# Data Analysis and Modelingof Claim Amounts of Car Insurance using Big Data: A Study for Pakistan

**S. M. Aqil Burney [a*], Laiq Muhammad Khan [a], Shumaila Burney [a] and Muhammad Humayoun [a]**

[a] *College of Computer Science and Information Systems, Institute of Business Management, Karachi, Pakistan.*

*Authors' contributions*

*This work was carried out in collaboration among all authors. All authors read and approved the final manuscript.*

*Original Research Article*

## Abstract

Modelling of data of claim amount is of paramount importance to manage risk reserve for payment of claims. Actuaries model uncertainty using probability distributions.

In this research paper claim amount distribution of the data of an insurance concern has been estimated and analysis was performed on big-data of claim amounts for better understanding and fitting of various probability distribution using R.

It was noticed that the claim amounts distribution is highly positive skewed, therefore we have studied Exponential distribution, Gamma distribution and Weibull distribution as possible candidates for modelling the claim amount data. Chi Square test has been used as goodness of fit technique to decide suitable statistical model to representing the claim amounts under study.

Exponential distribution is found suitable for modelling the data under study.

Proposed model is usefulto estimate claim amount on aggregate for insurance concern when total loss is required to be computed to manage the risk reserve for the payments of claims.

_____

*Corresponding author: Email: aqil.burney@iobm.edu.pk;*

# 1 Introduction

Modeling of claim amount is veryimportant and is of great interest for actuaries. Actuaries measure the degree of uncertainty on the basis of models. It is used to solve many problems in actuarial science as well as predicting insurance cost. The model could be used to decide when a claim be made and how much be paid [1]. Therefore, modeling of claim amount is an important technique for actuaries to estimate parameters of the data foe the proposed model and making decision for loses and premium calculation's [2].

Study of claim amount pattern using probability distribution approach when relevant data is available is an important technique foresting price of insurance policies, in order to estimate the liabilities of insurance companies. Modeling of claims amount and frequency can be used for better understanding of the implications of claims to the solvency of the company [3].

In order to study information about claim amount in insurancecompanies, which is very important for making decisions about premium levels, estimation of reservesobtained from premium and the profitability of insurance portfolios, loss modelling of claim amount plays an important role [4]. Claim amount collected in insurance arepositivelyskewed, therefore, probability distributions exhibit this characteristic are used for modeling [5].

Most commonlyused rightly skewed distributions in actuaries for modeling claim amount are Gamma, Lognormal Weibull and Pareto distribution, Beta, Pareto, Burr, Weibull, Lognormal [6], Normal distribution and many other distributions [7].

# 2 Some Claim Amount Distributions

Below aregivenfew claim amount distributions used for modeling [8] of claim amounts and their cumulative density functions in Table 1. Commonly used measures of central tendency and dispersion of each distribution are also given for reference as these are required to study useful quantities in insurance in Table 2.

**Table 1. Probability density functions and cumulative density functions**

| Probability distribution | Probability density function | Cumulative density function |
|---|---|---|
| Exponential distribution | $f(x) = \frac{1}{\beta} e^{-\frac{x}{\beta}}$ where $0 \le x < \infty$, and $\beta > 0$ | $F(x) = \int_0^x \frac{1}{\beta} e^{\frac{-x}{\beta}} \, dx$ <br> $F(x) = 1 - e^{\frac{-x}{\beta}}$ where $x \ge 0$, $\beta > 0$ |
| Gamma distribution | $f(x) = \frac{\beta^\alpha}{\Gamma\alpha} x^{\alpha-1} e^{-\beta x}$ <br> where $0 < x < \infty$ | $F(x) = \int_0^x \frac{\beta^\alpha}{\Gamma\alpha} x^{\alpha-1} e^{-\beta x}$ |
| Weibull Distribution | $f(x) = \frac{\gamma}{\alpha}\left(\frac{x-\mu}{\alpha}\right)^{\gamma-1} e^{-\left(\frac{(x-\mu)}{\alpha}\right)^\gamma}$ $x \ge 0$; $\mu, \gamma, \alpha > 0$ | $F(x) = \int_0^x \gamma x^{(\gamma-1)} e^{-x^\gamma} dx$ <br> $F(x) = 1 - e^{-x^\gamma}$ $x \ge 0$; $\gamma > 0$ |

**Table 2. Measures of exponential, gamma and weibull probability distributions**

| Measures | Exponential | Gamma | Weibull |
|---|---|---|---|
| Mean | β | $\frac{\alpha}{\beta}$ | $\Gamma\left(\frac{\gamma+1}{\gamma}\right)$ |
| Mode | 0 | $\frac{\alpha-1}{\beta}$ | $\left(1-\frac{1}{\gamma}\right)^{\frac{1}{\gamma}}$ |
| Range | 0 to ∞ | 0 to ∞ | 0 to ∞ |
| Standard deviation | β | $\frac{\sqrt{\alpha}}{\beta}$ | $\sqrt{\Gamma\left(\frac{\gamma+2}{\gamma}\right) - \left(\Gamma\left(\frac{\gamma+1}{\gamma}\right)\right)^2}$ |

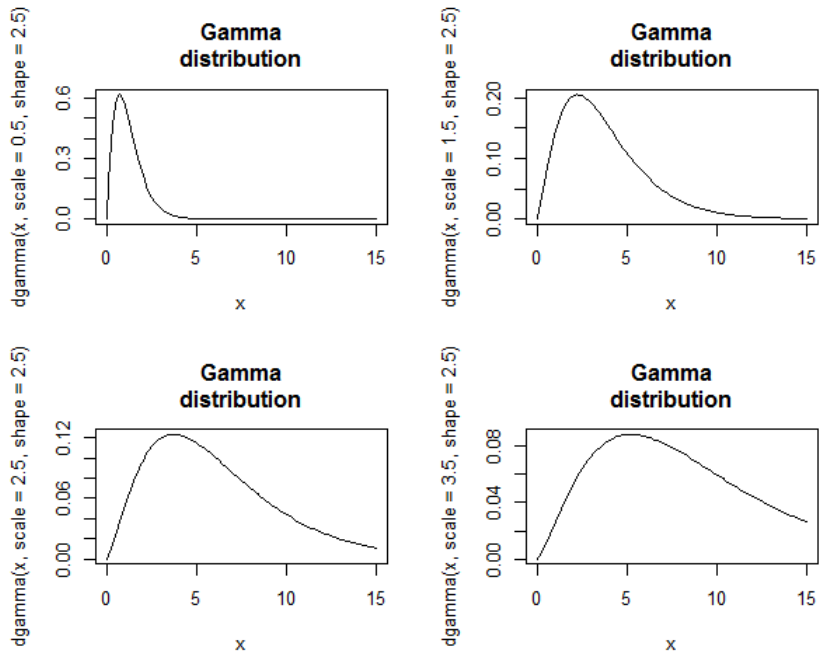| Measures | Exponential | Gamma | Weibull |
|---|---|---|---|
| **Coefficient of variation** | 1 | $\dfrac{1}{\sqrt{\alpha}}$ | $\sqrt{\dfrac{\Gamma\left(\frac{\gamma+2}{\gamma}\right)}{\left(\Gamma\left(\frac{\gamma+1}{\gamma}\right)\right)^2} - 1}$ |



**Fig. 1. Plot of Gamma distribution with different parameters**
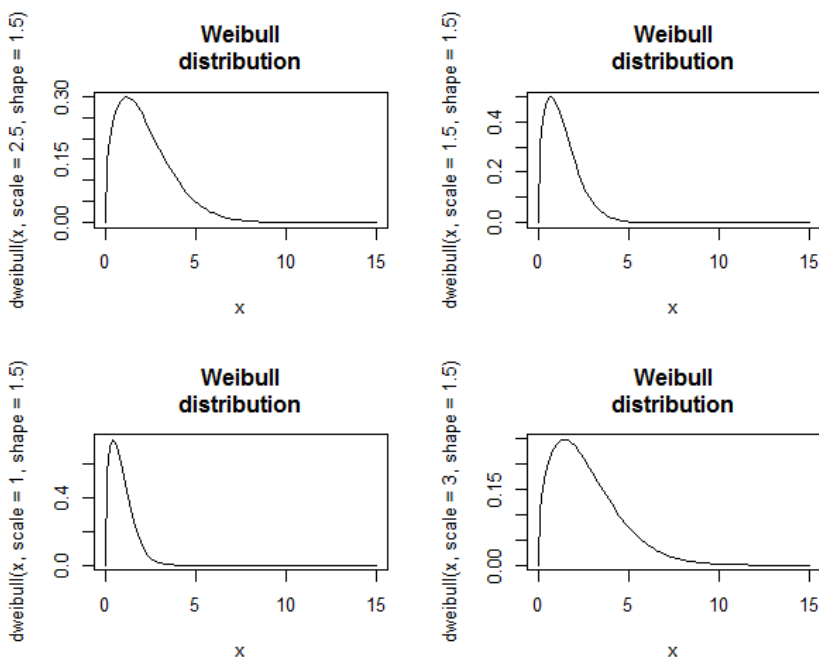


**Fig. 2. Plot of Weibull distribution with different parameters**

# 3 Literature Review

The claim amount data collected from financial records of a state owned major general insurance company in Iran. They fitted Skew Normal, Skew Laplace, generalized logistic, generalized hyperbolic, variance Gamma, normal inverse Gaussian, Marshal- Olkin, Log-Logistic, and Kumaraswamy Marshal-Olkinlog-Logistic distributions to the collected data and concluded that Kumaraswamy Marshal-Olkin log-Logistic distributions is better for modelling the claim amount [9].

The methodological frame work to select models for representing data of claim frequency and claim amount in insurance. They carried out this study on the basis of built-in data from R package insurance data. In this study researchers found lognormal distribution better model for representing claim amounts and negative binomial and geometric distribution as better distributions for modeling claim frequency [4].

The secondary data of claim amounts obtained from certain Insurance company in Nairobi, regarding their motor comprehensive policy. These researchers fitted Exponential, Gamma, Weibull and lognormal probability distributions and concluded that lognormal distribution is suitable for modeling the data under study [10].

Mohamed et al studied a model for claim amounts based on simulation of claim amounts. They fitted different probability distributions for claim amounts and found thatParetodistribution is suitable. This model was used for estimating insurance premiums for retention limit [11].

Burney and Hashmi have discussed different claim amount distributions as well as selection methods of distribution functions for claim amounts [12].

Talangtam, et al studied in order to model the data set of claim amounts of motor insurance using finite mixture lognormal distributions, and estimating parameters by EM algorithm. To decide best fitted model Kolmogorov Smirnoff (K-S) and A-D tests were used [1].

Meyres studied the data of 250 claims to decide suitable statistical probability distribution which could be used for modeling the data of claim amounts. The researcher fitted Gamma, Weibull and lognormal probability distributions to the data of claim amounts under study. The parameters of the fitted distribution were estimated by the method of maximum likelihood [13].

# 4 Methodology

## 4.1 Data analysis

The data for this research based on 133, 255 claimed amounts in Pk. Rupees. for the period Jan. 2010 to Dec 2015 from a well-known car insurance company in Pakistan. The name of company is not being mentioned here for confidentiality. The claim amounts of all types of vehicle showedan average of 29238.

The histogram of the data under study is drawn to decide the suitable probability distribution likely to fit the data.

# 5 Results and Discussion

Histogram shows highly skewed behavior. It can be seen that most of the claim amounts are small and there are few very high claims amounts. Therefore, it was decided to use highly skewed probability distributions for modeling the given data. Among rightly skewed continuous probability distributions, we used Exponential, Gammaand Weibull probability distributions for modeling the data under study using R package.

For detail estimation procedure using M.L.E and relation of exponential distribution with gamma distribution and Weibull distribution see (Saraless Nadarjah & Firoozeh Haghighi 2011) [14] thus in Table 4, we have computed expected frequency for exponential and gamma distribution.
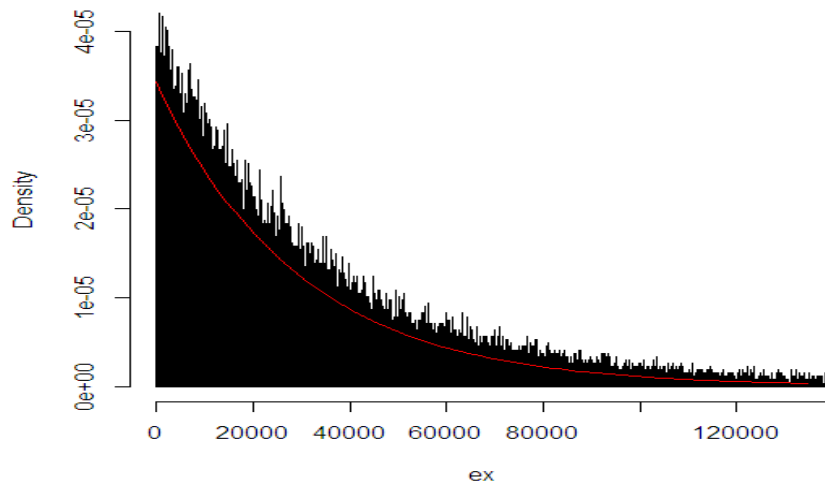
Following are the hypotheses to be tested about the modeling of claim amount data.

1. $H_0^*$: Exponential distribution is suitable candidate for modeling the data of claim amount under study.
2. $H_0^{**}$: Gamma distribution is suitable candidate for modeling the data of claim amount under study.
3 $H_0^{***}$: Weibull distribution is suitable candidate for modeling the data of claim amount under study.



**Fig. 3. Histogram of claim amounts for period Jan. 2010 to Dec 2015 in Pakistan**

Maximum Likelihood Estimates of the parameters of fitted distributions are given in Table 3.

**Table 3. Estimates of parameters of exponential and gamma distribution**

| Probability distribution | Estimates | |
|---|---|---|
| Exponential | $\lambda = 0.0000342$ | |
| Gamma | $\alpha = 0.00789$ | $\beta = 0.0000027$ |

**Table 4. Expected frequencies for exponential and gamma distribution**

| Claim Amount In 100,000 (PKR) | Observed Frequency | Expected Frequency (Exponential Distribution) | Expected Frequency (Gamma Distribution) |
|---|---|---|---|
| 0-1 | 129378 | 66257 | 122702 |
| 1-2 | 1618 | 33249 | 4879 |
| 2-3 | 514 | 16659 | 2276 |
| 3-4 | 341 | 8419 | 1257 |
| 4-5 | 229 | 4374 | 790 |
| 5-6 | 242 | 2165 | 486 |
| 6-7 | 231 | 1018 | 269 |
| 7-8 | 102 | 566 | 184 |
| 8-9 | 81 | 274 | 132 |
| 9-10 | 62 | 145 | 76 |
| 10-11 | 54 | 62 | 65 |
| 11-12 | 43 | 42 | 44 |
| 12-13 | 62 | 12 | 24 |
| 13-14 | 64 | 9 | 32 |
| 14-15 | 76 | 2 | 15 |
| 15-16 | 48 | 0 | 11 |
| 16-17 | 51 | 1 | 11 |
| 17-18 | 36 | 0 | 3 |
| 18-19 | 23 | 1 | 1 |

**Fig. 4. Graph of original data along with fitted distribution**

### 5.1 Goodness of fit test for claim amount distribution

Chi Square test was used for goodness of fit of the assumed probability distributions Chi Square values and p values are presented in Table 5. Observed frequencies and expected frequencies of Exponential and Gamma are presented in Table 4.

**Table 5. Chi square test statistic**

|  | Chi Square statistic | p-values |
|---|---|---|
| Exponential Distribution | $\chi^2 = 18.75859$, | 0.2813803 |
| Gamma Distribution | $\chi^2 = 43.02885$ | 0.006871283 |

For Table 5, it is observed that the p value for Gamma distribution is very small (p <0.01) which indicates that Gamma distribution is not a suitable candidate for modeling the data of claim amount under study. Exponential distribution on the other hand has large p value. Therefore, Exponential distribution is better candidate for modeling the claim amount data under study. We also fitted Weibull distribution and the relevant estimates are given in table 6, analysis of the fitting is as under:

**Table 6. Fitting of the distribution of Weibull distribution by the method of MLE**

|  | Estimate | Standard Error |
|---|---|---|
| Shape | 7.704193e-01 | 0.001263357 |
| Scale | 2.273621e+04 | 60.443552362 |
| Log likelihood: -1488209 | AIC: 2976422 | BIC: 2976442 |

Large value of Log Likelihood and small values of AIC and BIC indicates better fit(see Achieng, O. M., TRACK: ASTIN). But in our case value of LL is very small and values of AIC and BIC are very large which indicates that Weibull distribution is not suitable for the data of claim amounts under study.

Hence it is concluded that Exponential distribution as compare to Gamma and Weibull distributions is suitable probability distribution for modeling claim amount data.

## 6 Conclusion

The objective of this study was to decide suitable probability distribution among the three skewed probability distributions. Findings on the basis of empirical analysis of the data indicate that exponential distribution is a suitable statistical model for claim amounts.

# 7 Recommendations

Insurance companies need an accurate pricing system which could make sufficient space for estimation of contingencies, expenses, losses and profits. In result of occurrence of claim there is loss on part of insurance company, therefore estimation of losses which are likely to occur in future is very important for insurance companies. Estimation of such losses is not possible without modeling of the data of claim amount or losses.

This research provides a basis for car insurance companies to develop suitable models for the data of claim amount of respective companies. It is suggested to the insurance companies should make necessary adjustments in the probability distributions on the basis of their own claim amount data.

# Disclaimer

This manuscript was presented in a Conference.
**Conference name:** 11th International Conference on, Mathematics and Statistics Computer Science Actuarial Science Oct 27-28 2017At: IoBM Karachi.
**Available link:**
https://www.researchgate.net/publication/359199984_Data_Analysis_and_Modeling_of_Claim_Amounts_of_Car_Insurance_using_Big_Data_A_Study_for_Pakistan.

# Competing Interests

Authors have declared that no competing interests exist.

# References

[1]    R. "Fitting of finite mixture distributions to motor insurance claims. J. Math. Stat. 2012;8(1):49–56.
DOI: 10.3844/jmssp.2012.49.56

[2]    Rafal J. An empirical comparison of alternate regime-switching models or electricity spot prices. Munich Pers. RePEc Arch; 2010, [Online].
Available:https://mpra.ub.uni-muenchen.de/id/eprint/20546

[3]    Frees EW, Valdez EA. Hierarchical insurance claims modeling. J. Am. Stat. Assoc. 2008;103(484):1457–1469.
DOI: 10.1198/016214508000000823

[4]    Omari CO, Nyambura SG, Mwangi JMW. Modeling the frequency and severity of auto insurance claims using statistical distributions. J. Math. Finance. 2018;08(01):137–160.
DOI: 10.4236/jmf.2018.81012

[5]    COTOR Challenge Round 2. [Online Video].
Available:www.casact.org/cotor

[6]    Klugman SA, Panjer HH, Willmot GE. Loss models: From data to decisions, 1st ed. Wiley; 2008.
DOI: 10.1002/9780470391341

[7]    Klugman SA, Panjer HH, Willmot GE. Loss models: from data to decisions, Fifth edition. Hoboken, NJ: Wiley; 2019.

[8]    Boland PJ. Statistical methods in general insurance. Ireland, Dublin: Boland National University of Ireland, Dublin; 2006. [Online].
Available:https://www.stat.auckland.ac.nz/~iase/publications/17/5G1_BOLA.pdf

[9]     Akram Kohansa RK. Fitting skew distributions to iranian auto insurance claim data. Appl. Appl. Math. 2017;12(2):790–803.

[10]    Nairobi Kenya OMA. Actuarial modeling for insurance claim severity in motor comprehensive policy using industrial statistical distributions.

[11]    Mohamed. Approximation of aggregate losses using simulation. J. Math. Stat. 2010;6(3):233–239. DOI: 10.3844/jmssp.2010.233.239.

[12]    Aqil Burney SM. Risk theory and insurance an invited lecture. [Online]. Available:https://www.researchgate.net/publication/342624357_Risk_Theory_and_Insurance_an_Invited _Lecture

[13]    Vito Ricci, Fitting Distribution with R. Free software foundation; 2005. [Online]. Available:https://www.researchgate.net/publication/228791072_Fitting_Distributions_with_R

[14]    Nadarajah S, Haghighi F. An extension of the exponential distribution. Statistics. 2011;45(6):543–558. DOI: 10.1080/02331881003678678

_____