

Statistical Modelling for Health Insurance Data

K. Navatha¹, V V Hara Gopal²

¹Department of Statistics, Loyola Degree and PG College, Hyderabad, India
Email Id: navatha021[at]gmail.com

²Department of Statistics, Osmania University, Hyderabad - 7, India
Email: haragopal_vajjha[at]yahoo.com

Abstract: *Uncertainty refers to the randomness and is different from a lack of predictability or market inefficiency. An emergent research view holds that financial markets are both uncertain and predictable. Also, markets can be efficient but also uncertain. Insurance companies typically face two major problems when they want to forecast future premiums paid by using past or present behavior of premiums paid. For this, one has to find an appropriate statistical Probability distribution for the premiums paid. Then after test how well this statistical distribution fits the claims data. In modeling insurance claims, when there are extreme observations in the data, the commonly used loss distributions often are able to fit the bulk of the data well but fail to do so at the tail. One approach to overcome this problem is to focus on the extreme observations only and model them with the generalized Pareto distribution, supported by extreme value theory. The objective of this paper is to obtain an appropriate statistical Probability distribution for the insurance premium amounts and to test how well the chosen statistical distribution fits the premiums data. The modeling process will ascertain a statistical distribution that could capable model the claim amounts, and then the goodness of fit test was done mathematically using graphically using the Probability - Probability Plots (P - P plots) and Quantile - Quantile Plots (Q - Q plots). Finally, the study gives a summary, conclusion and recommendations that can be used by insurance companies to improve their results concerning future premium inferences.*

Keywords: premiums, extreme value theory, generalized Pareto distribution, generalized extreme value distribution, P - P plot, Q - Q Plot

1. Introduction

This Paper basically discusses the various distributions of premiums paid from an insurance company in India for the year 2009 - 2010. Here the data contains 50, 000 observations and 38 variables like age, gender, Type of policy, Date of birth, Type of disease, Bonus, Floater amount, sum insured, premium, claims paid etc. we are interested in fitting of distribution for the variable Premiums paid. The paper used exploratory data analysis (histogram, mean, variance skewness, kurtosis maximum value, minimum value, standard deviation, and 1st and 3rd quartile) to help in the identification of the family of distribution which the data might follow. Probability plot was used to graphically demonstrate goodness of fit to different distributions. Different Goodness - of - Fit tests were used to test fitness of the distributions.

1) Descriptive Statistics

A total of 156 Premiums paid data on Health insurance for 2009 - 2010 was used for the modelling. Table 2.1 below summarizes the result of the descriptive data analysis of the premiums paid

Table 2.1

Statistic	Value	Percentile	Value
Sample Size	156	Min	1.00E+05
Range	1.84E+07	5%	1.08E+05
Mean	1.08E+06	10%	1.20E+05
Variance	5.74E+12	25% (Q1)	2.01E+05
Std. Deviation	2.40E+06	50% (Median)	3.90E+05
Coef. of Variation	2.2242	75% (Q3)	8.79E+05
Std. Error	1.92E+05	90%	2.02E+06
Skewness	5.3031	95%	4.40E+06
Kurtosis	32.055	Max	1.85E+07

The average of premium was computed. The standard deviation, skewness, minimum and maximum value as well as the quantiles are also shown. This summary was necessary because it helps to identify key features of the data.: Descriptive statistics for the Health claim data From table 2.1 above summarizes the result of the descriptive data analysis of the Premium data.1, 00, 000 is the minimum Premium amount that was paid. The maximum Premium amount is 1, 84, 697. This indicates that within that period, the highest premium amount paid by the policy holder to the insurer. The 25th quartile of premium paid was 2, 01, 155 and that of the 75th quartile was 8, 79, 165. The mean premium paid was 1, 07, 726 and the variance was 5, 74, 012 and The standard deviation is 1, 91, 820 and that of the coefficient of skewness is 5.3031. The skewness was measuring the symmetric nature of the claim. The value 5.3031 informs how the premium amount was positively skewed. Kurtosis has a value of 9.3289, kurtosis measures whether the data is heavy - tailed or light tailed. The value 32.055 indicates that the data heavy tailed also and Leptokurtic curve. Before choosing one or more models for the data, it is necessary to choose good model among a predetermine set of models. This can be done with the help of fitting different distributions such as Gamma, Normal distribution, Log Normal, Weibull other distributions were chosen as a family of models for the study.

Before fitting any distribution one has to know the parameters. These are estimated by Statistical Software.

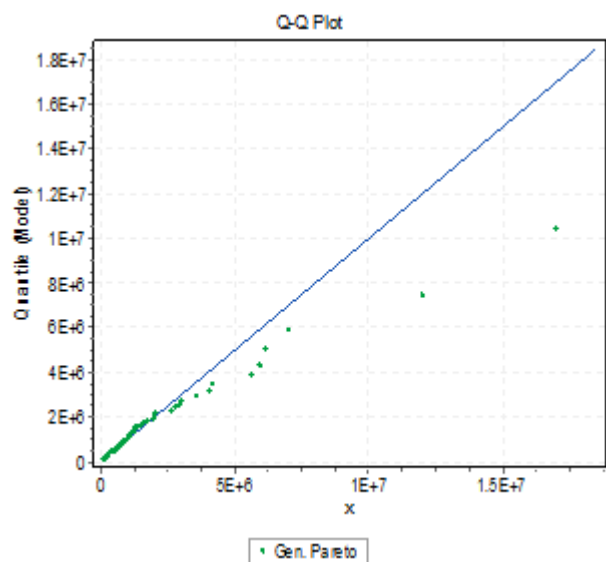
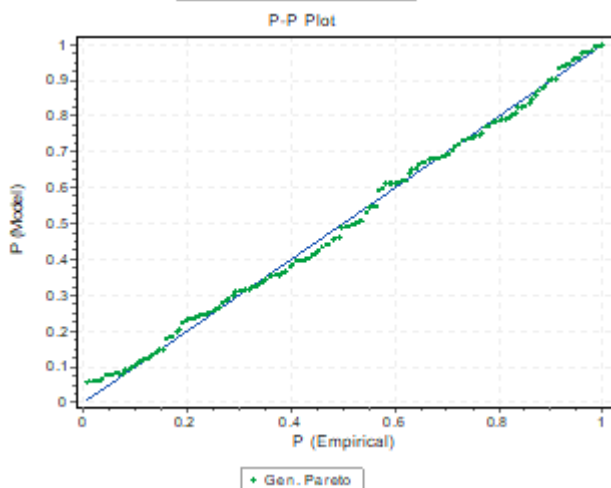
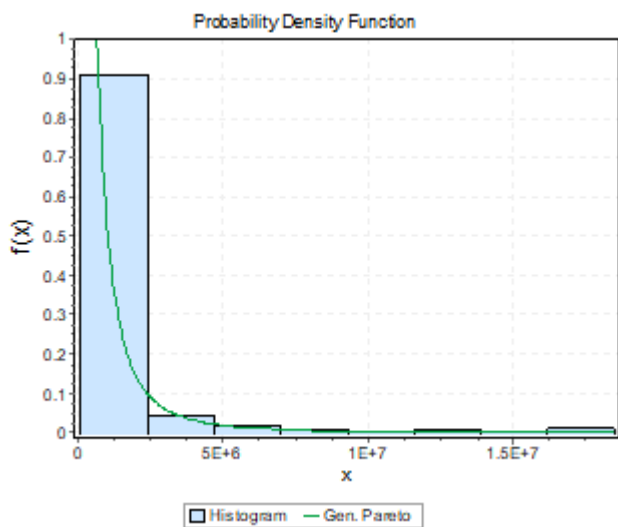
After fitting distributions, goodness of fit by different tests. with this one can know that which distribution fits well for the claim data. Not only manual we can test graphically also. that is done with of P - P plot and Q - Q plot Histogram.

P - P plot and Q - Q plots for Generalized Pareto Distribution is shown below.

Volume 11 Issue 11, November 2022

www.ijsr.net

Licensed Under Creative Commons Attribution CC BY



2) The Parameters of all the distributions fitted to the data are listed in the are listed in the following

Table

S. no	Distribution	Parameters
1	Exponential	$\lambda=9.2835E - 7$
2	Exponential (2P)	$\lambda=1.0234E - 6$ $a=1.0000E+5$
3	Fatigue Life	$a=1.3494$ $b=5.9898E+5$
4	Fatigue Life (3P)	$\lambda=1.9977$ $a=3.4076E+5$ $b=85908.0$
5	Gamma	$k=0.20214$ $a=5.3288E+6$
6	Gamma (3P)	$k=0.43725$ $a=2.2433E+6$ $b=1.0000E+5$
7	Gen. Gamma	$k=1.2243$ $a0.34906$ $b=5.3288E+6$ $k=0.82692$ $a=0.46809$
8	Gen. Gamma (4P)	$b=1.5858E+6$ $c=1.0000E+5$
9	Gen. Pareto	$k=0.62466$ $a=3.7472E+5$ $b=78809.0$
10	Log - Gamma	$a=136.53$ $b=0.09562$
11	Lognormal	$a=1.1137$ $b=13.055$
12	Lognormal (3P)	$k=1.6834$ $a=12.5$ $b=96619.0$
13	Normal	$a=2.3958E+6$ $b=1.0772E+6$
14	Pareto	$a=0.64853$ $b=1.0000E+5$
15	Pareto 2	$a=2.076$ $b=1.1007E+6$

k - threshold, a - Shape parameter and b - scale parameter

3) Goodness of Fit – Summary

After treating the data fits with various distributions, we now test these with KS - test, Anderson Darling and Chi - squared test and found that, the Generalized Pareto Distribution fits well with all the three tests for the claim data which is listed under Table 2.3, Serial Number 9 in the table below which ranks “1” under the three tests. Therefore, one can conclude that, Premiums paid follows Generalized Pareto Distribution with authenticity and while Log Normal and Log Gamma falls under second and so on.

SNo	Distribution	Kolmogorov Smirnov		Anderson Darling		Chi - Squared	
		Statistic	Rank	Statistic	Rank	Statistic	Rank
1	Exponential	0.22317	11	15.514	11	42.67	11
2	Exponential (2P)	0.26812	12	29.335	13	59.385	12
3	Fatigue Life	0.15115	8	7.254	9	21.042	8
4	Fatigue Life (3P)	0.10811	5	2.1152	5	20.277	5
5	Gamma	0.48636	16	38.12	15	179.17	15
6	Gamma (3P)	0.12746	6	5.7533	7	24.312	9
7	Gen. Gamma	0.27377	13	20.019	12	60.569	13

8	Gen. Gamma (4P)	0.15929	9	7.2611	10	26.057	10
9	Gen. Pareto	0.05406	1	0.55121	1	4.2852	1
10	Log - Gamma	0.07931	3	1.5886	3	7.9485	2
11	Lognormal	0.08788	4	2.0716	4	11.564	3
12	Lognormal (3P)	0.06082	2	0.72204	2	14.465	4
13	Normal	0.34169	14	33.099	14	72.663	14
14	Pareto	0.13543	7	6.0151	8	20.477	6
15	Pareto 2	0.16517	10	4.0693	6	20.661	7

2. Conclusions and Future Study

With this empirical study of the claim data suggest that the premium paid for the beneficiaries follows the Log Normal Distribution, Thus, with this we can estimate the number of persons paying premium above any premiums paid by using the above Log Normal Distribution. Also, from this we know that estimation is possible by the Log Normal

Distribution that how many are paying premium below or above 3 lakhs as compensation from the data on premium paid, This study can further be studied for various data sets on actuaries which will be helpful to process the premium and other parameters of the data. Thus, we can explore this procedure of fitting and estimating the parameters for other than premiums paid data for different set of variables such as claimspaid, floater amount, Sum assured, etc.