

A New Flexible Discrete Distribution: Theory and Empirical Evidence

Giovanni Pollio, Giovanni De Luca
 The Parthenope University of Naples
 Italy

Abstract

A new discrete distribution with two parameters is introduced and discussed. We present its derivation after identifying the generator mechanism of the mass probability function. Then we show its flexibility in terms of dispersion index and show how to estimate the parameters by maximum likelihood. Finally, we compare it to traditional as well as flexible discrete distributions using some popular insurance datasets. The AIC and BIC criteria strongly suggest that the new distributions is able to provide a fit to discrete data in a very satisfactory way.

1. Introduction

The prediction of insurance claims is one of the most important problems faced in the insurance industry. Companies need statistical models able to provide a probability law to the number of claims based on a discrete distribution, given the discrete nature of claims. The choice of a discrete statistical distribution is a very critical point, because the sensitiveness of the prediction to the statistical model can be highly considerable. The Poisson and the Negative Binomial are the most important discrete distributions able to represent the number of claims. However, in many cases, these distributions do not provide satisfactory results, especially under the circumstance that a high number of zeros occurs. The search for alternative discrete distributions has enriched the literature in a remarkable way. The article is organized as follows. Section 2 presents the most popular discrete distributions, while in Section 3 we propose a novel discrete distribution which is theoretically compared to the most traditional alternatives. Finally, in Section 4 an application to some popular insurance data shows the relevance of the proposal. Some concluding remarks will close the article.

2. Discrete distributions

The analysis of the most important discrete random variables is carried out following a unified approach considering the convergent numerical series with positive terms as the generator mechanism of the mass probability function. In particular, this approach can shed light on the relationship between the probability laws for discrete random variables with infinite support and the numerical series.

Actually, a general probability function for a discrete random variable Y could be written as

$$P(Y = y) = \frac{f(y; \theta)}{\sum_{j=0}^{\infty} f(j; \theta)}$$

where θ is the parameter vector, possibly constrained to ensure that $f(j; \theta)$ be convergent with positive values, and f denotes a mathematical function.

2.1 Poisson distribution

The Poisson distribution is certainly the most popular discrete distribution. Its probability function is given by

$$P(Y = y) = \frac{f(y; \mu)}{\sum_{j=0}^{\infty} f(j; \mu)} = \frac{\frac{\mu^y}{y!}}{\sum_{j=0}^{\infty} \frac{\mu^j}{j!}} = \frac{e^{-\mu} \mu^y}{y!}$$

In this case, the series $\sum_{j=0}^{\infty} \frac{\mu^j}{j!}$ with $\mu > 0$ has positive values and it is straightforward to show that converges to e^μ . It is characterized by an expected value equal to the variance, so that the dispersion index D

$$D = \frac{\sigma^2}{\mu}$$

that is the ration between variance and expected value, is always unit. This constraint is a serious drawback in practical applications.

2.2 Poisson-Lindley distribution

The Poisson-Lindley distribution has been introduced by Sankaran (1970) as a discrete distribution for count data. It comes from the Poisson distribution with the parameter $\mu > 0$ following the Lindley distribution with density function

$$f(\mu; \theta) = \frac{\theta^2}{\theta + 1} (1 + \mu)e^{-\theta\mu}$$

with $\theta > 0$.

So it belongs to the class of mixed Poisson distributions.

It is interesting to note that an alternative genesis of this distribution is given by the following convergent series with positive values, function of the parameter $\theta > 0$:

$$\begin{aligned} \sum_{j=0}^{\infty} \frac{\theta + 2 + j}{(\theta + 1)^j} &= \sum_{j=0}^{\infty} \left(\frac{\theta}{(\theta + 1)^j} + \frac{2}{(\theta + 1)^j} + \frac{j}{(\theta + 1)^j} \right) = \\ &= \left[\theta \sum_{j=0}^{\infty} \frac{1}{(\theta + 1)^j} \right] + \left[2 \sum_{j=0}^{\infty} \frac{1}{(\theta + 1)^j} \right] + \left[\sum_{j=0}^{\infty} \frac{j}{(\theta + 1)^j} \right] = \\ &= \left(1 + \frac{1}{\theta} \right) \theta + 2 \left(1 + \frac{1}{\theta} \right) + \left(\frac{1}{\theta} + \frac{1}{\theta^2} \right) = \frac{\theta^3 + \theta^2 + 2\theta^2 + 2\theta + \theta + 1}{\theta^2} = \frac{(\theta + 1)^3}{\theta^2} \end{aligned}$$

Then,

$$P(Y = y) = \frac{f(y; \theta)}{\sum_{j=0}^{\infty} f(j; \theta)} = \frac{\frac{\theta + 2 + y}{(\theta + 1)^y}}{\sum_{j=0}^{\infty} \frac{\theta + 2 + j}{(\theta + 1)^j}} = \frac{\frac{\theta + 2 + y}{(\theta + 1)^y}}{\frac{(\theta + 1)^3}{\theta^2}} = \frac{\theta^2(\theta + 2 + y)}{(\theta + 1)^{y+3}}.$$

The dispersion index is given by

$$D = \frac{\theta^3 + 4\theta^2 + 6\theta + 2}{\theta^2(\theta + 1)^2} \cdot \frac{\theta(\theta + 1)}{\theta + 2} = \frac{\theta^3 + 4\theta^2 + 6\theta + 2}{\theta(\theta + 1)(\theta + 2)}$$

As a result, the index is never below 1, that is $D \geq 1, \forall \theta > 0$, so the Poisson-Lindley can never adaptin case of under-dispersion.

2.3 Conway-Maxwell Poisson distribution

The Conway-Maxwell Poisson (CMP) is a generalization of the Poisson distribution introduced by Conway and Maxwell (1962) in a queue theory context. It depends on two parameters, $\lambda > 0$ and $v > 0$. The probability function is defined as

$$P(Y = y) = \frac{f(y; \theta)}{\sum_{j=0}^{\infty} f(j; \theta)} = \frac{\frac{\lambda^y}{(y!)^v}}{\sum_{j=0}^{\infty} \frac{\lambda^j}{(j!)^v}}$$

where $\theta = [\lambda, v]'$. In order to compute this probability function, the normalizing constant $\sum_{j=0}^{\infty} f(j; \theta)$ has to be computed. A practical approach is its truncation after the k -th term, obtaining

$$\sum_{j=0}^{\infty} f(j; \theta) = \sum_{j=0}^k \frac{\lambda^j}{(j!)^v} + R_k$$

where $R_k = \sum_{j=k+1}^{\infty} \frac{\lambda^j}{(j!)^v}$ is the absolute error due to the truncation.

The dispersion index of the CMP distribution can allow for over-dispersion (in this case $0 < v < 1$), equi-dispersion ($v = 1$, and the CMP collapses to a Poisson distribution with parameter λ) and under-dispersion ($v > 1$).

2.4 Negative Binomial Distribution

The Negative Binomial distribution depends on two parameters, N and P . The probability function is

$$P(Y = y) = \frac{\Gamma(N + y) \left(\frac{P}{P+1}\right)^y}{\Gamma(y+1)\Gamma(N) \left(\frac{P}{P+1}\right)^y} = \frac{\Gamma(N + y) \left(\frac{P}{P+1}\right)^y}{\Gamma(y+1)\Gamma(N) \left(\frac{P}{P+1}\right)^y} = \frac{\Gamma(N + y) \left(\frac{P}{P+1}\right)^y \left(\frac{1}{P+1}\right)^N}{\sum_{j=0}^{\infty} \Gamma(N + j) \left(\frac{P}{P+1}\right)^j \left(1 - \frac{P}{P+1}\right)^{-N}}$$

with $N > 0$ and $P > 0$.

Its ability to adapt to real data finds a limitation in the dispersion index

$$D = \frac{P + 1}{P}$$

which is greater than 1, only allowing for the case of over-dispersion.

3 A new distribution

A new discrete distribution is proposed starting from the series with generic term given by

$$f(j) = \frac{(1 + j)^a}{(1 + c)^j + \frac{1}{1+j}}$$

with $a \in \mathcal{R}$ and $c > 0$. It has positive values and is convergent. In fact, using the ratio criterion, we have

$$\lim_{j \rightarrow \infty} \frac{f(j+1)}{f(j)} = \frac{(2 + j)^a}{\left[(1 + c)^{(1+j)} + \frac{1}{2 + j}\right]} \cdot \frac{\left[(1 + c)^j + \frac{1}{1 + j}\right]}{(1 + j)^a} = \lim_{j \rightarrow \infty} \frac{(2 + j)^a}{(1 + j)^a} \cdot \frac{\left[(1 + c)^j + \frac{1}{1 + j}\right]}{\left[(1 + c)^{(1+j)} + \frac{1}{2 + j}\right]} = \frac{1}{1 + c} < 1.$$

Using this convergent series, a novel probability distribution can be generated, that is

$$P(Y = y) = \frac{f(y; \theta)}{\sum_{j=0}^{\infty} f(j; \theta)} = \frac{1}{C(a, c)} \cdot \frac{(1 + y)^a}{\left[(1 + c)^y + \frac{1}{1+y}\right]}$$

where $\theta = [a, c]'$ and

$$C(a, c) = \sum_{j=0}^{\infty} \frac{(1 + j)^a}{\left[(1 + c)^j + \frac{1}{1+j}\right]}$$

The value of $C(a, c)$ depends on the parameters a and c . In particular, when c increases, the function $C(a, c)$ decreases, while the relationship with the parameters a depends on the sign of a . When $a > 0$, fixed c , $C(a, c)$ increases when a increases. When a is negative, the inverse relationship holds.

Figures 1 and 2 show the shape of the probability function:

- varying a for a given value of c ;
- varying c for a given value of a .

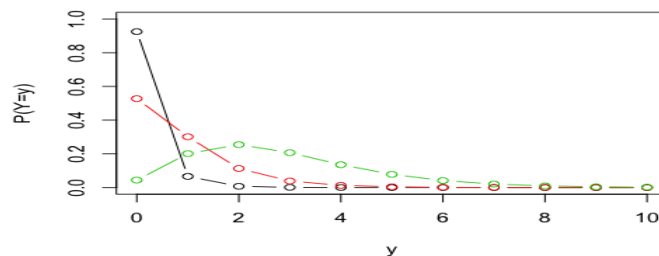


Figure 1 – Probability function with $c = 2, a = -3$ (black), $a = 0$ (red), $a = 3$ (green).

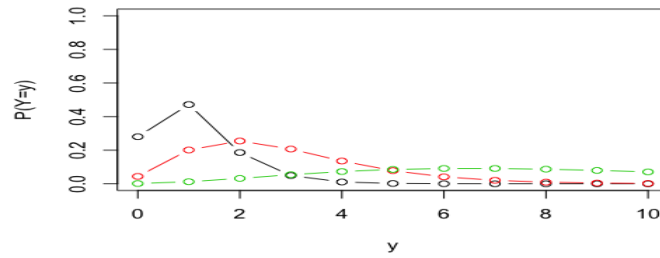


Figure 2 – Probability function with $a = 3, c = 8$ (black), $c = 2$ (red), $c = 0.5$ (green).

The computation of $C(a, c)$ compels to truncate the series at a high value J (in the application the sum will be computed up to $J = 50000$).

The expected value is given by

$$E[Y] = \sum_{y=0}^{\infty} \frac{y}{C(a, c)} \cdot \frac{(1+y)^a}{\left[\frac{1}{1+y} + (1+c)^y \right]} = \frac{1}{C(a, c)} \sum_{y=0}^{\infty} \frac{y \cdot (1+y)^a}{\left[\frac{1}{1+y} + (1+c)^y \right]}$$

It has no closed form, but can be expressed by a convergent series. In fact, keeping in mind the ratio criterion, the series with general term given by

$$\frac{y \cdot (1+y)^a}{C(a, c) \cdot \left[\frac{1}{1+y} + (1+c)^y \right]}$$

converges. Moreover, it is easy to show that

$$\lim_{y \rightarrow \infty} \frac{y+1}{y} \cdot \frac{(2+y)^a}{(1+y)^a} \cdot \frac{\left[(1+c)^y + \frac{1}{1+y} \right]}{\left[(1+c)^{(1+y)} + \frac{1}{2+y} \right]} = \frac{1}{1+c}.$$

Considering that the variance is

$$Var[Y] = \frac{1}{C(a, c)} \sum_{y=0}^{\infty} \frac{(y - E[Y])^2 \cdot (1+y)^a}{\left[\frac{1}{1+y} + (1+c)^y \right]},$$

the dispersion index is given by

$$D = \frac{\sum_{y=0}^{\infty} \frac{(y - E(Y))^2 \cdot (1+y)^a}{\left[\frac{1}{1+y} + (1+c)^y \right]}}{\sum_{y=0}^{\infty} \frac{y \cdot (1+y)^a}{\left[\frac{1}{1+y} + (1+c)^y \right]}}$$

The index D can assume values both greater than 1 and lower than 1. The proposed distribution is then a very flexible probabilistic model that can model both over-dispersion and under-dispersion. Figure 3 shows the flexible behavior of the dispersion index when parameter c varies, given some fixed values of a .

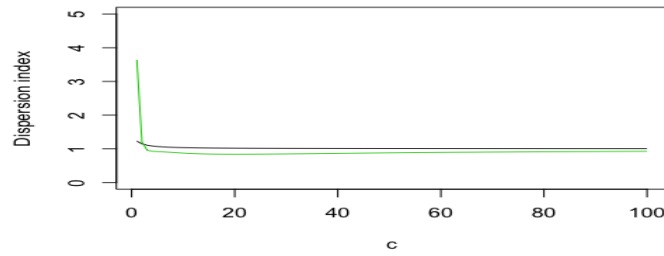


Figure 3 – Dispersion index with $a = -3$ (black), $a = 3$ (green).

In the estimation step, we have to take into account a parameter constraint, which concerns the parameter c , which has to be positive. The log-likelihood given by

$$\begin{aligned} \ln(L) &= \sum_{i=1}^N \left\{ a \ln(y_i + 1) - \ln[C(a, c)] - \ln \left[(1 + c)^{y_i} + \frac{1}{1 + y_i} \right] \right\} = \\ &= a \sum_{i=1}^N \ln(y_i + 1) - N \cdot \ln[C(a, c)] - \sum_{i=1}^N \ln \left[(1 + c)^{y_i} + \frac{1}{1 + y_i} \right], \end{aligned}$$

includes the quantity $C(a, c)$ which has to be considered after truncating the series. The system of equations for the log-likelihood solutions is

$$\begin{aligned} \frac{\partial \ln(L)}{\partial a} &= \sum_{i=1}^N \ln(y_i + 1) - \frac{N}{C(a, c)} \cdot \sum_{j=0}^{\infty} \frac{(1 + j)^a \ln(1 + j)}{\left[(1 + c)^j + \frac{1}{1 + j} \right]} \\ \frac{\partial \ln(L)}{\partial c} &= \frac{N}{C(a, c)} \cdot \sum_{j=0}^{\infty} \frac{(j + 1)^a j (1 + c)^{j-1}}{\left[(1 + c)^j + \frac{1}{1 + j} \right]} - \sum_{j=0}^{\infty} \frac{y_i (1 + c)^{y_i-1}}{\left[(1 + c)^{y_i} + \frac{1}{1 + y_i} \right]} \end{aligned}$$

and is solved using iterative methods.

4 Applications to real datasets

The various discrete distributions here considered have been evaluated in their ability to adequately fit counting variables. The analysis has been carried out using six automobiles insurance datasets relating to different countries / years and concerning the annual number of claims per policy. The datasets has already been used in Denuit (1997). It contains the data of four countries (Belgium for years 1958, 1975/1976 and 1994, Germany for year 1960, Switzerland for year 1961 and finally Zaire for year 1974). The main features of these datasets are summarized in Table 1. Two remarkable features are the over-dispersion and the presence of a high percentage of zeros, from 82% to over 92%, for all the datasets. The range, equal to the maximum value, is low for each dataset. Considering the enormous importance in the industrial countries of motor vehicles insurance for civil liability towards third parties, in the statistical and actuarial literature special attention has to be paid to the search for an appropriate probabilistic model to model the distribution of the number of road accidents in which a motorist has incurred in a given period of time or the annual number of claims per policy presented to the insurance company. The performance of the different distributions has been evaluated using statistical indices able to take into account the number of parameters to be estimated, the Akaike Information Criterion (AIC) and the Bayesian Information Criterion (BIC). The AIC proposed in Akaike (1974) is given by

$$AIC = -2 \ln(L) + 2r$$

where $\ln(L)$ denotes the log-likelihood and r is the number of parameters, while the BIC (Schwarz, 1978) is

$$BIC = -2 \ln(L) + r \ln(n)$$

where n is the number of observations. The best model is the model presenting the lowest value of AIC or BIC.

In Tables 2-7, θ_1 denotes the unique parameter in Poisson and Poisson-Lindely probability functions. For CMP, Negative Binomial and novel distribution θ_1 denotes parameter λ , N and a , respectively, while θ_2 denotes the second parameter, that is v , P and c .

Table 2 (Belgian data, 1958) shows that the Poisson distribution has a very bad performance and is very distant from the remaining distributions in terms of AIC and BIC. According to the two indices, the novel distribution provides the best fit to the data and its principal competitor is the Negative Binomial. Table 3 (German data, 1960) indicates that according both to the AIC and to the BIC the winner probabilistic model turns out to be the novel distribution. Table 4 (Swiss data, 1961) confirms what we have found for the two datasets, that is the prominent role of the novel distribution which again ensures the most satisfactory fit of the data in terms of AIC and BIC. In Table 5 (Zairian data, 1974) the distribution here proposed provides better results for both the criteria. Tables 6 and 7 (Belgian data, 1975/1976 and 1994) provide evidences that do not deviate from the main previous results, that is the best fit is reached using the novel distribution, followed by the Negative Binomial distribution. Finally, Tables 8 and 9 summarize the AIC and BIC values.

	Belgium (1958)	Germany (1960)	Switzerland (1961)	Zaire (1974)	Belgium (1975/76)	Belgium (1994)
Observations	9461	23589	119853	4000	106974	131182
Average	0,2143	0,1442	0,1551	0,0865	0,1011	0,1036
Variance	0,2889	0,1639	0,1793	0,1225	0,1074	0,1115
Dispersion index	1,3478	1,1362	1,1558	1,4164	1,0630	1,0764
Percentage of 0's	0,8287	0,8729	0,8653	0,9297	0,9066	0,9048
Range	7	6	6	5	4	4

Table 1 – Descriptive statistics for the six datasets

	$\hat{\theta}_1$	s. e. ($\hat{\theta}_1$)	$\hat{\theta}_2$	s. e. ($\hat{\theta}_2$)	AIC	BIC
Poisson	0,2144	0,0048	-	-	10983,56	10990,71
Poisson-Lindley	5,3998	0,1182	-	-	10714,56	10721,71
CMP	0,1765	0,0049	0,0000	0,0715	10713,36	10727,67
NB	0,7014	0,0628	0,3056	0,0284	10696,08	10714,39
New distribution	-1,8856	0,1377	1,7722	0,2388	10693,83	10708,14

Table 2 – Estimation results for Belgium data (1958)

	$\hat{\theta}_1$	s. e. ($\hat{\theta}_1$)	$\hat{\theta}_2$	s. e. ($\hat{\theta}_2$)	AIC	BIC
Poisson	0,1442	0,0025	-	-	20597,69	20605,76
Poisson-Lindley	7,7279	0,1297	-	-	20449,76	20457,83
CMP	0,1274	0,0026	0,0599	0,0684	20451,22	20467,36
NB	1,1175	0,1194	0,1291	0,0140	20450,84	20466,98
New distribution	4,2311	0,0495	0,7316	0,0271	20447,28	20453,35

Table 3 – Estimation results for German data (1960)

	$\hat{\theta}_1$	s. e. ($\hat{\theta}_1$)	$\hat{\theta}_2$	s. e. ($\hat{\theta}_2$)	AIC	BIC
Poisson	0,1551	0,0011	-	-	110218,90	110228,59
Poisson-Lindley	7,2292	0,0519	-	-	109231,40	109243,09
CMP	0,1346	0,0012	0,0119	0,0276	109235,00	109254,39
NB	1,0328	0,0436	0,1502	0,0065	109234,60	109253,99
New distribution	4,0613	0,0192	0,7013	0,0112	109223,18	109242,57

Table 4 – Estimation results for Swiss data (1961)

	$\hat{\theta}_1$	s. e. ($\hat{\theta}_1$)	$\hat{\theta}_2$	s. e. ($\hat{\theta}_2$)	AIC	BIC
Poisson	0,0865	0,0046	-	-	2494,15	2500,45
Poisson-Lindley	12,4367	0,6547	-	-	2417,30	2423,60
CMP	0,0796	0,0047	0,0000	0,2589	2418,85	2431,44
NB	0,2166	0,0364	0,3993	0,0717	2371,10	2383,69
New distribution	4,5349	0,0971	0,2277	0,0981	2372,23	2384,82

Table 5 – Estimation results for Zairian data (1974)

	$\hat{\theta}_1$	s.e. ($\hat{\theta}_1$)	$\hat{\theta}_2$	s.e. ($\hat{\theta}_2$)	AIC	BIC
Poisson	0,1011	0,0010	-	-	72378,51	72388,09
Poisson-Lindley	10,7345	0,1010	-	-	72247,06	72256,64
CMP	0,0951	0,0010	0,2982	0,0500	72213,01	72232,17
NB	1,6295	0,1474	0,0620	0,0056	72212,20	72231,36
New distribution	5,1799	0,0404	0,8682	0,0164	72211,52	72230,68

Table 6 – Estimation results for Belgian data (1975/1976)

	$\hat{\theta}_1$	s.e. ($\hat{\theta}_1$)	$\hat{\theta}_2$	s.e. ($\hat{\theta}_2$)	AIC	BIC
Poisson	0,1036	0,0009	-	-	90455,12	90474,69
Poisson-Lindley	10,4888	0,0880	-	-	90184,50	90194,28
CMP	0,0963	0,0009	0,2030	0,0425	90164,34	90183,91
NB	1,3812	0,0992	0,0750	0,0054	90163,16	90182,73
New distribution	5,0560	0,0387	0,8325	0,0160	90160,54	90180,11

Table 7 – Estimation results for Belgia data (1994)

	Belgium (1958)	Germany (1960)	Switzerland (1961)	Zaire (1974)	Belgium (1975/1976)	Belgium (1994)
Poisson	10983,56	20597,69	110218,90	2494,15	72378,51	90455,12
Poisson-Lindley	10714,56	20449,76	109233,40	2417,30	72247,06	90184,50
CMP	10713,46	20451,22	109235,00	2418,85	72213,01	90164,34
NB	10700,08	20450,84	109234,60	2371,10	72212,20	90163,16
New distribution	10693,83	20447,76	109223,96	2370,95	72211,68	90160,44

Table 8 – Summary of AIC criterion

	Belgium (1958)	Germany (1960)	Switzerland (1961)	Zaire (1974)	Belgium (1975/1976)	Belgium (1994)
Poisson	10990,71	20605,76	110228,59	2500,45	72388,09	90464,90
Poisson-Lindley	10721,71	20457,83	109243,09	2423,60	72256,64	90194,28
CMP	10727,77	20467,36	109254,39	2431,44	72232,17	90183,91
NB	10714,39	20466,98	109253,99	2383,69	72231,36	90182,73
New distribution	10708,14	20453,35	109242,57	2383,54	72230,84	90180,01

Table 9 – Summary of BIC criterion

5 Conclusions

In this article a new discrete distribution is proposed. Its main feature is the flexibility to adapt to many datasets using two parameters. It allows for a dispersion index taking into account both underdispersion and overdispersion. The distribution is compared to traditional (Poisson and Negative Binomial) and flexible (Poisson-Lindley and CMP) discrete distribution, using popular insurance datasets already analyzed in literature. The best fit in terms of AIC and BIC is always achieved with our proposed distribution which appears to be very promising for any discrete dataset.

References

Akaike H. (1974), A new look at the statistical model identification, *IEEE Transactions on Automatic Control*, 19, 716–723

Conway R.W. and Maxwell W.L. (1962), A queuing model with state dependent service rates, *Journal of Industrial Engineering*, 12, 132-136

DenuitM. (1997), A new distribution of Poisson-type for the number of claims. *ASTIN Bulletin*, 27, 229–242.

Sankaran M. (1970), The discrete Poisson-Lindley distribution, *Biometrics*, 26, 145-149

Schwarz G.E. (1978), Estimating the dimension of a model, *Annals of Statistics*, 6, 461-464.