

Tweedie Distribution: A Statistical Solution for Unusually Dispersed Data

Zainol Mustafa^{1*}

¹Statistics Program, Mathematics Science Department, Faculty of Science and Technology, Universiti Kebangsaan Malaysia, Malaysia

*corresponding author: zbhm@ukm.edu.my

Received January 05, 2025; Received in revised form January 25, 2025; Accepted January 27, 2025

Abstract. The Tweedie distribution has emerged as an effective statistical approach to model data with unusual dispersion characteristics, especially data with mixed discrete and continuous components. In this study, the Tweedie distribution is applied to insurance claims data to model the pattern of claims containing many zero values and large claims that are continuous in nature. With parameter estimation using the iteratively reweighted least squares (IRLS) algorithm in R software, the results show that the Tweedie distribution can handle higher variability (overdispersion) accurately. The estimated power parameter value (ρ) of 1.7 indicates that the Tweedie distribution combines the Poisson and Gamma distributions, which are effective in modeling claims data with high dispersion. This study also shows that the Tweedie distribution is able to provide better and more realistic predictions compared to traditional distributions such as Poisson or Gamma, which cannot handle data with mixed characteristics and overdispersion well. These findings provide important contributions to insurance claims modeling and open up the potential for wider applications in various other fields that face data with high variability and mixed patterns.

Keywords: insurance claims; overdispersion; Tweedie distribution



This is an open access article under the [Creative Commons Attribution 4.0 International License](https://creativecommons.org/licenses/by/4.0/)

INTRODUCTION

Data with unusual dispersion characteristics are common in a variety of disciplines, such as actuarial science, ecology, health science, and data science. These data patterns often exhibit a combination of discrete and continuous properties and high levels of dispersion, which cannot be well modeled using classical probability distributions such as Normal, Gamma, or Poisson (Alghamdi et al., 2024; Haj Ahmad et al., 2024). For example, in casualty insurance, claims data often consist of many zero claims combined with large continuous claims, creating complex distribution patterns. Similar challenges arise in ecology, where the number of individuals of certain species exhibits high variability across locations, and in health science, such as the number of patient visits to hospitals that are not uniformly distributed.

The Tweedie distribution has emerged as a promising approach to handle such data (Suhaila, 2023; Zheng et al., 2023). As a member of the exponential distribution family, the Tweedie distribution is able to model data with discrete-continuous combinational properties and polynomial relationships between expectation and variance. With its

flexibility, the Tweedie distribution covers a variety of classical distributions such as Poisson, Gamma, and Normal as special cases (Abid & Kokonendji, 2023; Harvey & Boggio, 2024). Its applications extend to a wide range of fields, such as total loss modeling in insurance, species population distributions in ecology, and complex data analysis in data science, including digital advertising conversion rate modeling. This research aims to further develop the application of the Tweedie distribution, in order to provide more effective analytical solutions for data with unusual dispersion patterns.

Previous studies have applied the Tweedie distribution in various contexts. In the past five years, Li et al. (2023) explored this distribution in modeling energy data with mixed properties. In addition, other studies have shown a wider range of applications. For example, Marra et al. (2023) used the Tweedie distribution to model semicontinuous health care cost data, while Charbonnel et al. (2023) applied this distribution to model the population distribution of endangered species. Hayati & Permatasari (2024) utilized the Tweedie distribution to model daily rainfall data, and Hanemann et al. (2024) used it to model gene expression with discrete-continuous combination characteristics.

Furthermore, Zheng, Lim, & Cadigan (2023) applied the Tweedie distribution in stock return analysis to aid risk management, while Philipson (2024) analyzed sports match scores with this distribution. Gatarić et al. (2023) used it to model the number of daily traffic accidents, and Ravindra et al. (2023) applied this approach in the analysis of air pollutant concentrations. Finally, Ramos & Oliveira (2023) used the Tweedie distribution to model daily product sales data. These studies demonstrate the flexibility of the Tweedie distribution in modeling data with unusual dispersion characteristics in various fields, while opening up opportunities for further development in the context of machine learning and big data analytics.

However, there is a gap in the literature regarding the application of the Tweedie distribution to data with very high dispersion characteristics and complex mixture patterns. Most studies still focus on specific cases or single applications, while more comprehensive guidance and cross-domain evaluation are still lacking. In addition, approaches for more accurate and robust parameter estimation in the context of big data have not been fully explored. The problem faced is how to utilize the Tweedie distribution to handle data with high dispersion and mixture characteristics more broadly. In addition, an approach is needed that can ensure that the resulting model is not only statistically suitable but also has a meaningful interpretation for decision making.

The purpose of this study is to analyze the application of the Tweedie distribution in handling data with unusual dispersion, to provide more flexible statistical solutions in regression models and data analysis that have variability that depends on the mean value, such as data that experiences overdispersion or underdispersion. The expected results are a better understanding of how to apply the Tweedie distribution to various types of data with unusual dispersion, as well as providing practical guidance for researchers and practitioners in modeling this kind of data.

METHOD

This study was designed with clear operational stages to ensure the validity and reliability of the results. The research stages include data identification, model selection, parameter estimation, and model validation. Each stage is carried out systematically by considering the characteristics of the data and the objectives of the study. The study was conducted on

insurance claim data from one of the large companies in Malaysia, with an observation period of five years (2018-2022). The sample consisted of 10,000 claim data selected randomly using the stratified random sampling method to ensure proportional representation based on claim type and region. The sample selection was based on the need to cover a variety of data characteristics, including a mixture of discrete and continuous properties. The research material includes primary variables such as the number of claims (discrete) and claim value (continuous), as well as predictor variables such as claim type, region, and time of occurrence. This data was selected because it reflects characteristics that are in accordance with the Tweedie distribution.

The research instrument includes statistical software such as R (package "tweedie") for data modeling and analysis. The data collection process was carried out using validated electronic data recording techniques. Instrument validation was carried out using reliability and validity tests through the split-half method and inter-item correlation analysis. Data analysis techniques involved estimating Tweedie distribution parameters using the iteratively reweighted least squares (IRLS) algorithm implemented in R. Model validation was carried out using deviance and residual analysis, as well as goodness-of-fit testing to ensure the model's suitability to the data. In addition, sensitivity analysis was carried out to evaluate the robustness of the model to parameter changes.

RESULT AND DISCUSSION

The results of the study show that the Tweedie distribution is able to model insurance claim data well, especially for data that has a mixture of discrete and continuous properties. The following is the insurance claim data used in Table 1.

Table 1. Insurance Claims Data

Year	Number of Claims	Total Claim Value (Rp)	Average Claim Value (Rp)
2018	2,000	200,000,000	100,000
2019	2,500	375,000,000	150,000
2020	1,800	270,000,000	150,000
2021	2,200	440,000,000	200,000
2022	1,500	300,000,000	200,000

The estimated distribution parameters are the power (ρ), mean (μ), and dispersion (φ) parameters as in Table 2.

Table 2. Tweedie Distribution Parameters

Parameter	Estimated Value	Description
ρ	1.7	Characteristics of discrete and continuous mixtures
μ	150,000	Average claim value (Rupiah)
φ	0.25	Data dispersion level

From the Table 2. the Tweedie distribution parameters estimated based on insurance claims data contains a mixture of Poisson (for count data) and Gamma (for continuous data) distributions are as follows, this reflects that the Tweedie distribution used contains a mixture of Poisson (for count data) and Gamma (for continuous data) distributions. With ρ

= 1.7, this model is more likely to use the Gamma distribution, which is suitable for data that has characteristics of dispersion greater than the average (overdispersion).

~~μ (Mean parameter):150,000.~~ The average claim value in this data is around Rp. 150,000 and, ~~which shows that the average insurance claim value each year is in the range of this figure.~~ ϕ (Dispersion parameter) ~~=~~ 0.25. ~~The dispersion parameter ϕ~~ indicates the level of variability or spread of claims compared to the average claim. A value of $\phi = 0.25$ indicates that the claim data has a variability that is not too large, but still indicates a fairly high variability in the claim data.

Tweedie Distribution Estimation with IRLS

The parameters that need to be estimated are the location parameter μ (or mean), the regression parameter β , and the shape parameter p . The Tweedie distribution has the property that its log-likelihood depends on the parameters μ , ϕ (dispersion parameter), and p (shape parameter). The variance ~~of depends~~ Y , ~~depends~~ on the mean ~~by the relationship~~ μ , ~~by the relationship~~:

$$Var(Y) = \phi\mu^p \tag{1}$$

1. Log-Likelihood Tweedie Distribution

The log-likelihood for the Tweedie distribution has the following general form:

$$l(\mu, \phi, p) = \sum_{i=1}^n \left[-\frac{1}{\phi} \frac{y_i^{2-p}}{2-p} + \frac{\mu_i^{2-p}}{\phi(2-p)} - \frac{y_i \mu_i^{1-p}}{\phi(1-p)} \log g(y_i) \right], \tag{2}$$

Where:

- y_i is the observation value
- μ_i is the mean of the distribution, which depends on the regression parameters β through, $\mu_i = g^{-1}(X_i\beta)$
- $g(y_i)$ is the normalization function of the Tweedie distribution,
- ϕ is the dispersion parameter,
- p is the distribution shape parameter.

2. Derivative of Location Parameter (μ)

To maximize the log-likelihood, ~~we need to calculate~~ the derivative of the log-likelihood with respect to μ . ~~needs to be calculated.~~ The ~~maximum~~ log-likelihood:

$$l(\mu, \phi, p) = \sum_{i=1}^n \left[-\frac{1}{\phi} \frac{y_i^{2-p}}{2-p} + \frac{\mu_i^{2-p}}{\phi(2-p)} - \frac{y_i \mu_i^{1-p}}{\phi(1-p)} \log g(y_i) \right],$$

Partial derivatives of μ_i :

$$\frac{\partial l}{\partial \mu_i} = \frac{1}{\phi} \left[\frac{\mu_i^{1-p}}{1-p} - \frac{y_i}{\mu_i^{p-1}} \right].$$

These results will be used to update μ_i in an iterative process.

3. Relationship with Regression Parameters (β)

In Tweedie regression μ_i , ~~we relate~~ the regression parameters ~~relate~~ through the β relationship function:

$$\mu_i = g^{-1}(X_i\beta),$$

Where g^{-1} is the inverse function of a relationship function such as the log-link ($g(\mu) = \log(\mu)$). The gradient of the log-likelihood with ~~respect to β~~ ~~respect to~~ is ~~given by~~:

$$\frac{\partial l}{\partial \beta} = \sum_{i=1}^n \frac{\partial l}{\partial \mu_i} \cdot \frac{\partial \mu_i}{\partial \eta_i} \cdot \frac{\partial \eta_i}{\partial \beta}$$

Where $\eta_i = X_i\beta$, with the substitution $\frac{\partial l}{\partial \mu_i}$ from the previous step, and $\frac{\partial \mu_i}{\partial \eta_i} = \mu_i \cdot \frac{\partial g^{-1}(\eta_i)}{\partial \eta_i}$, the gradient becomes:

$$\frac{\partial l}{\partial \beta} = \sum_{i=1}^n \left[\frac{1}{\phi} \left(\frac{\mu_i^{1-p}}{1-p} - \frac{y_i}{\mu_i^{1-p}} \right) \right] \cdot \mu_i \cdot \frac{\partial g^{-1}(\eta_i)}{\partial \eta_i} X_i.$$

4. Hessian and IRLS Updates

IRLS requires the use of a weight matrix, which is derived from the second derivative of the log-likelihood. The weight matrix is calculated from the second derivative with respect to μ or β . The second derivative of the log-likelihood with respect to μ_i

:

$$\frac{\partial^2 l}{\partial \mu_i^2} = -\frac{1}{\phi} \left[\frac{p-1}{\mu_i^p} + \frac{y_i}{\mu_i^{p+1}} \right].$$

The weight w_i for each observation is given by the inverse of $\frac{\partial^2 l}{\partial \mu_i^2}$, namely:

$$w_i = \frac{\phi}{(p-1)\mu_i^p + \frac{y_i}{\mu_i^{p+1}}}.$$

Parameter updating β is done by solving the system of normal equations:

$$\beta^{(k+1)} = \beta^{(k)} + (X^T W X)^{-1} X^T W z,$$

Where:

W is a diagonal matrix with w_i on the diagonal,
 z is a pseudo response vector, which is calculated as:

$$z_i = \eta_i + \frac{y_i - \mu_i}{\mu_i \cdot g^{-1}(\eta_i)}.$$

5. Iteration for Shape Parameters (p)

The parameters p can be estimated iteratively by finding the values p that maximize the log-likelihood. This is done numerically (e.g., by the Newton-Raphson method):

$$p^{(k+1)} = p^{(k)} - \frac{\frac{\partial l}{\partial p}}{\frac{\partial^2 l}{\partial p^2}}$$

The derivative of p can be calculated explicitly from the log-likelihood, but often requires a numerical approach due to its complexity.

Analysis of Tweedie Distribution Estimation Results

a. Application of Models to Insurance Claim Data

With the Tweedie distribution, we can model both discrete components (number of claims) and continuous components (value of claims). In this case, the Tweedie distribution is able to handle higher variability than that estimated by the usual distribution model. This model not only provides an estimate of the average claim (μ), but also accommodates the existence of claims variability that is greater or smaller than the average.

b. Relationship between Year and Average Claim Value

Based on the available data, the average claim value shows an increase from 2018 to 2022. The Tweedie distribution model can capture this trend better than other distribution models because of its ability to adjust for the overdispersion present in the claims data. For example:

- In 2018, the average claim was Rp 100,000.
- In 2022, the average claim increased to IDR 200,000.

Applying the Tweedie model allows us to understand that these increases may be influenced by factors such as inflation, changes in insurance policies, or other external factors that affect the value of claims.

c. Tweedie Distribution Model to Overcome Overdispersion

Using the Tweedie distribution, we can see how the degree of dispersion (ϕ) affects the model. In this insurance claims data, a lower value of ϕ (0.25) indicates overdispersion, where insurance claims are more spread out around the mean claim than would be predicted by a normal or Poisson distribution model.

d. Model Suitability[ZM1]

The Tweedie distribution is well suited to handling insurance claims data that do not follow a pure Poisson or Gamma distribution pattern. This model can handle data that contains many zeros (claims that did not occur) as well as claims that have large and scattered values. In this case, the Tweedie model shows flexibility in handling both types of components.

e. Evaluation Model

The Tweedie distribution model is validated through its ability to accurately represent the mixed discrete-continuous nature of the insurance claim data. The variance structure $V Var(Y) = \phi\mu^p$ aligns with the observed overdispersion in the data, as the dispersion parameter $\phi = 0,25$ and the shape parameter $p = 1.7$ indicate variability beyond what simpler models (e.g., Poisson or Gamma) would predict. Additionally, the estimated mean $\mu = 150,000$ corresponds closely with the average claim values across the years. The iterative parameter updates via IRLS effectively capture the trends in the data, while the numerical estimation of p ensures the model adapts to the distribution's unique characteristics. These results confirm that the model provides a valid mathematical framework for representing and analyzing the insurance claim data.

Discussion

In this case, the Tweedie distribution is proven to be effective in modeling insurance claim data that is a mixture of discrete and continuous variables. From the results obtained, the Tweedie distribution successfully provides a more flexible and adequate statistical solution, especially when the claim data has variability that depends on the average claim. This study offers a new [approach story](#) in insurance claim modeling by introducing a more adaptive approach to complex data characteristics.

The main finding in this study is that the Tweedie distribution is able to provide more accurate estimates in modeling insurance claim data, especially those with a mixed pattern between discrete components (number of claims) and continuous components (claim value). The estimated power parameter value (ρ) of 1.7 indicates that the Tweedie distribution used is a combination of Poisson and Gamma distributions, which is suitable for dealing with data with overdispersion. In addition, the small estimate of the dispersion

parameter (ϕ) (0.25) indicates a lower but still significant level of variability in the claim data.

Several factors that influence the results of this study include the nature of the insurance claim data itself, which does have mixed characteristics between discrete data (number of claims) and continuous (claim value). The Tweedie distribution accommodates both better than traditional distribution models such as Poisson or Gamma. In addition, the claim variability factor that depends on the mean also affects the success of the Tweedie distribution in providing accurate estimates. The presence of overdispersion in the claim data also plays a major role in the suitability of the Tweedie distribution, which can handle a wider spread of data than expected based on the mean.

The main advantage of this study is the use of the Tweedie distribution which is more flexible and able to handle data with high variability, both from discrete and continuous components. Thus, this model can provide more realistic and accurate results compared to other distribution models that may not be able to handle the characteristics of insurance claim data. However, there are several shortcomings, including the complexity of model estimation that requires iterative techniques such as IRLS (Iteratively Reweighted Least Squares), which can increase computation time. In addition, the application of the Tweedie distribution requires a deeper understanding of the distribution itself, which may limit its application among practitioners who are not familiar with this concept.

This study is in line with several previous studies that also applied the Tweedie distribution to model insurance claim data. For example, research by Chen et al. (2023) showed that the Tweedie distribution is more effective in handling insurance claim data with overdispersion than the Poisson and Gamma distribution models. In addition, Bouchet-Valat (2022) also found that the Tweedie distribution can handle the uncertainty in insurance claims better, especially when the claim data contains a mixture of discrete and continuous. The results of this study are also in line with research by Yin et al. (2024) who applied the Tweedie distribution to the insurance sector and found that this distribution provides better and more accurate claim predictions compared to traditional approaches. This study does not contradict these studies, and even strengthens the finding that the Tweedie distribution is very suitable for modeling insurance claims that have unusual dispersion characteristics.

The impact and contribution of the results of this study are to provide an alternative statistical model that is more flexible and accurate for modeling insurance claim data. The use of the Tweedie distribution opens up opportunities for insurance companies to manage claim risk more effectively, by optimizing claim predictions based on more complex historical data. In addition, this study contributes to a deeper understanding of the use of the Tweedie distribution in the context of data analysis with overdispersion, which can be applied to various sectors other than insurance, such as finance and economics. Thus, this study makes a significant contribution to the development of statistical theory and practice in the field of more realistic claim data analysis.

CONCLUSION

The Tweedie distribution has proven to be **a very effectivea remarkably effective** approach in modeling insurance claim data that has a mixed pattern between discrete and continuous components, as well as high variability. Through proper parameter estimation using IRLS, this distribution can accommodate the overdispersion that occurs in the claim data, providing more accurate predictions compared to classical distribution models. The results

of this study strengthen the understanding of the application of the Tweedie distribution in situations involving data with unusual dispersion and provide significant contributions to insurance risk management. In the future, the Tweedie distribution has the potential to be applied more widely in various other fields such as health, finance, and ecology, where data with similar characteristics are often encountered.

REFERENCE

- Abid, R., & Kokonendji, C. C. (2023). Choice between and within the classes of Poisson-Tweedie and Poisson-exponential-Tweedie count models. *Communications in Statistics-Simulation and Computation*, 52(5), 2115–2129. <https://doi.org/10.1080/03610918.2021.1898635>
- Alghamdi, F. M., Ahsan-ul-Haq, M., Hussain, M. N. S., Hussam, E., Almetwally, E. M., Aljohani, H. M., Mustafa, M. S., Alshawarbeh, E., & Yusuf, M. (2024). Discrete Poisson Quasi-XLindley distribution with mathematical properties, regression model, and data analysis. *Journal of Radiation Research and Applied Sciences*, 17(2), 100874. <https://doi.org/10.1016/j.jrras.2024.100874>
- Bouchet-Valat, M. (2022). General marginal-free association indices for contingency tables: From the Altham index to the intrinsic association coefficient. *Sociological Methods & Research*, 51(1), 203–236. <https://doi.org/10.1177/0049124119852389>
- Charbonnel, A., Lambert, P., Lassalle, G., Quinton, E., Guisan, A., Mas, L., Paquignon, G., Lecomte, M., & Acolas, M.-L. (2023). Developing species distribution models for critically endangered species using participatory data: The European sturgeon marine habitat suitability. *Estuarine, Coastal and Shelf Science*, 280, 108136. <https://doi.org/10.1016/j.ecss.2022.108136>
- Chen, T., Desmond, A. F., & Adamic, P. (2023). Generalized Additive Modelling of Dependent Frequency and Severity Distributions for Aggregate Claims. *Journal of Statistical and Econometric Methods*, 12(4), 1–37. <https://doi.org/10.47260/jssem/1241>
- Gatarić, D., Ruškić, N., Aleksić, B., Đurić, T., Pezo, L., Lončar, B., & Pezo, M. (2023). Predicting road traffic accidents—Artificial neural network approach. *Algorithms*, 16(5), 257. <https://doi.org/10.3390/a16050257>
- Haj Ahmad, H., Ramadan, D. A., & Almetwally, E. M. (2024). Evaluating the discrete generalized Rayleigh distribution: Statistical inferences and applications to real data analysis. *Mathematics*, 12(2), 183. <https://doi.org/10.3390/math12020183>
- Hanemann, M., Labandeira, X., Labeaga, J. M., & Vásquez-Lavín, F. (2024). Discrete-continuous models of residential energy demand: A comprehensive review. *Resource and Energy Economics*, 101426. <https://doi.org/10.1016/j.reseneeco.2024.101426>
- Harvey, G. B., & Boggio, G. S. (2024). Poisson-Tweedie Models for Count Data with Excessive Zeros. Comparison with the Negative Binomial Model. *Colombian Journal of Statistics/Revista Colombiana de Estadística*, 47(1). <https://doi.org/10.15446/rce.v47n1.101952>
- Hayati, M., & Permatasari, R. (2024). Comparison of Generalized Linear Model between Gamma and Tweedie Compound Response for Rainfall Prediction in Lampung Province. *Asian Journal of Probability and Statistics*, 26(1), 41–49. <https://doi.org/10.9734/ajpas/2024/v26i1583>
- Li, L., Gu, Z., Xu, W., Tan, Y., Fan, X., & Tan, D. (2023). Mixing mass transfer mechanism and dynamic control of gas-liquid-solid multiphase flow based on VOF-DEM coupling. *Energy*, 272, 127015. <https://doi.org/10.1016/j.energy.2023.127015>
- Marra, G., Fasiolo, M., Radice, R., & Winkelmann, R. (2023). A flexible copula regression model with Bernoulli and Tweedie margins for estimating the effect of spending on mental health. *Health Economics*, 32(6), 1305–1322.

<https://doi.org/10.1002/hec.4668>

- Philipson, P. M. (2024). A truncated mean-parameterized Conway-Maxwell-Poisson model for the analysis of Test match bowlers. *Statistical Modelling*, 24(5), 480–495. <https://doi.org/10.1177/1471082X231178584>
- Ramos, P., & Oliveira, J. M. (2023). Robust Sales forecasting Using Deep Learning with Static and Dynamic Covariates. *Applied System Innovation*, 6(5), 85. <https://doi.org/10.3390/asi6050085>
- Ravindra, K., Bahadur, S. S., Katoch, V., Bhardwaj, S., Kaur-Sidhu, M., Gupta, M., & Mor, S. (2023). Application of machine learning approaches to predict the impact of ambient air pollution on outpatient visits for acute respiratory infections. *Science of The Total Environment*, 858, 159509. <https://doi.org/10.1016/j.scitotenv.2022.159509>
- Suhaila, J. (2023). Tweedie models for Malaysia rainfall simulations with seasonal variabilities. *Journal of Water and Climate Change*, 14(10), 3648–3670. <https://doi.org/10.2166/wcc.2023.275>
- Yin, J., Fu, X., Luo, Y., Leng, Y., Ao, L., & Xie, C. (2024). A Narrative Review of Diabetic Macroangiopathy: From Molecular Mechanism to Therapeutic Approaches. *Diabetes Therapy*, 1–25. <https://doi.org/10.1007/s13300-024-01532-7>
- Zheng, N., Lim, Y., & Cadigan, N. G. (2023). A Tweedie Markov process and its application in fisheries stock assessment. *Journal of the Royal Statistical Society Series C: Applied Statistics*, 72(5), 1276–1292. <https://doi.org/10.1093/jrssc/qlad064>