail = TRUE, log.p = FALSE)
	Quantile Funtion	qbeta(p, shape1, shape2, ncp = 0, lower.tail = TRUE, log.p = FALSE)
	Random Generation	rbeta(n, shape1, shape2, ncp = 0)
	Details	<p>Before defining the Beta distribution, it is convenient to define the Beta function by:</p> $B(a, b) = \int_0^1 t^{a-1} (1-t)^{b-1} dt \text{ for } a>0, b>0. \text{ It turns out that } B(a, b) = \frac{\Gamma(a) \Gamma(b)}{\Gamma(a+b)} .$ <p>Then $f(x) = \frac{1}{B(a, b)} x^{a-1} (1-x)^{b-1}$ for $a > 0, b > 0$ and $0 \leq x \leq 1$, is the density for the Beta distribution with shape1 = a and shape2 = b. The boundary values at $x=0$ and $x=1$ are defined by continuity (as limits).</p> <p>The mean is $\frac{a}{a+b}$ and the variance is $\frac{ab}{(a+b)^2(a+b+1)}$.</p> <p>Pbeta is the cumulative distribution function $F(x)$ and is closely related to the incomplete beta function, defined by $B(x; a, b) = \int_0^x t^{a-1} (1-t)^{b-1} dt$.</p> $pbeta(x, a, b) = F(x) = \frac{B(x; a, b)}{B(a, b)}$
<code>plot(function(x) pbeta(x, 1.25,12.0, log=FALSE), main = "Beta CDF", col="red")</code>	<p>CDF Sample: <code>plot(function(x) pbeta(x, 1.25,12.0, log=FALSE), main = "Beta CDF", col="red")</code></p> 	

Modeling Loss Emergence and Settlement Processes

Name	R Representation and Explanation		
Exponential	Density	dexp(x, rate = 1, log = FALSE)	
	Probability Function	pexp(q, rate = 1, lower.tail = TRUE, log.p = FALSE)	
	Quantile Function	qexp(p, rate = 1, lower.tail = TRUE, log.p = FALSE)	
	Random Generation	rexp(n, rate = 1)	
	Details	<p>If rate is not specified, it assumes the default value of 1.</p> <p>The exponential distribution with rate λ has density</p> $f(x) = \lambda e^{-\lambda x} \text{ for } x \geq 0.$	
	<p>CDF Sample: <code>plot(function(x) pexp(x, rate=1.2, log=FALSE), main = "Exponential CDF", col="red")</code></p> <div style="text-align: center;"> <p>Exponential CDF</p> </div>		

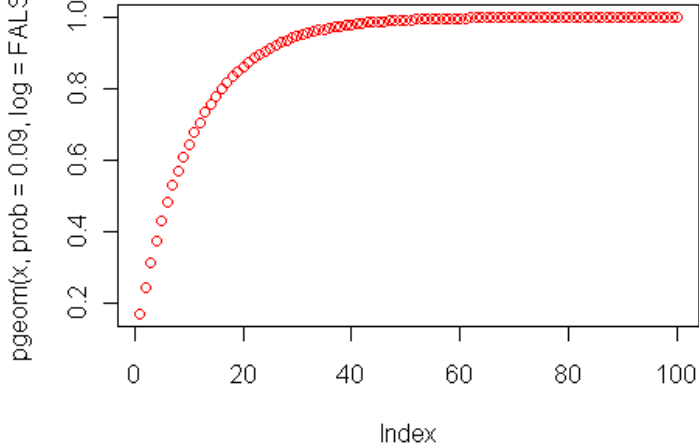
Modeling Loss Emergence and Settlement Processes

Name	R Representation and Explanation	
Gamma	Density	dgamma(x, shape, rate = 1, scale = 1/rate, log = FALSE)
	Probability Function	pgamma(q, shape, rate = 1, scale = 1/rate, lower.tail = TRUE, log.p = FALSE)
	Quantile Function	qgamma(p, shape, rate = 1, scale = 1/rate, lower.tail = TRUE, log.p = FALSE)
	Random Generation	rgamma(n, shape, rate = 1, scale = 1/rate)
	Details	<p>The Gamma distribution with parameters shape = a and scale = s has density</p> $f(x) = \frac{x^{a-1} \exp(-x/s)}{s^a \Gamma(a)}, \text{ for } x \geq 0, a > 0 \text{ and } s > 0.$ <p>(Here $\Gamma(a)$ is the function implemented by R's gamma() and defined in its help. Note that $a=0$ corresponds to the trivial distribution with all mass at point 0.)</p> <p>The mean and variance are $E(X) = a*s$ and $Var(X) = a*s^2$.</p> <p>The cumulative hazard $H(t) = -\log(1 - F(t))$ is pgamma(t, ..., lower = FALSE, log = TRUE).</p> <p>Note that for small values of shape (and moderate scale) a large part of the mass of the Gamma distribution is on values of x so near zero that they will be represented as zero in computer arithmetic. So rgamma can well return values which will be represented as zero. (This will also happen for very large values of scale since the actual generation is done for scale=1.)</p>
<p>CDF Sample: plot(function(x) pgamma(x, shape=2.5, scale==3, log=FALSE), main = "Gamma CDF", col="red")</p> <div style="text-align: center;"> <p>Gamma CDF</p>  </div>		

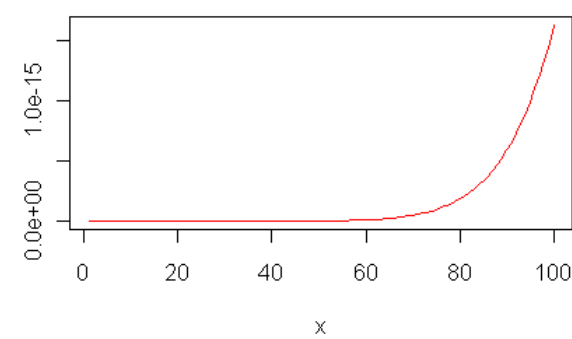
Modeling Loss Emergence and Settlement Processes

Name	R Representation and Explanation	
Gaussian (Normal)	Density	dnorm(x, mean = 0, sd = 1, log = FALSE)
	Probability Function	pnorm(q, mean = 0, sd = 1, lower.tail = TRUE, log.p = FALSE)
	Quantile Function	qnorm(p, mean = 0, sd = 1, lower.tail = TRUE, log.p = FALSE)
	Random Generation	rnorm(n, mean = 0, sd = 1)
	Details	<p>The normal distribution has density</p> $f(x) = \frac{1}{\sigma\sqrt{2\pi}} \exp(-z^2 / 2), \text{ where } z = \frac{x - \mu}{\sigma}$ <p>where μ is the mean of the distribution and σ the standard deviation.</p>
<p>CDF Sample: <code>plot(function(x) pnorm(x), -5, 5, main = "Normal CDF", col="red")</code></p> <div style="text-align: center;"> <p>Normal CDF</p> </div>		

Modeling Loss Emergence and Settlement Processes

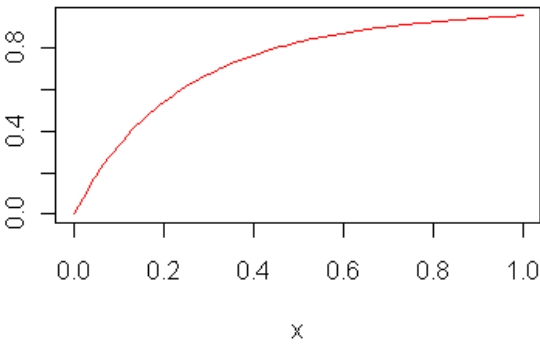
Name	R Representation and Explanation	
Geometric	Density	dgeom(x, prob, log = FALSE)
	Probability Function	pgeom(q, prob, lower.tail = TRUE, log.p = FALSE)
	Quantile Function	qgeom(p, prob, lower.tail = TRUE, log.p = FALSE)
	Random Generation	rgeom(n, prob)
	Details	<p>The geometric distribution with prob = p has density</p> $p(x) = p(1-p)^x$ <p>for $x = 0, 1, 2, \dots, 0 < p \leq 1$.</p> <p>If an element of x is not integer, the result of pgeom is zero, with a warning.</p>
<p>CDF Sample: x<-1:100 plot(pgeom(x, prob=0.09, log=FALSE), main = "Geometric CDF", col="red")</p> <div style="text-align: center;"> <p>Geometric CDF</p>  </div>		

Modeling Loss Emergence and Settlement Processes

Name	R Representation and Explanation	
Lognormal	Density	dlnorm(x, meanlog = 0, sdlog = 1, log = FALSE)
	Probability Function	plnorm(q, meanlog = 0, sdlog = 1, lower.tail = TRUE, log.p = FALSE)
	Quantile Funtion	qlnorm(p, meanlog = 0, sdlog = 1, lower.tail = TRUE, log.p = FALSE)
	Random Generation	rlnorm(n, meanlog = 0, sdlog = 1)
	Details	<p>A lognormal random variable X is one for which $Y = \ln(X)$ is normally distributed. The parameters μ and σ are the mean and standard deviation of $\ln(X)$. We can think of X as $\exp(Y)$, where Y has the normal distribution. X has density</p> $f(x) = \frac{1}{x\sigma\sqrt{2\pi}} \exp(-z^2 / 2), \text{ where } z = \frac{\ln x - \mu}{\sigma}$ <p>The mean is $EX = \exp(\mu + 1/2 \sigma^2)$, the median is $med(X) = \exp(\mu)$, and the variance $Var(X) = \exp(2*\mu + \sigma^2)*(\exp(\sigma^2) - 1)$ and hence the coefficient of variation is $\sqrt{\exp(\sigma^2) - 1}$ which is approximately σ when that is small (e.g., $\sigma < 1/2$).</p> <p>For example, the user enters the meanlog (μ) and sdlog (s) parameters for all lognormal distributions in the model. Thus, if you want to define the size of loss to be lognormal with mean 100,000 and standard deviation 100,000, you need to do a conversion calculate first.</p> <p>To achieve the given size of loss distribution above, you can set the parameter $\mu = \text{meanlog} = 11.16636357$ and $s = \text{sdlog} = 0.832549779$</p> <p>since this produces a size of loss X with $EX = \exp\left(\mu + \frac{1}{2}\sigma^2\right) \approx 100,000$.</p> <p>and second moment $EX^2 = \exp(2\mu + 2\sigma^2) \approx 2 \times 10^{10}$</p> <p>This results in $Var(X) \sim 2 \times 10^{10} - 10^{10}$ and $s(X) \sim 100,000$.</p>
<pre>rm(x, meanlog = 11.16636357, sdlog = 0.832549779)</pre>	<p>CDF Sample <code>plot(x, plnorm(x, meanlog = 11.16636357, sdlog=0.832549779), type="l", main = "Lognormal CDF", col="red")</code></p>	
		

Name	R Representation and Explanation	
Negative Binomial	Density	dnbinom(x, size, prob, mu, log = FALSE)
	Probability Function	pnbinom(q, size, prob, mu, lower.tail = TRUE, log.p = FALSE)
	Quantile Function	qnbinom(p, size, prob, mu, lower.tail = TRUE, log.p = FALSE)
	Random Generation	rnbinom(n, size, prob, mu)
	Details	<p>The negative binomial distribution with size = n and prob = p has density</p> $\binom{x+n-1}{n-1} p^n (1-p)^x = \frac{\Gamma(x+n)}{x! \Gamma(n)} p^n (1-p)^x$ <p>for $x = 0, 1, 2, \dots$, where n is a real number, $n > 0$ and $0 < p \leq 1$.</p> <p>If n is a positive integer, this represents the number of failures before the n^{th} success in a series of independent Bernoulli trials with probability of success p. The mean is $\frac{n(1-p)}{p}$ and the variance is $\frac{n(1-p)}{p^2}$. A negative binomial distribution results from a Poisson distribution whose mean has a gamma prior distribution with scale $(1-p)/p$ and shape n. (This definition allows non-integer values of size.)</p> <p>An alternative parametrization (often used in ecology) is by the <i>mean</i> μ, and size, the <i>dispersion parameter</i>, where $\text{prob} = \text{size}/(\text{size}+\mu)$. The variance is $\mu + \mu^2/\text{size}$ in this parametrization or $n(1-p)/p^2$ in the first one.</p>
	CDF Sample: <code>x<-(1:40)</code> <code>plot(pnbinom(x, size=50, prob=0.8, log=FALSE), main = "Negative Binomial CDF", col="red")</code>	

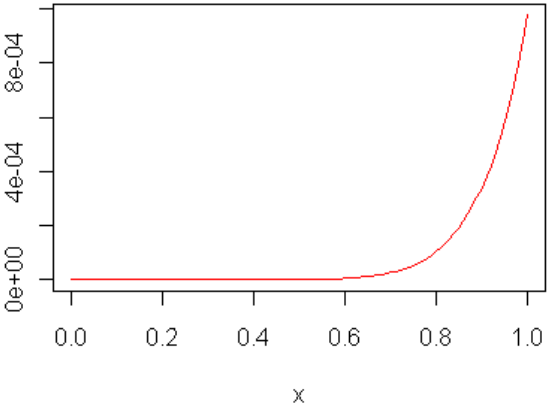
Modeling Loss Emergence and Settlement Processes

Name	R Representation and Explanation	
Pareto	Density	dpareto(x, shape, scale, log = FALSE)
	Probability Function	ppareto(q, shape, scale, lower.tail = TRUE, log.p = FALSE)
	Quantile Funtion	qpareto(p, shape, scale, lower.tail = TRUE, log.p = FALSE)
	Random Generation	rpareto(n, shape, scale)
	Details	<p>The Pareto distribution with parameters shape = a and scale = θ has density:</p> $f(x) = \alpha \theta^\alpha (x + \theta)^{-(\alpha+1)}$ <p>for $x > 0$, $a > 0$ and $\theta > 0$.</p>
<p>CDF Sample: plot(function(x) ppareto(x, shape=5, scale=1.2), main = "Pareto CDF", col="red")</p> <div style="text-align: center;"> <p>Pareto CDF</p>  </div>		

Modeling Loss Emergence and Settlement Processes

Name	R Representation and Explanation	
Poisson	Density	<code>dpois(x, lambda, log = FALSE)</code>
	Probability Function	<code>ppois(q, lambda, lower.tail = TRUE, log.p = FALSE)</code>
	Quantile Funtion	<code>qpois(p, lambda, lower.tail = TRUE, log.p = FALSE)</code>
	Random Generation	<code>rpois(n, lambda)</code>
	Details	<p>The Poisson distribution with parameter λ has density</p> $p(x) = e^{-\lambda} \frac{\lambda^x}{x!}$ <p>for $x = 0, 1, 2, \dots$. The mean and variance are $E(X) = Var(X) = \lambda$.</p> <p>If an element of x is not integer, the result of <code>dpois</code> is zero, with a warning. $p(x)$ is computed using Loader's algorithm.</p> <p>The quantile is right continuous: <code>qpois(p, lambda)</code> is the smallest integer x such that $P(X \leq x) \geq p$.</p> <p>Setting <code>lower.tail = FALSE</code> allows to get much more precise results when the default, <code>lower.tail = TRUE</code> would return 1, see the example below.</p>
<p>CDF Sample: <code>x<-(1:80)</code> <code>plot(ppois(x, lambda=120, log=FALSE), main = "Poisson CDF", col="red")</code></p> <div style="text-align: center;"> </div>		

Modeling Loss Emergence and Settlement Processes

Name	R Representation and Explanation	
Weibull	Density	dweibull(x, shape, scale = 1, log = FALSE)
	Probability Function	pweibull(q, shape, scale = 1, lower.tail = TRUE, log.p = FALSE)
	Quantile Function	qweibull(p, shape, scale = 1, lower.tail = TRUE, log.p = FALSE)
	Random Generation	rweibull(n, shape, scale = 1)
	Details	<p>The Weibull distribution with shape parameter τ and scale parameter θ has density $f(x)$ and c.d.f. $F(x)$ given by</p> $f(x) = \frac{\tau x^{\tau-1}}{\theta^\tau} \exp[-(x/\theta)^\tau], \quad F(x) = 1 - \exp[-(x/\theta)^\tau] \quad \text{for } x > 0.$ <p>The mean is $E(X) = \theta \Gamma(1 + 1/\tau)$, and $Var(X) = \theta^2 * [\Gamma(1 + 2/\tau) - \{\Gamma(1 + 1/\tau)\}^2]$</p>
<p>CDf Sample: plot(function(x) pweibull(x, shape=10, scale=2.0, log=FALSE), main = "Weibull Cumulative", col="red")</p> <p>ction(x) pweibull(x, shape = 10, scale = 2, log = FALSE</p>	<p style="text-align: center;">Weibull Cumulative</p> 	

APPENDIX B

6.2.1 Test of Elementary Frequencies, Trend, and Zero-modification

This is a complete description of the parameters for this run:

Test run 10/27/2009

Project name: Frequency Test

Purpose: Test frequency with trend. Two types within one line.

- One Line with annual frequency Poisson(120)
- Set claim/acc distribution matrix as follows:
 - Prob =75% that one Type 1 claim is generated.
 - Prob =25% that one Type 2 claim is generated.
- Freq Trend: 1.02 constant throughout
- $P(0) = 0.4$, $EstP(0) = 0.4$ for each Type.
- Accident Years: 2000-2002
- Random Seed: 16807
- Frequency correlation copula: normal Correlation=c() Dim = 1

Other frequency parameters

Monthly exposure: (1)

Seasonality: 1.0

The severity parameters were not used in this run. They are presented here only for completeness.

Type 1 and Type 2 severity are lognormal with different parameters:

Type 1: Lognormal, mean=100,000, std. dev. 100,000, max 1,000,000.

Type 2: Lognormal, mean 10,000, std. dev. 5,000, max 1,000,000.

Zero trend for severity.

Lags: Irrelevant for this run. We are not testing.

Report lag: Exponential with Rate = 4/365, mean=365/4 days. Max=365.

Payment lag: Exponential with Rate = 1/365, mean=365 days. Max = 700.

Inter-valuation lag: Exponential with Rate = 4/365, mean = 365/4 days. Max=365.

Correlation of Amount with lag: normal Correlation=c() Dim=2

Reserve adequacy: Irrelevant, leave at default values.

Recovery lag box: Not checked.

Alpha = 0 Inertia=0.2

Interpolation parameters: irrelevant

P(1) not set because recovery lag was not checked.

Minimum change = Min Rel Chg = 0.

Initial payment adequacy not set because recovery lag was not checked.

Modeling Loss Emergence and Settlement Processes

Run: 1,000 simulations.

Actuaries need tools that will enable them to better understand the underlying loss development process and will aid them in determining what methods and models work best in different reserving situations.

Section 6.2.1 discussed a series of GLM runs that investigated the effect of predictors on claim counts. This section provides more detail on the S-PLUS statements used.

Following is the full output from the S-PLUS “glm” command defining the “reduced” model5x:

```
> model5x<- glm(count ~ + Type + Status,
+ data = temp.dataacc.stack,
+ family = poisson,
+ x=T)
>
>
> summary(model5x,correlation=F)

Call: glm(formula = count ~ + Type + Status, family = poisson, data =
temp.dataacc.stack, x = TRUE)
Deviance Residuals:
    Min       1Q   Median       3Q      Max
-3.0487245 -0.817678553 -0.0481310816  0.498450005  4.53059039

Coefficients:
            Value      Std. Error    t value
(Intercept)  1.126272675  0.00276627012  407.144865
      Type -1.097685774  0.00378962189 -289.655751
      Status  0.410026757  0.00335189741  122.326762

(Dispersion Parameter for Poisson family taken to be 1 )

Null Deviance: 273221.814 on 143999 degrees of freedom
Residual Deviance: 160969.366 on 143997 degrees of freedom
Number of Fisher Scoring Iterations: 5
```

Following is the result of the “glm” command for the full model model6x:

```
> #####
> ##### The only possible variable to add is Type*Status
> ##### Just create the bigger model and compare
> ##### No need to use stepAIC here.
> #####
> model6x<- glm(count ~ Type + Status + Type*Status
+ data = temp.dataacc.stack,
+ family = poisson,
+ x=T)
>
>
> summary(model6x)

Call: glm(formula = count ~ Type + Status + Type * Status, family = poisson, data =
temp.dataacc.stack, x = TRUE)
Deviance Residuals:
    Min       1Q   Median       3Q      Max
-3.04838755 -0.817233217 -0.0487132512  0.497838102  4.52978802
```

Modeling Loss Emergence and Settlement Processes

```
Coefficients:
              value      Std. Error      t value
(Intercept)  1.12660568075  0.00300060918  375.458986266
Type         -1.09901753526  0.00600212647 -183.104694918
Status       0.40947269733  0.00387066987  105.788587162
Type:Status  0.00221507211  0.00773994871   0.286186924

(Dispersion Parameter for Poisson family taken to be 1 )

Null Deviance: 273221.814 on 143999 degrees of freedom
Residual Deviance: 160969.284 on 143996 degrees of freedom
Number of Fisher Scoring Iterations: 5

Correlation of Coefficients:
              (Intercept)      Type      Status
Type         -0.499924351
Status       -0.775217024  0.387549868
Type:Status  0.387678174 -0.775473675 -0.500089861
>
```

Results of the test to see whether interactive variable can be dropped:

```
>
> anova(model5x,model6x,test="Chi")
Analysis of Deviance Table

Response: count

  Terms Resid. Df Resid. Dev      Test Df
  1      + Type + Status    143997 160969.366
  2 Type + Status + Type * Status  143996 160969.284 +Type:Status  1

      Deviance      Pr(Chi)
  1
  2 0.0819088429 0.774727081
>
>
```

One way to test whether a Poisson GLM is appropriate is to see how close the dispersion parameter is to 1.0. Following is an execution of this process¹³.

```
> #####
> ##### Estimate dispersion parameter - see Faraway p 60.
> #####
> temp.pearson <- residuals(model5x,type="pearson") ### (y-fitted y)/(y^0.5)
> model5x$df.residual
[1] 143997
> tempdp <- sum(temp.pearson^2)/model5x$df.residual
> cat("\n Estimate of dispersion parameter using Pearson residuals.
disp=",tempdp,"\n")

Estimate of dispersion parameter using Pearson residuals. disp= 1.00653317880341
>
>
> temp.pearson <- residuals(model6x,type="pearson") ### (y-fitted y) / (y^0.5)
> model6x$df.residual
[1] 143996
```

¹³ [6] , p. 60.

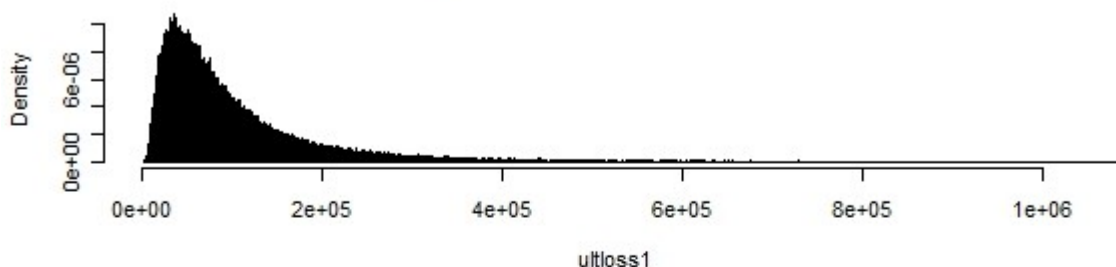
Modeling Loss Emergence and Settlement Processes

```
> tempdp <- sum(temp.pearson^2)/model6x$df.residual
> cat("\n Estimate of dispersion parameter using Pearson residuals.
disp=",tempdp,"\n")
  Estimate of dispersion parameter using Pearson residuals.  disp= 1.00659609377679
>
> #####
> ##### Dispersion in both cases very close to 1.0, supporting Poisson model
> #####
>
```

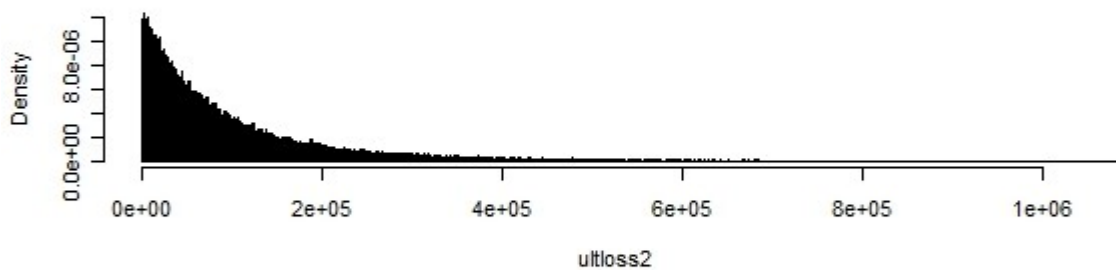
Appendix to section 6.2.2 Test of Severity Distributions

Following is the set of histograms, empirical densities and log densities, and empirical cdfs for the size of loss for each line.

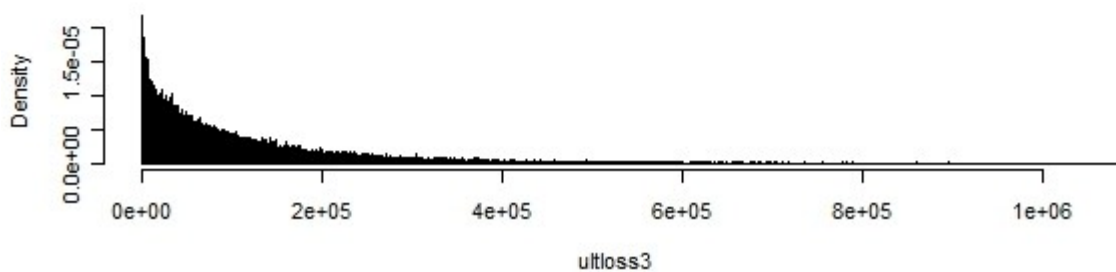
Histogram of observed data of Line 1



Histogram of observed data of Line 2



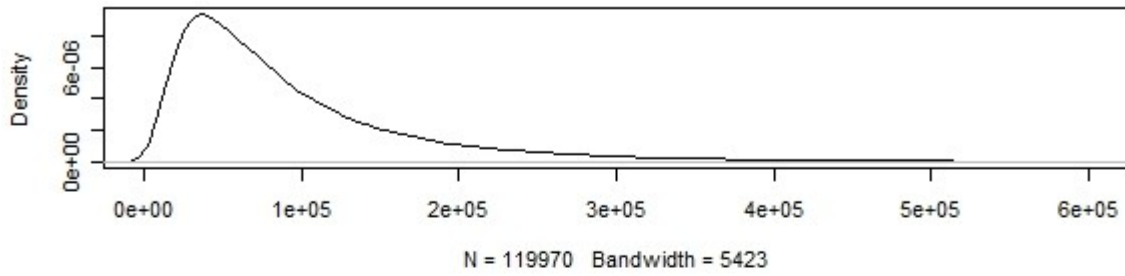
Histogram of observed data of Line 3



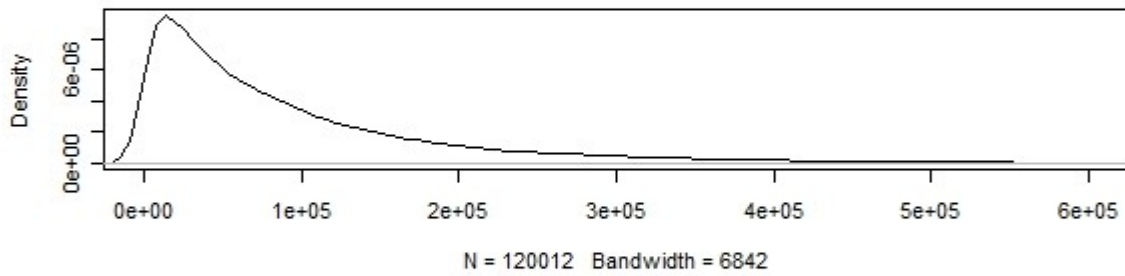
Modeling Loss Emergence and Settlement Processes

Density estimates:

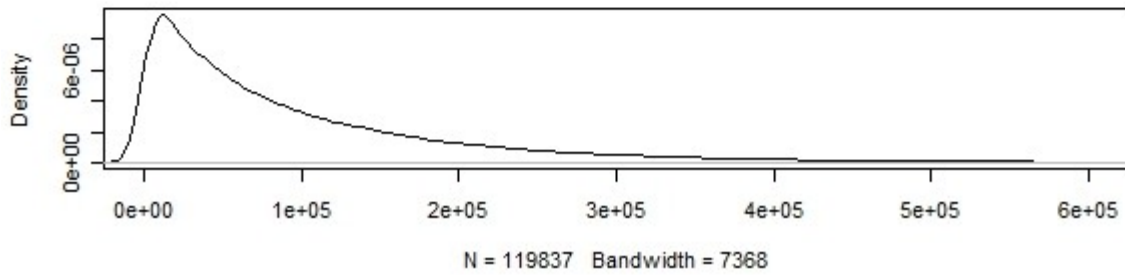
Density estimate of Line 1



Density estimate of Line 2

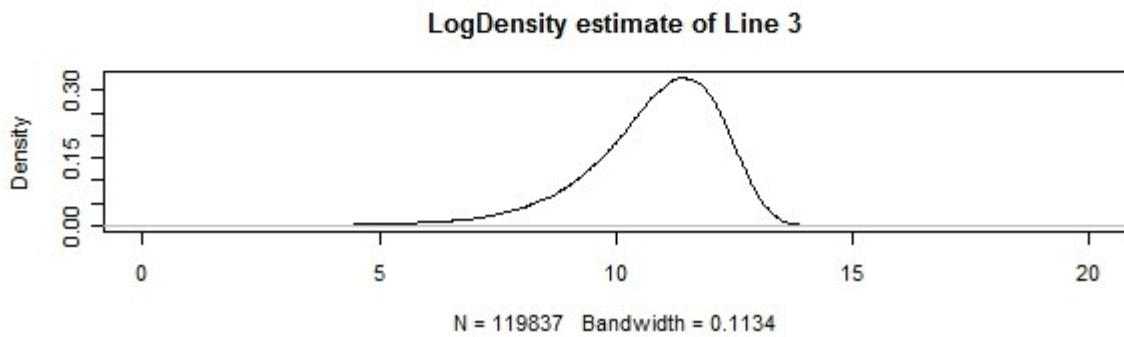
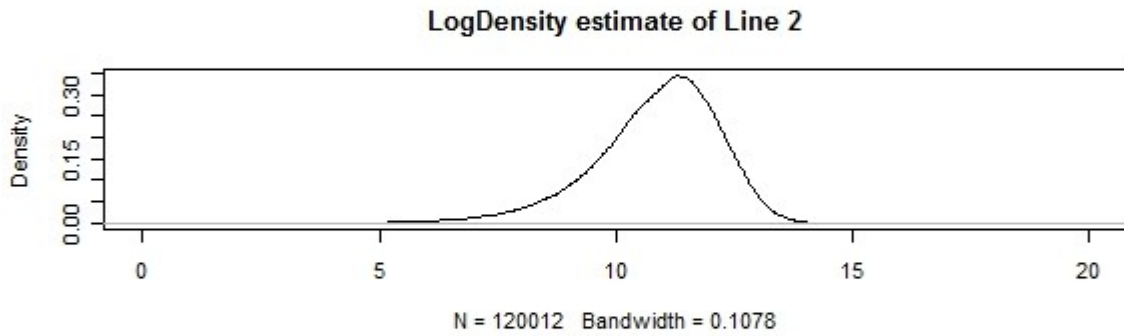
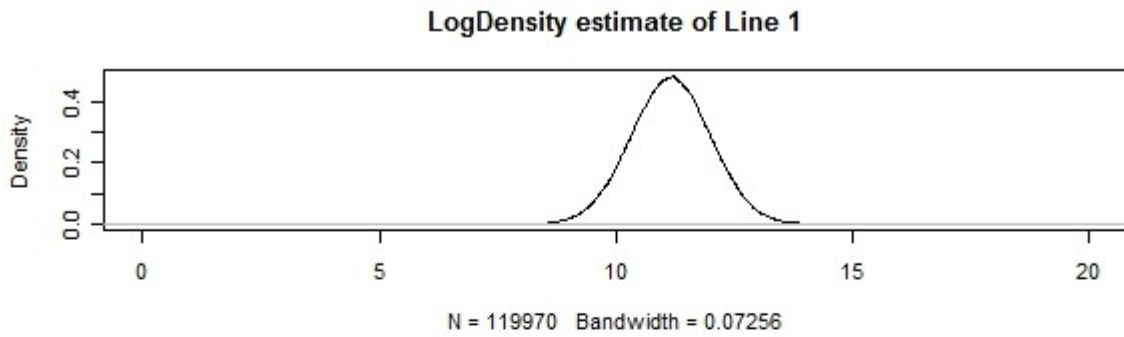


Density estimate of Line 3



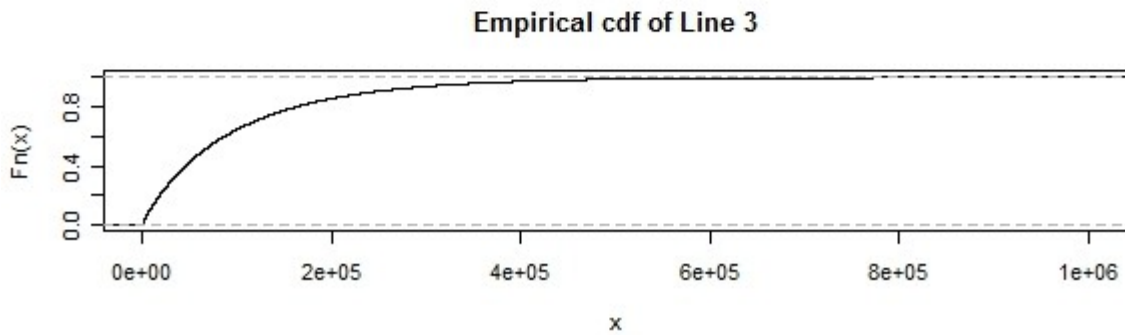
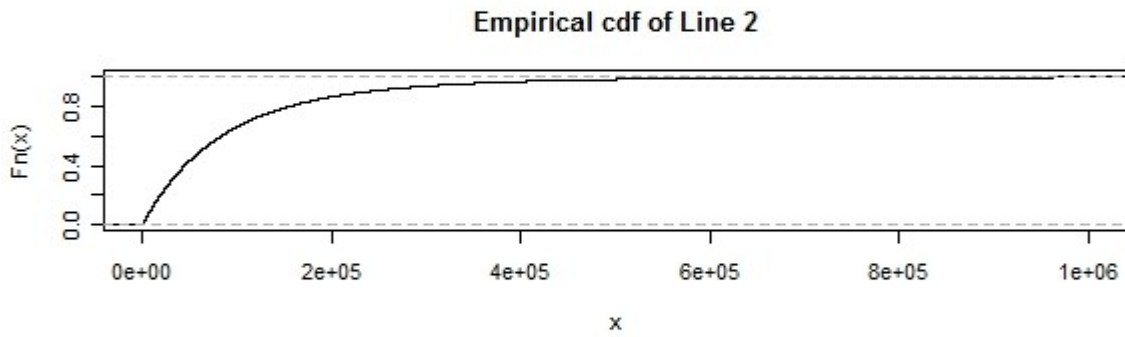
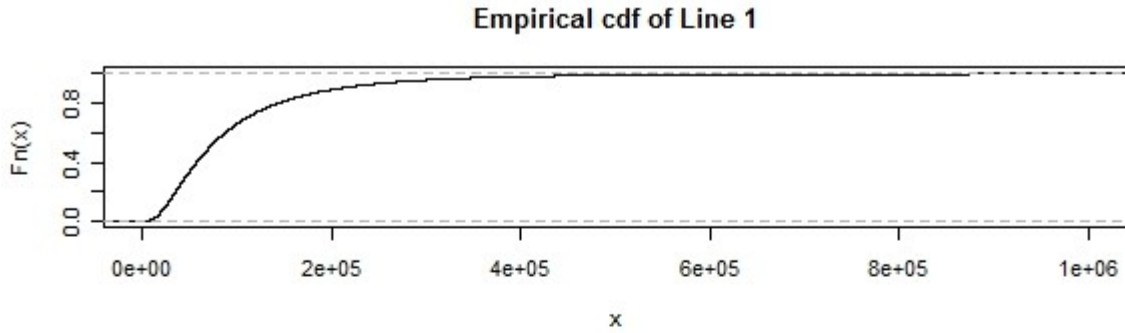
Modeling Loss Emergence and Settlement Processes

Log densities:



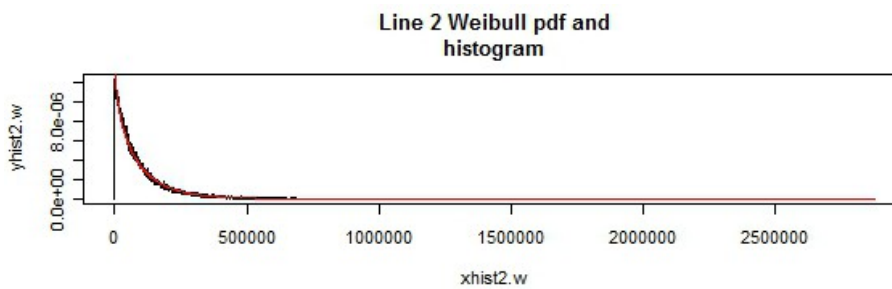
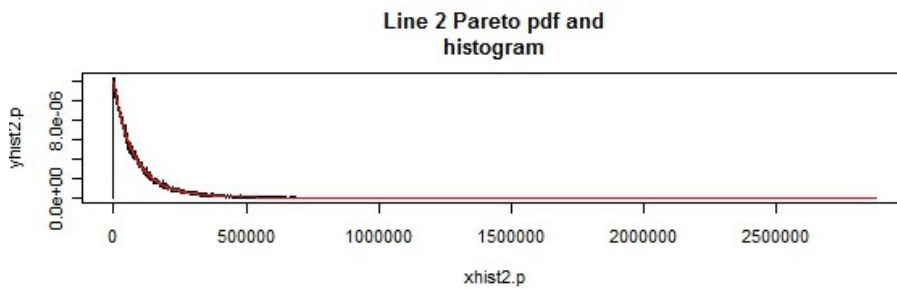
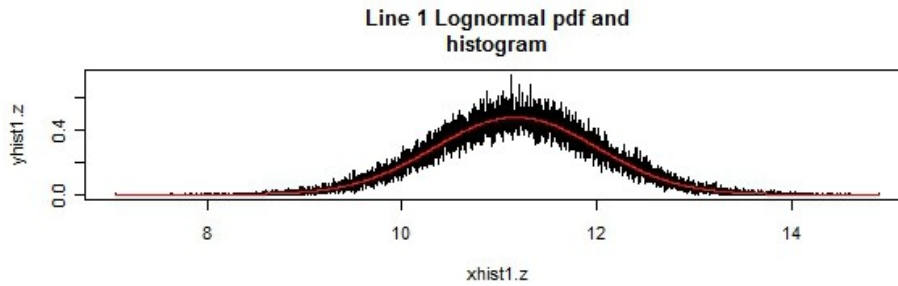
Modeling Loss Emergence and Settlement Processes

Empirical distribution functions:

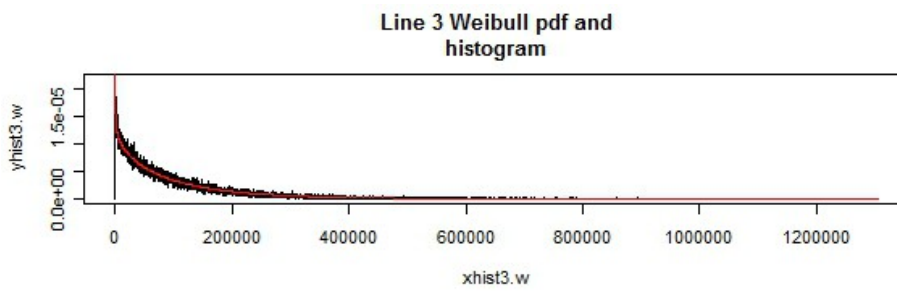
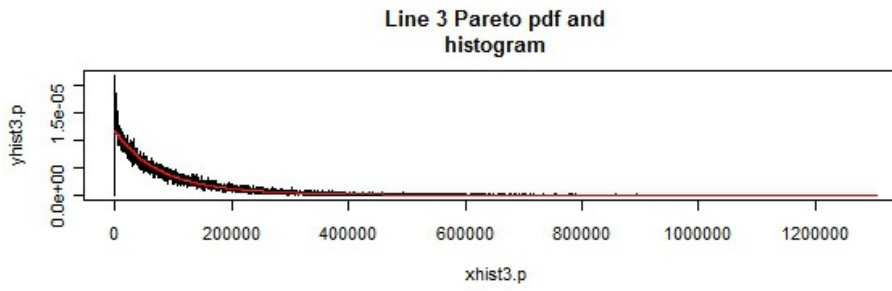


Modeling Loss Emergence and Settlement Processes

We next show the histograms and density functions for each model that produced a reasonable fit to the observed data.



Modeling Loss Emergence and Settlement Processes



Modeling Loss Emergence and Settlement Processes

Following is the complete “R” code for the severity testing.

```
library(stats4)
library(MASS)
library(actuar)
library(graphics)

##### Part I: Exploratory Data Analysis #####

ultloss1<-as.vector(scan("line 1 ultloss.txt"))
ultloss2<-as.vector(scan("line 2 ultloss.txt"))
ultloss3<-as.vector(scan("line 3 ultloss.txt"))

mean1 <- mean(ultloss1)
var1 <- var(ultloss1)

mean2 <- mean(ultloss2)
var2 <- var(ultloss2)

mean3 <- mean(ultloss3)
var3 <- var(ultloss3)

n1<-length(ultloss1)
n2<-length(ultloss2)
n3<-length(ultloss3)

#####
# Histograms, empirical density/log density, empirical cdf
#####

par(mfrow=c(3,1)) # Histograms
hist(ultloss1,main="Histogram of observed data of Line
1",freq=FALSE,breaks=10000,xlim=c(0,1050000))
hist(ultloss2,main="Histogram of observed data of Line
2",freq=FALSE,breaks=10000,xlim=c(0,1050000))
hist(ultloss3,main="Histogram of observed data of Line
3",freq=FALSE,breaks=10000,xlim=c(0,1050000))

par(mfrow=c(3,1)) # Density
plot(density(ultloss1),main="Density estimate of Line 1",xlim=c(-1000,600000))
plot(density(ultloss2),main="Density estimate of Line 2",xlim=c(-1000,600000))
plot(density(ultloss3),main="Density estimate of Line 3",xlim=c(-1000,600000))

par(mfrow=c(3,1)) # Log Density
plot(density(log(ultloss1)),main="LogDensity estimate of Line 1",xlim=c(0,20))
plot(density(log(ultloss2)),main="LogDensity estimate of Line 2",xlim=c(0,20))
plot(density(log(ultloss3)),main="LogDensity estimate of Line 3",xlim=c(0,20))

par(mfrow=c(3,1)) # Empirical Cumulative Distribution Function
plot(ecdf(ultloss1),main="Empirical cdf of Line 1",xlim=c(0,1e+06)) ## lightest tail
plot(ecdf(ultloss2),main="Empirical cdf of Line 2",xlim=c(0,1e+06)) ## heaviest tail
plot(ecdf(ultloss3),main="Empirical cdf of Line 3",xlim=c(0,1e+06)) ## heavier tail

##### Part II: Maximum Likelihood Estimates of Parameters #####

#####
# Pareto (shape=alpha, scale=theta) #
# pdf f(x) = alpha*theta^alpha/((x+theta)^(alpha+1)) #
#####

#####
# weibull(shape=tao, scale=lambda) #
# pdf f(x) = tao*x^(tao-1)/(lambda^tao)*exp(-(x/lambda)^(tao)) #
#####
```

Modeling Loss Emergence and Settlement Processes

```

# nll=-log(tao/(lambda^tao))-(tao-1)*log(x)+(x/k)^tao #
#####

#####
# lognormal(miu,sigma) #
# pdf f(x)=1/(sqrt(2*pi)*sigma)*exp(-(x-miu)^2/(2*sigma^2)) #
# nll = log(sqrt(2*pi)*sigma)+(x-miu)^2/(2*sigma^2) #
#####

ultloss1.0 <- ultloss1[ultloss1!=0]
n1.0 <- length(ultloss1.0)

ultloss2.0 <- ultloss2[ultloss2!=0]
n2.0 <- length(ultloss2.0)

ultloss3.0 <- ultloss3[ultloss3!=0]
n3.0 <- length(ultloss3.0)

#####
# Line 1 #
#####

## From the exploratory analysis, it is clear that Line 1 comes from lognormal.

fit1.ln <- fitdistr(log(ultloss1),"normal")
fit1.ln$estimate # mean sd
#11.1659376 0.8361509
-fit1.ln$loglik #148761.9

fit.ln1 <- fitdistr(log(ultloss1.0),dnorm,list(mean=miu1,sd=sigma1))
fit.ln1$estimate # mean sd
#11.1659362 0.8361455
-fit.ln1$loglik #148761.9

## QQ plot ##
qqnorm(log(ultloss1),main="Line 1, Lognormal")
abline(0,1,col="red")

## Chi-Square Test ##

ult1.cut <- cut(log(ultloss1),breaks = seq(0,9,10,11,12,13,22)) ## binning data
table.ult1 <- table(ult1.cut) ## binned data table
ult1.os <- c(as.vector(table.ult1)) ## vectorization

b = length(ult1.os)

labs <- levels(ult1.cut) ## extract the breakpoints
break.1 <- cbind(lower = as.numeric(sub("\\((.+),.*", "\\1", labs)),upper =
as.numeric(sub("[^,]*,([^)]*)\\)", "\\1", labs)))

ult1.ln <- pnorm(break.1,mean=fit.ln1$estimate[1],sd=fit.ln1$estimate[2])

ult1.prob <- (ult1.ln[,2]-ult1.ln[,1])[1:b-1] ## cut off at the last bin to absorb
all tail prob
ult1.ex <- n1*c(ult1.prob,1-sum(ult1.prob))

#cbind(ult1.ex,ult1.os)

E.1 = ult1.ex
O.1 = ult1.os

x.sq.1 = (E.1-O.1)^2/E.1

```

Modeling Loss Emergence and Settlement Processes

```

#cbind(E.1,0.1,x.sq.1)
##chi-square test statistic##
df=length(E.1)-1-2    ## df = 3
chi.sq.1 <- sum(x.sq.1)    ## test statistic
chi.sq.1              ## 3.594133
qchisq(.95,df)        # 7.814728
1-pchisq(chi.sq.1,df)   # 0.308757

## Chi-Square Test in R ##
chisq.test(0.1,p=E.1/n1)

## Chi-squared test for given probabilities
## X-squared = 3.5941, df = 5, p-value = 0.6092

#####
# Calculation for the test statistic
#
#-----
#           E.1    0.1    x.sq.1
#[1,]    575.073    531 3.377707e+00
#[2,]   9213.966   9254 1.739471e-01
#[3,]  40759.725  40740 9.545298e-03
#[4,]  50315.005  50316 1.969319e-05
#[5,]  17410.264  17434 3.236081e-02
#[6,]   1695.968   1695 5.528125e-04
#####

#####
# Line 2 #
#####

## 2.1-Pareto ##
fit.p2<-fitdistr(ultloss2.0,dpareto,list(shape=6,scale=500000)) ## list() provides
initial values for optimization
fit.p2$estimate # shape      scale
                #5.97635e+00 5.00000e+05
-fit.p2$loglik  #1500363

## 2.2-weibull (second method slightly better) ##
fit2.w <- fitdistr(ultloss2.0,"weibull")
fit2.w$estimate # shape      scale
                # 9.056193e-01 9.750673e+04
fit2.w$loglik  # -1500950

fit.w2 <- fitdistr(ultloss2.0,dweibull,list(shape=.9097626,scale=95000))
fit.w2$estimate # shape      scale
                # 9.009281e-01 9.500000e+04
-fit.w2$loglik  # 1500926

## qq plot ##
par(mfrow=c(2,1))
thqua.p2 <- rpareto(n2,shape=fit.p2$estimate[1],scale=fit.p2$estimate[2])
qqplot(ultloss2,thqua.p2,xlab="Sample Quantiles", ylab="Theoretical
Quantiles",main="Line 2, Pareto")
abline(0,1,col="red")

thqua.w2 <- rweibull(n2,shape=fit.w2$estimate[1],scale=fit.w2$estimate[2])
qqplot(ultloss2,thqua.w2,xlab="Sample Quantiles", ylab="Theoretical Quantiles",
main="Line 2, Weibull")
abline(0,1,col="red")

```

Modeling Loss Emergence and Settlement Processes

```
## 2.3 Chi-Square Test ##

#2.3.1 Pareto Chi-Square #
m = mean(ultloss2)
s = sqrt(var(ultloss2))

ult2.cut <- cut(ultloss2.0,breaks = c(0,m-
s/2,m,m+s/4,m+s/2,m+s,m+2*s,2*max(ultloss2))) ##binning data
table.ult2 <- table(ult2.cut) ## binned data table
ult2.os <- c(as.vector(table.ult2)) ## vectorization

b = length(ult2.os)

labs.2 <- levels(ult2.cut) ## extract the breakpoints
break.2 <- cbind(lower = as.numeric(sub("\\((.+),.*", "\\1", labs.2)),upper =
as.numeric(sub("[^,]*,([^)]*)\\)", "\\1", labs.2))) See note below14

ult2.p <- ppareto(break.2,shape=fit.p2$estimate[1],scale=fit.p2$estimate[2]) ## Pareto
cdf values at break points

ult2.prob <- (ult2.p[,2]-ult2.p[,1])[1:b-1] ## Probabilities of each interval
ult2.ex <- n2.0*c(ult2.prob,1-sum(ult2.prob)) ## Expected frequency of each interval
and the "excess" interval

#cbind(ult2.ex,ult2.os) ## expected and observed frequencies

E.2 = ult2.ex
O.2 = ult2.os

x.sq.2 = (E.2-O.2)^2/E.2

#cbind(E.2,O.2,x.sq.2) ## expected, observed, and chi-square of each interval
after full adjustment

##chi-square test statistic##

df=length(E.2)-1-2 ## df = 4

chi.sq.2 <- sum(x.sq.2) ## test statistic
## 6.155374

qchisq(.95,df) ## critical value ## 9.487729
```

¹⁴ For example, suppose that labs.2 consists of seven intervals:

```
"(0,3.97e+04]"
"(3.97e+04,1e+05]"
"(1e+05,1.31e+05]"
"(1.31e+05,1.61e+05]"
"(1.61e+05,2.22e+05]"
"(2.22e+05,3.43e+05]"
"(3.43e+05,5.75e+06]"
```

Then break.2 looks like:

	lower	upper
[1,]	0	39700
[2,]	39700	100000
[3,]	100000	131000
[4,]	131000	161000
[5,]	161000	222000
[6,]	222000	343000
[7,]	343000	5750000

Modeling Loss Emergence and Settlement Processes

```

1-pchisq(chi.sq.2,df)      ## p-value  ## 0.1878414
## chi-square goodness-of-fit test from R ##
chisq.test(0.2,p=E.2/n2.0)

#####
#   chi-squared test for given probabilities
# data:  0.2
# X-squared = 6.1554, df = 6, p-value = 0.406
#####

##### Test Statistic Calculation #####
#
#-----
#          E.2    0.2    x.sq.2
#[1,] 43993.890 44087 0.19705959
#[2,] 35651.989 35680 0.02200752
#[3,] 10493.758 10323 2.77864169
#[4,]  7240.583  7269 0.11152721
#[5,]  9277.383  9164 1.38570182
#[6,]  8063.576  8176 1.56743997
#[7,]  5289.820  5312 0.09299630
#####

### 2.3.2 weibull ###
m = mean(ultloss2)
s = sqrt(var(ultloss2))

ult2.cut      <-      cut(ultloss2.0,breaks      =      c(0,m-
s/2,m,m+s/4,m+s/2,m+s,m+2*s,2*max(ultloss2))) ##binning data
table.ult2 <- table(ult2.cut)      ## binned data table
ult2.os <- c(as.vector(table.ult2))      ## vectorization

b = length(ult2.os)

labs.2 <- levels(ult2.cut)      ## extract the breakpoints
break.2 <- cbind(lower = as.numeric(sub("\\((.+),.*", "\\1", labs.2)),upper =
as.numeric(sub("[^,]*,([^)]*)\\)", "\\1", labs.2)))

ult2.w      <-      pweibull(break.2,shape=fit.w2$estimate[1],scale=fit.w2$estimate[2])
##weibull cdf values at break points

ult2.prob <- (ult2.w[,2]-ult2.w[,1])[1:b-1] ## Probabilities of each interval
ult2.ex <- n2.0*c(ult2.prob,1-sum(ult2.prob)) ## Expected frequency of each interval
and the "excess" interval

#cbind(ult2.ex,ult2.os)      ## expected and observed frequencies before adjustment

E.2 = ult2.ex
O.2 = ult2.os

x.sq.2 = (E.2-O.2)^2/E.2

#cbind(E.2,O.2,x.sq.2)      ## expected, observed, and chi-square of each interval
after full adjustment

##chi-square test statistic##

df=length(E.2)-1-2      ## df = 4

chi.sq.2 <- sum(x.sq.2)      ## test statistic
chi.sq.2      ## 270.3838

```

Modeling Loss Emergence and Settlement Processes

```
qchisq(.95,df)          ## critical value ## 9.487729
1-pchisq(chi.sq.2,df)   ## p-value ## 0
# chi-square goodness-of-fit test from R #
chisq.test(O.2,p=E.2/n2.0)

#####
# Chi-squared test for given probabilities
#data: 0.2
#X-squared = 270.3838, df = 6, p-value < 2.2e-16
#####

## Test Statistic Calculation #####
#-----#
#          E.2      0.2      x.sq.2
#[1,] 43917.914 44087 0.6509914
#[2,] 33982.992 35680 84.7434487
#[3,] 10551.628 10323 4.9538155
#[4,] 7532.284 7269 9.2028627
#[5,] 10023.960 9164 73.7763860
#[6,] 9007.970 8176 76.8402384
#[7,] 4994.251 5312 20.2160763
#####

#####
# Line 3 #
#####

## 3.1-Pareto ##
fit.p3<-fitdistr(ultloss3.0,dpareto,list(shape=7,scale=6.026793e+05))
fit.p3$estimate # shape scale
#6.966806e+00 6.026793e+05
-fit.p3$loglik #1499343

## 3.2-weibull (first method slightly better) ##
fit.w3 <- fitdistr(ultloss3.0,"weibull")
fit.w3$estimate # shape scale
# 9.052532e-01 9.907429e+04
-fit.w3$loglik # 1498920

fit3.w <- fitdistr(ultloss3.0,dweibull,list(shape=0.9,scale=100000))
fit3.w$estimate # shape scale
#9.067305e-01 9.999992e+04
-fit3.w$loglik #1498955

## QQ plot ##
par(mfrow=c(2,1))
thqua.p3 <- rpareto(n3,shape=fit.p3$estimate[1],scale=fit.p3$estimate[2])
qqplot(ultloss3,thqua.p3,xlab="Sample Quantiles", ylab="Theoretical
Quantiles",main="Line 3, Pareto")
abline(0,1,col="red")

thqua.w3 <- rweibull(n3,shape=fit.w3$estimate[1],scale=fit.w3$estimate[2])
qqplot(ultloss3,thqua.w3,xlab="Sample Quantiles", ylab="Theoretical
Quantiles",main="Line 3, weibull")
abline(0,1,col="red")

## 3.3 Chi-Square Test ##

# 3.3.1 Pareto #
m = mean(ultloss3)
s = sqrt(var(ultloss3))
```

Modeling Loss Emergence and Settlement Processes

```

M=max(ultloss3)

ult3.cut <- cut(ultloss3,breaks = c(seq(0,1000000,200000),2000000)) ##binning data
table.ult3 <- table(ult3.cut) ## binned data table
ult3.os <- c(as.vector(table.ult3)) ## vectorization

b=length(ult3.os)

labs.3 <- levels(ult3.cut) ## extract the breakpoints
break.3 <- cbind(lower = as.numeric(sub("\\((.+),.*", "\\1", labs.3)),upper =
as.numeric(sub("[^,]*,([^,]*)\\)", "\\1", labs.3)))

ult3.p <- ppareto(break.3,shape=fit.p3$estimate[1],scale=fit.p3$estimate[2])

ult3.probab <- (ult3.p[,2]-ult3.p[,1])[1:b-1]
ult3.ex <- n3*c(ult3.probab,1-sum(ult3.probab))

#cbind(ult3.ex,ult3.os)

E.3 = ult3.ex
O.3 = ult3.os

x.sq.3 = (E.3-O.3)^2/E.3

#cbind(E.3,O.3,x.sq.3)

# Chi-Square Test Statistic #
df=length(E.3)-1-2 # df = 3

chi.sq.3 <- sum(x.sq.3) ## chi-square test statistic
chi.sq.3 ## 275.7469

qchisq(.95,df) ##critical value 7.814728
1-pchisq(chi.sq.3,df) ##p-value ## 0

# chi square test in R #
chisq.test(O.3,p=E.3/n3)

#####
# Chi-squared test for given probabilities
# data: 0.3
# X-squared = 275.7469, df = 5, p-value < 2.2e-16
#####

##### Test Statistic Calculation #####
#
#-----
# [1,] 103561.6008 102680 7.504906
# [2,] 12820.7894 13997 107.908432
# [3,] 2481.6241 2535 1.148034
# [4,] 639.8169 491 34.613767
# [5,] 201.5347 102 49.158580
# [6,] 131.6340 32 75.413160
#####

# 3.3.2 weibull
m = mean(ultloss3)
s = sqrt(var(ultloss3))
M=max(ultloss3)

ult3.cut <- cut(ultloss3,breaks = c(seq(0,1000000,200000),2000000)) ##binning data

```

Modeling Loss Emergence and Settlement Processes

```

table.ult3 <- table(ult3.cut)                ## binned data table
ult3.os <- c(as.vector(table.ult3))        ## vectorization

b=length(ult3.os)

labs.3 <- levels(ult3.cut)                ## extract the breakpoints
break.3 <- cbind(lower = as.numeric(sub("\\((.+),.*", "\\1", labs.3)),upper =
as.numeric(sub("[^,]*,([^)]*)\\)", "\\1", labs.3)))

ult3.w <- pweibull(break.3,shape=fit.w3$estimate[1],scale=fit.w3$estimate[2])

ult3.prob <- (ult3.w[,2]-ult3.w[,1])[1:b-1]
ult3.ex <- n3*c(ult3.prob,1-sum(ult3.prob))

#cbind(ult3.ex,ult3.os)

E.3 = ult3.ex
O.3 = ult3.os

x.sq.3 = (E.3-O.3)^2/E.3

#cbind(E.3,O.3,x.sq.3)

# Chi-Square Test Statistic #

df=length(E.3)-1-2

chi.sq.3 <- sum(x.sq.3)    ## chi-square test statistic
chi.sq.3                ## 70.21185

qchisq(.95,df)          ##critical value ## 7.814728

1-pchisq(chi.sq.3,df)  ##p-value 3.885781e-15

# chi square test in R #

chisq.test(O.3,p=E.3/n3)
#####
#      Chi-squared test for given probabilities
# data:  O.3
# X-squared = 70.2119, df = 5, p-value = 9.259e-14
#####

### Test Statistic Calculation #####
#
#-----
#      E.3      O.3      x.sq.3
#[1,] 101709.59203 102680  9.2586314
#[2,] 14641.17032 13997  28.3416822
#[3,]  2759.98729  2535  18.3404039
#[4,]   567.24209   491  10.2475765
#[5,]   122.92003   102  3.5604243
#[6,]    36.08824    32  0.4631346
#####

```

Appendix to section 6.2.3 -- Testing Correlated Frequencies

The following material is the "R Code" used to produce the results in 6.2.3.

```
library(MASS)
library(methods)
library(mvtnorm)
library(scatterplot3d)
library(mnormt)
library(sn)
library(pspline)
library(copula)

# import data
#annual frequency for each line and each simulation
datar<-read.csv("D:/LSMWP/byyear.csv")
summary(datar)
n<-length(datar$Line.1)
set.seed(123)
x<- sapply(datar, rank, ties.method = "random") / (n + 1)
x12<-subset(x,select=-Line.3)
x13<-subset(x,select=-Line.2)
x23<-subset(x,select=-Line.1)
plot(x12)
plot(x13)
plot(x23)
write.csv(x12,"D:/LSMWP/x12ry.csv")
write.csv(x23,"D:/LSMWP/x23ry.csv")
write.csv(x13,"D:/LSMWP/x13ry.csv")

#Set up copula object for copula distribution and goodness-of-fit test later
normal.cop <- normalCopula(c(0,0,0),dim=3,dispstr="un")

#Copula fit with prespecified type.

date()
fit.normal<-fitCopula(normal.cop,x,method="ml")
fit.normal

fit.normal<-fitCopula(normal.cop,x,method="itau")
fit.normal

date()

#Copula Goodness-of-fit test

date()

normal2.cop <- normalCopula(c(0),dim=2,dispstr="un")
gofCopula(normal2.cop, x12, N=100, method = "mp1")
gofCopula(normal2.cop, x13, N=100, method = "mp1")
gofCopula(normal2.cop, x23, N=100, method = "mp1")

#gofCopula(normal.cop, x, N=100, method = "mp1")
#gofCopula(normal.cop, x, N=100, method = "itau")
date()

#K-S test.
normal.fit12<-normalCopula(0, dim=2)
normal.fit13<-normalCopula(0.99,dim=2)
```

Modeling Loss Emergence and Settlement Processes

```
normal.fit23<-normalCopula(-0.01,dim=2)
y12<-rcopula(normal.fit12,n)
y23<-rcopula(normal.fit23,n)
y13<-rcopula(normal.fit13,n)
ks.test(x12,y12)
ks.test(x13,y13)
ks.test(x23,y23)
```

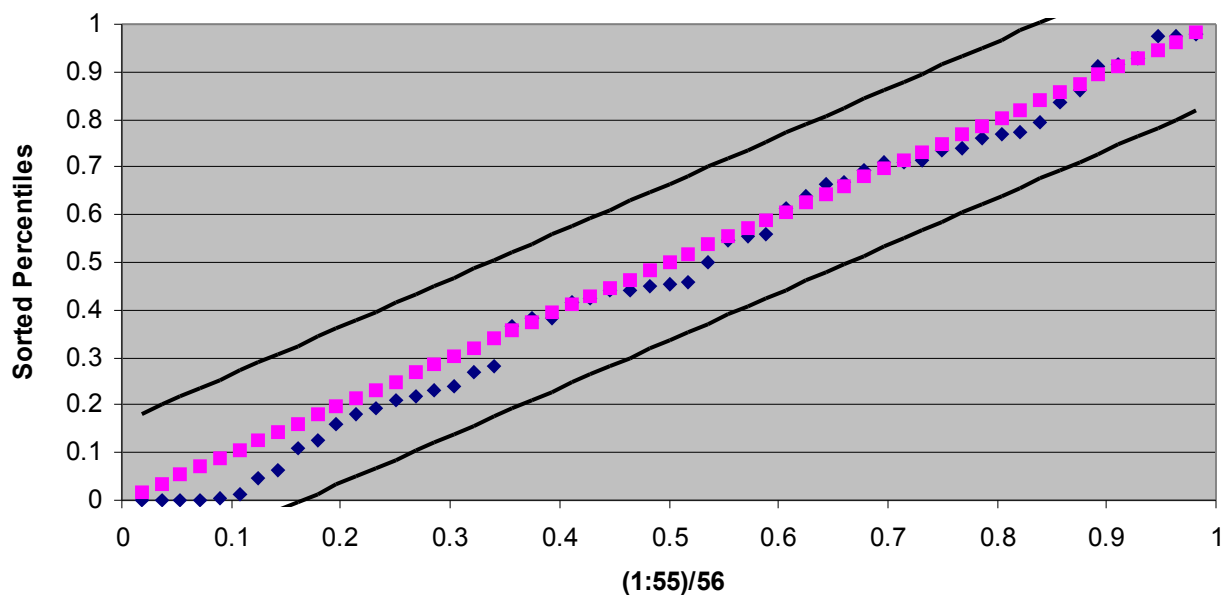
The following table shows the first three records and last two records from "D:/LSMWP/byyear.csv", the dataset used in the R code above.,

Line 1	Line 2	Line 3
114	95	117
89	85	90
94	119	99
....
94	113	94
105	97	105

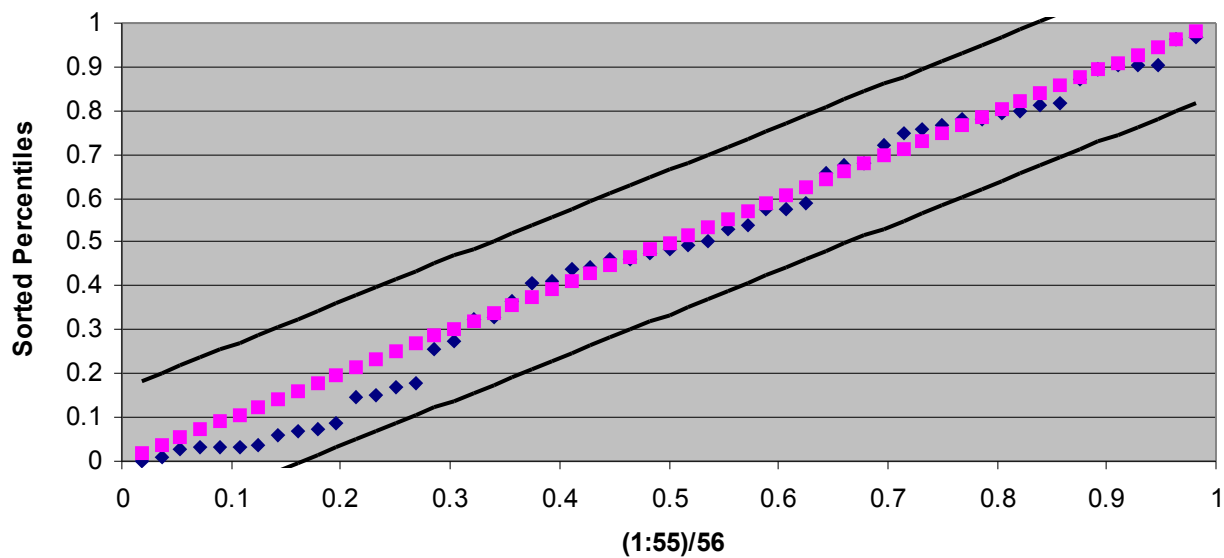
Each row represents the number of claims by line for one simulation.

Appendix C

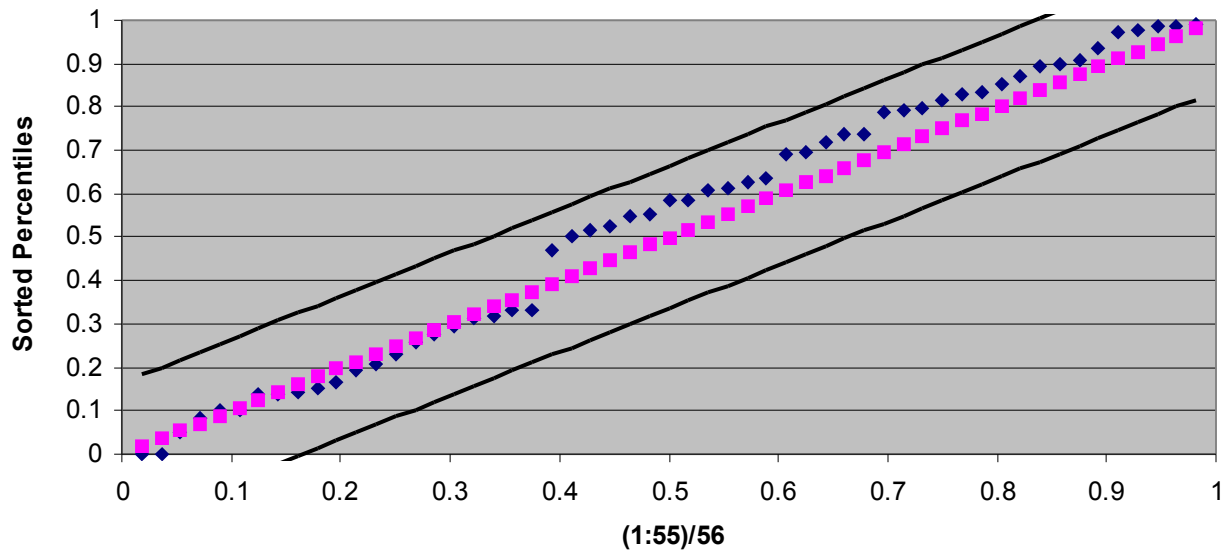
Candidate 1: Appendix C.1



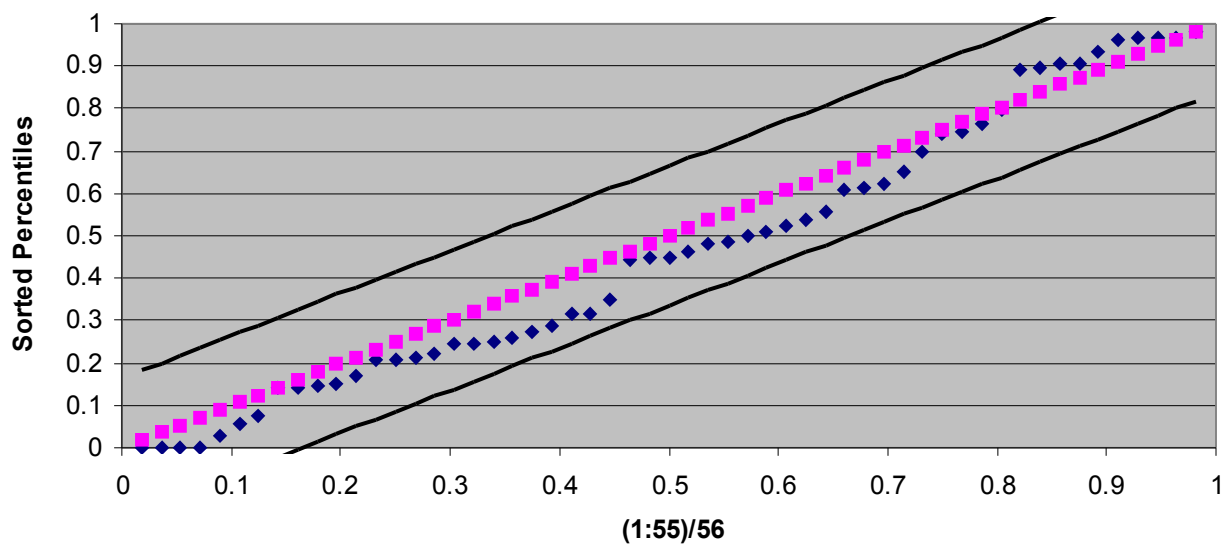
Candidate 2: Appendix C.2



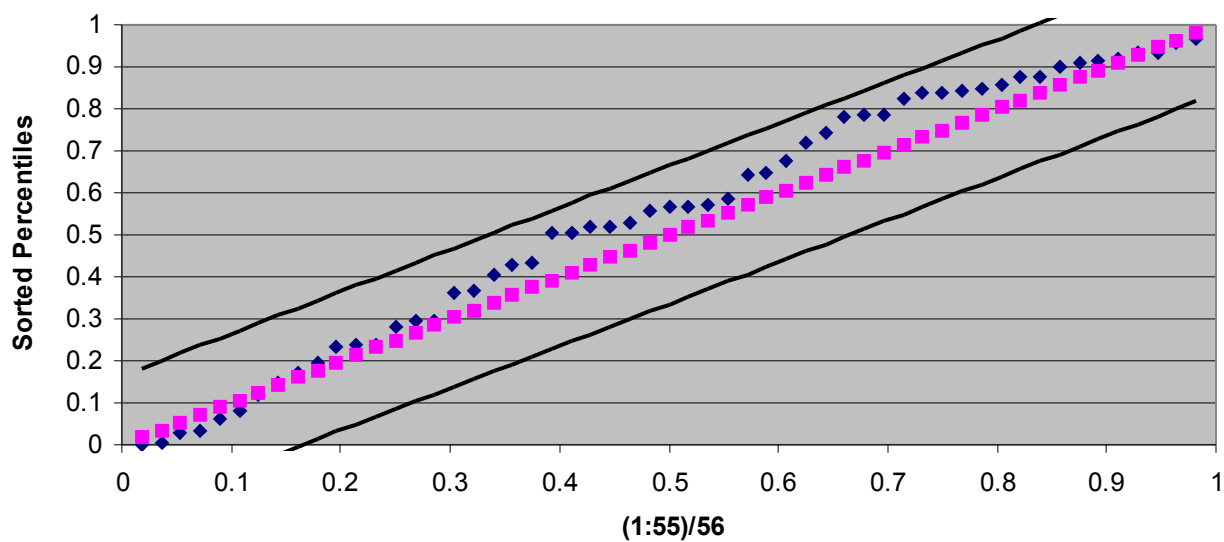
Candidate 3: Appendix C.3



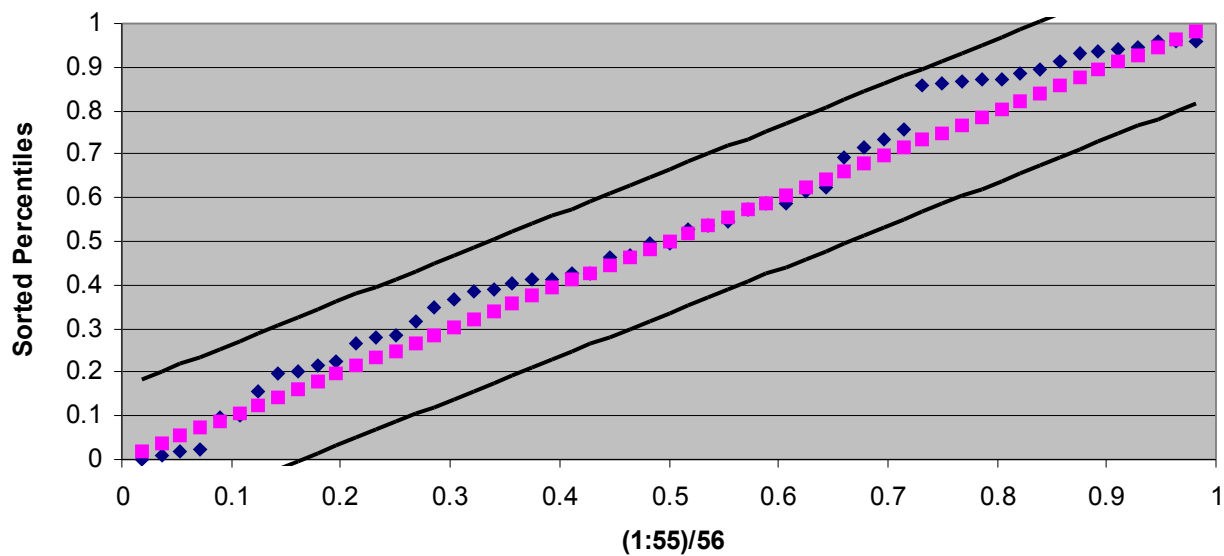
Candidate 4: Appendix C.4



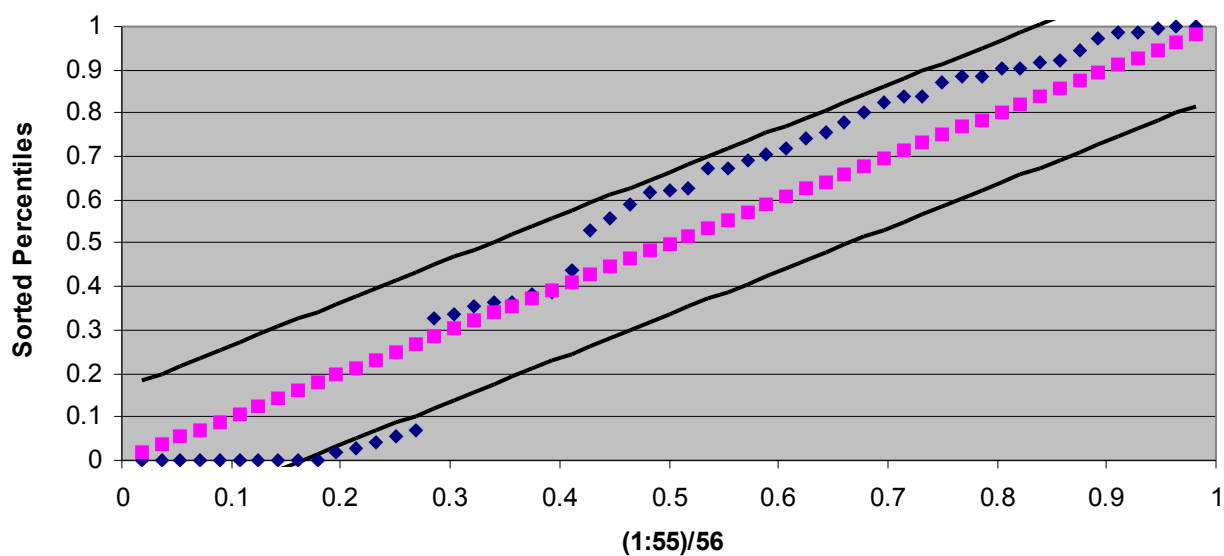
Candidate 5: Appendix C.5



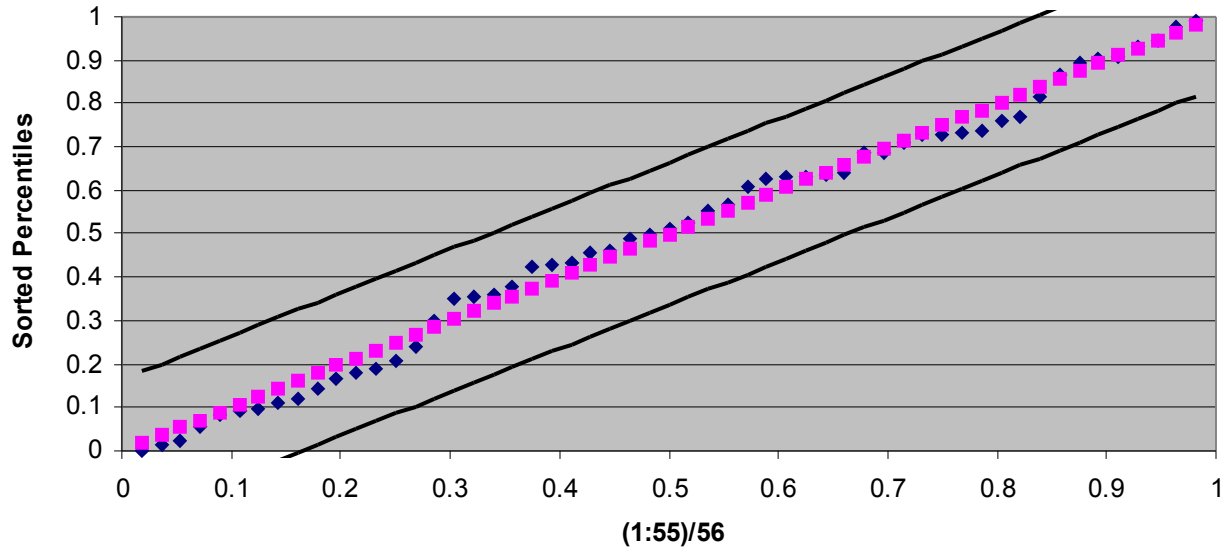
Candidate 6: Appendix C.6



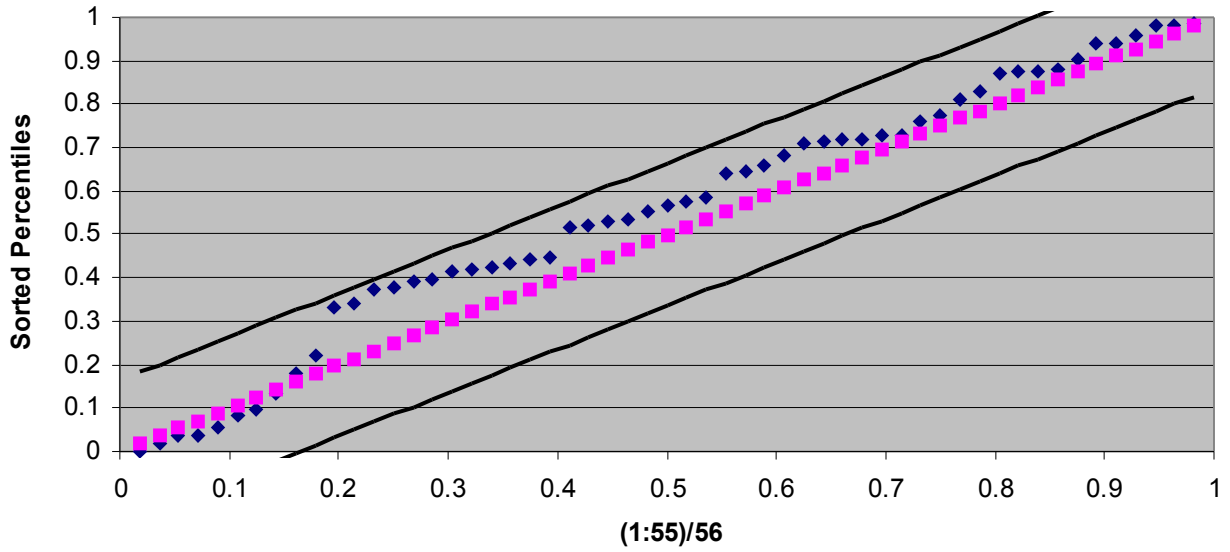
Candidate 7: Appendix C.7



Candidate 8: Appendix C.8



Candidate 9: Appendix C.9



Candidate 10: Appendix C.10

