

Modeling COVID-19 Pandemic Data with New Pareto Model

ABSTRACT

This paper aims to find a statistical model for modeling the COVID-19 data. We deduced a robust and effective model for fitting the COVID 19 mortality. This model is a new Extended-Pareto distribution (NE-P). The maximum likelihood method is utilized to obtain the estimator of the parameters. A simulation was carried out using different sample sizes and different values of the parameters. In addition, the goodness of fit test statistics was calculated for proposed model compared with the baseline model to find out that our new model is the best for modeling data COVID-19.

Keywords: A new Extended-Pareto distribution, COVID 19 mortality, The maximum likelihood method, Goodness of fit.

1 Introduction

Many important distributions are widely used in various statistical applications to model several life-time data in applied fields such as engineering, insurance, economics, medicine, and life testing, among others. The limitation of the standard distributions arouses the interest of finding new distributions by extending existing ones. One of the most important of these distributions is Pareto distribution (PD) that is well known in the literature for its capability in modelling the heavy-tailed distributions that are mostly common in data on income distribution, economics, survival analysis. So, many authors have been interested in proposing new generalized forms of Pareto distribution which expanded the applications of this distribution to include medicine and health. For example, Aldahlan et al. [1] introduced the Marshall–Olkin Pareto

27 Type-I (MOPTI) distribution. They studied the statistical properties of the
 28 MOPTI distribution. Also, they presented a simulation study and application on
 29 a real data set on breast cancer. Caeiro and Mateus [6] developed a new class
 30 of estimators for the parameters of Pareto type I distribution named the log-
 31 generalized probability-weighted moment (LGPWM). They found the proposed
 32 LGPWM estimators were capable of competing with the most commonly used
 33 estimation methods. Jayakumar et al. [7] presented a new four parameter
 34 distribution called New Generalized Pareto distribution, which was a
 35 generalization of the classical Pareto distribution. Boumaraf et al. [5] used
 36 nonlinear optimization methods to find the estimators of beta Pareto
 37 distribution.

38 The motivation behind this research is to expand the application areas of
 39 Pareto distribution to include medical data modeling. By introducing a new
 40 model as a generalization for Pareto distribution and demonstrate its flexibility
 41 in modeling COVID-19 data compared to the original Pareto distribution.

42 The probability density function (PDF) of Pareto distribution is given by:

$$43 \quad g(x) = \frac{\beta}{x^{\beta+1}} ; x \geq 1 ; \beta > 0 . \quad (1)$$

44 and the cumulative distribution function (CDF) is:

$$45 \quad G(x) = 1 - x^{-\beta} ; x \geq 1 , \beta > 0 . \quad (2)$$

46 where β is the scale parameter.
 47
 48
 49

50 Moreover, many researchers have focused on finding generators for new
 51 distributions by finding new families, for example: Sule et al. in [8] and Bantan
 52 et al. in [4]. In [10], the authors studied a new extended (NE-X) family of
 53 distributions which is the generator of our new model. The PDF of this NE-X
 54 distribution is given by:

$$55 \quad f(x) = \frac{2\theta^2 g(x)G(x)(1-G^2(x))^{\theta-1}}{(1-(1-\theta)G^2(x))^{\theta+1}} ; \theta > 0 ; x \in R . \quad (3)$$

56

57 The CDF of NE-X distribution is:

58

$$59 \quad F(x) = 1 - \left(\frac{1-G^2(x)}{1-(1-\theta)G^2(x)} \right)^\theta ; \theta > 0; x \in R. \quad (4)$$

60

61 Depending on Equations (1) , (2) and the family in [10], we deduced a new

62 distribution called it a New Extended-Pareto Distribution (NE-P). The PDF

63 and CDF of NE-P distribution with two parameters(θ, β) is obtained

64 respectively as:

$$65 \quad f(x) = \frac{2\theta^2(\beta x^{-(\beta+1)})(1-x^{-\beta})(1-(1-x^{-\beta})^2)^{\theta-1}}{(1-(1-\theta)(1-x^{-\beta})^2)^{\theta+1}} ; \theta > 0, \beta > 0; x \in R. \quad (5)$$

66

67

68 And

$$69 \quad F(x) = 1 - \left(\frac{1-(1-x^{-\beta})^2}{1-(1-\theta)(1-x^{-\beta})^2} \right)^\theta ; \theta > 0, \beta > 0; x \in R \quad (6)$$

70

71 We can rewrite the PDF & CDF of NE-P distribution, using the series

72 representation as follows

$$73 \quad f(x) = 2\theta^2\beta \sum_{i=0}^{\infty} \sum_{j=0}^{\theta-1} \sum_{k=0}^{2(i+j)+1} \binom{i+\theta}{\theta} \binom{\theta-1}{j} \binom{2(i+j)+1}{k} \\ 74 \quad \times (1-\theta)^i (-1)^{j+k} (x^{-\beta(k+1)-1}), \quad (7)$$

75 and

$$76 \quad F(x) = 1 - \left[\sum_{i=0}^{\infty} \sum_{j=0}^{\theta} \sum_{k=0}^{2(i+j)} \binom{i+\theta-1}{\theta-1} \binom{\theta}{j} \binom{2(i+j)}{k} \right. \\ 77 \quad \left. \times (1-\theta)^i (-1)^{j+k} x^{-\beta k} \right]. \quad (8)$$

78

79 In the article, we estimate the NE-P distribution parameters using the

80 maximum likelihood estimation and carry out by different complete samples

81 size of NE-P distribution. In addition, the goodness of fit test statistics calculate

82 for proposed models to find out the best of it for data of Coronavirus disease

83 (COVID-19) .

84

85 **2 Maximum Likelihood Estimation Method**

86 This section presents the maximum likelihood estimator of the NE-P
 87 distribution parameters (θ, β) . If x_1, x_2, \dots, x_n is a random sample from NE-
 88 P distribution, the Likelihood function is $L(\underline{x})$ can be obtained as:

$$89 \quad L(\underline{x}) = (2\theta^2\beta)^n \prod_{i=1}^n x_i^{-(\beta+1)} \prod_{i=1}^n (1 - x_i^{-\beta}) \prod_{i=1}^n \left(1 - (1 - x_i^{-\beta})^2\right)^{\theta-1} \\
 90 \quad \times \prod_{i=1}^n \left(1 - (1 - \theta)(1 - x_i^{-\beta})^2\right)^{-(\theta+1)}. \quad (9)$$

92 And the log-likelihood function is given as follows

$$93 \\
 94 \quad l = \log(L) = n \log(2) + n \log(\theta^2) + n \log(\beta) - (\beta + 1) \sum_{i=1}^n \log x_i \\
 95 \quad + \sum_{i=1}^n \log(1 - x_i^{-\beta}) + (\theta - 1) \sum_{i=1}^n \log(1 - (1 - x_i^{-\beta})^2) \\
 96 \quad - (\theta + 1) \sum_{i=1}^n \log(1 - (1 - \theta)(1 - x_i^{-\beta})^2). \quad (10)$$

97 Differentiating (10) with respect to each of the parameters θ and β gives

$$98 \\
 99 \quad \frac{\partial \log L}{\partial \theta} = \frac{2n}{\theta} + \sum_{i=1}^n \log(1 - (1 - x_i^{-\beta})^2) - (\theta + 1) \\
 100 \quad \times \sum_{i=1}^n \frac{(1 - x_i^{-\beta})^2}{(1 - (1 - \theta)(1 - x_i^{-\beta})^2)} - \sum_{i=1}^n \log(1 - (1 - \theta)(1 - x_i^{-\beta})^2), \quad (11)$$

$$101 \quad \frac{\partial \log L}{\partial \beta} = \frac{n}{\beta} - \sum_{i=1}^n \log x_i + \sum_{i=1}^n \frac{x_i^{-\beta} \ln x_i}{1 - x_i^{-\beta}} - 2(\theta - 1) \\
 102 \quad \times \sum_{i=0}^n \frac{x_i^{-\beta} \ln(x_i)(1 - x_i^{-\beta})}{1 - (1 - x_i^{-\beta})^2} + 2(1 - \theta^2) \sum_{i=0}^n \frac{x_i^{-\beta} \ln(x_i)(1 - x_i^{-\beta})}{(1 - (1 - \theta)(1 - x_i^{-\beta})^2)}. \quad (12)$$

103
 104 There isn't a closed form to solve this equations for (θ, β) . As a result, the
 105 equations can be solved numerically using the Newton-Raphson method and
 106 Mathematica program V. 11.0 to determine the Maximum Likelihood Estimate
 107 $\hat{\theta}_{MLE}$ and $\hat{\beta}_{MLE}$.

108
 109

110 3 Simulation Study

111 In this section, the simulation result for the ML method is given when two
 112 parameters are unknown based on complete samples for various sample sizes
 113 and proposed initial values for parameters. The parameter values are selected
 114 as $\beta = 2.8, \theta=1.5$ and $n = 25, 100, 250, 450,$ and 1000 . This process is
 115 repeated $N = 500$ times . Furthermore, performance of different estimators is
 116 considered in terms of their biases and mean square errors (MSEs) that given,
 117 respectively, by

118

119 $Bias(\hat{\lambda}) = E(\hat{\lambda}) - \lambda$ and $MSE(\hat{\lambda}) = E(\hat{\lambda} - \lambda)^2$.

120 Where λ any parameter.

121
 122
 123
 124

Table 1: Mean, MSEs and Bias for the parameter estimates when $\beta_0 = 2.8, \theta_0=1.5$.

N		MLE	MSE	BIAS
25	$\hat{\beta}$	3.37488	0.76493	0.574879
	$\hat{\theta}$	1.67533	0.627229	0.17533
100	$\hat{\beta}$	3.00401	0.638085	0.204005
	$\hat{\theta}$	1.67228	0.622801	0.17228
250	$\hat{\beta}$	3.00453	0.27286	0.204525
	$\hat{\theta}$	1.579	0.185657	0.0790023
450	$\hat{\beta}$	2.94519	0.339871	0.145187
	$\hat{\theta}$	1.55915	0.151468	0.059146
1000	$\hat{\beta}$	2.89381	0.069541	0.0938141
	$\hat{\theta}$	1.53192	0.05152	0.0319226

125

126 From Table (1) MSE and Bias are displayed. It can be illustrated clearly that
 127 these estimates are reasonably consistent and approaches to the true values
 128 of parameters as sample size increases. Furthermore, with increasing sample
 129 size the MSEs and Bias decrease for all parameter combinations. Therefore,

130 it has been concluded that MLE process performs well in estimating the
131 parameters of NE-P distribution.

132

133 **4 Real Data Applications**

134 In this section, we provide the application with real data sets to assess the
135 flexibility of NE-P distribution comparing with the base line Pareto distribution.

136 The parameters are estimated using maximum likelihood method.

137 Mathematica (V.11.0) is used for computation. Moreover, we consider the

138 model selection criteria, including Akaike information criterion (AIC), Bayesian

139 information criterion (BIC), consistent Akaike information criterion (CAIC), and

140 Hannan-Quinn information criterion (HQIC). They are defined as follows:

$$141 \quad AIC = -2l(\hat{\theta}) + 2k$$

$$142 \quad CAIC = AIC + \frac{2k(k+1)}{n-k-1}$$

$$143 \quad HQIC = -2\log l(\hat{\theta}) + 2k \log(\log(n))$$

144 (See Whittaker and Furlow [9]).

145

146 **5 Daily mortality cases of COVID-19:**

147 In this section, we will study the data number of daily mortality COVID-19

148 cases will be compared with (PD).

149

150 **5.1 Data set 1:**

151

152 The first data represents a COVID-19 mortality rates data belongs to Italy of

153 59 days, that is recorded from 27 February to 27 April 2020. The data is taken

154 from (Almongy et al. [3]) as follows:

155 4.571, 7.201, 3.606, 8.479, 11.410 ,8.961, 10.919, 10.908, 6.503 ,18.474,

156 11.010 ,17.337, 16.561, 13.226, 15.137 ,8.697 ,15.787 ,13.333 ,11.822

157 ,14.242 ,11.273, 14.330, 16.046, 11.950, 10.282, 11.775 ,10.138 ,9.037

158 ,12.396 ,10.644, 8.646 ,8.905, 8.906, 7.407, 7.445, 7.214, 6.194, 4.640 ,5.452

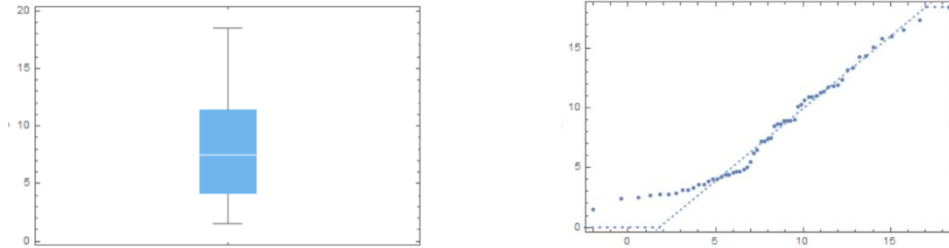
159 ,5.073, 4.416, 4.859 ,4.408 ,4.639 ,3.148 ,4.040 ,4.253 ,4.011 ,3.564, 3.827,
 160 3.134, 2.780, 2.881, 3.341, 2.686, 2.814, 2.508, 2.450 ,1.518.

161
 162

Table 2: Descriptive statistics for data set 1.

n	Min	Q_1	Median	Mean	Q_3	Max	Skewness	Kurtosis
59	1.518	4.04	7.44	8.15	11.41	18.474	0.45	2.12

163
 164
 165
 166
 167
 168



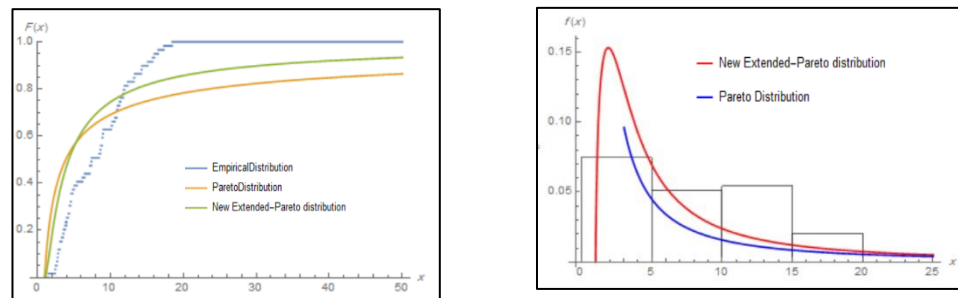
169 **Figure 1.** PP plot of the NE-P distribution and the box plot for data set 1.

170
 171

Table 3. Parameter estimation for various distributions depending on data set 1.

Model	Parameters		LL	AIC	CAIC	HQIC
	$\hat{\theta}$	$\hat{\beta}$				
NE-P	0.31	2.31	-190.47	384.94	385.15	386.56
PD		0.52	-211.25	424.50	424.57	425.31

172
 173
 174
 175
 176
 177
 178
 179
 180
 181



182
 183 **Figure 2.** Plots of the fitted CDF (left) and the histogram with fitted PDF (right) of the NE-P model
 184 for data set 1.

185
 186
 187
 188
 189

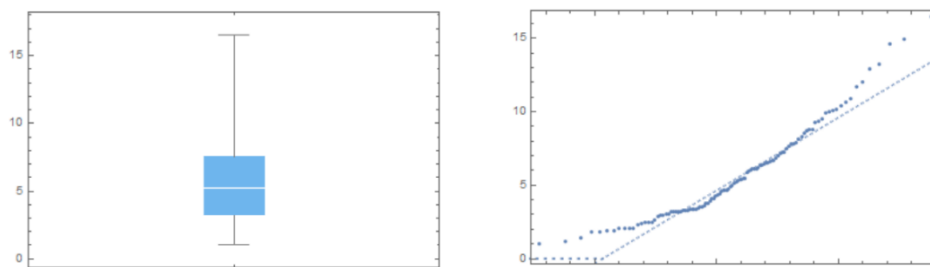
5.2 Data set 2:

190
 191 The second data set of data of COVID-19 mortality numbers in Mexico of 108
 192 days, that is recorded from 4 March to 20 July 2020. The data is taken from
 193 Nagy et al. The data is taken from (Almongy et al. [3]) as follows:
 194 8.826, 6.105, 10.383, 7.267, 13.220, 6.015, 10.855, 6.122, 10.685, 10.035,
 195 5.242, 7.630, 14.604, 7.903, 6.327, 9.391, 14.962, 4.730, 3.215, 16.498,
 196 11.665, 9.284, 12.878, 6.656, 3.440, 5.854, 8.813, 10.043, 7.260, 5.985,
 197 4.424, 4.344, 5.143, 9.935, 7.840, 9.550, 6.968, 6.370, 3.537, 3.286, 10.158,
 198 8.108, 6.697, 7.151, 6.560, 2.988, 3.336, 6.814, 8.325, 7.854, 8.551, 3.228,
 199 3.499, 3.751, 7.486, 6.625, 6.140, 4.909, 4.661, 1.867, 2.838, 5.392, 12.042,
 200 8.696, 6.412, 3.395, 1.815, 3.327, 5.406, 6.182, 4.949, 4.089, 3.359, 2.070,
 201 3.298, 5.317, 5.442, 4.557, 4.292, 2.500, 6.535, 4.648, 4.697, 5.459, 4.120,
 202 3.922, 3.219, 1.402, 2.438, 3.257, 3.632, 3.233, 3.027, 2.352, 1.205, 2.077,
 203 3.778, 3.218, 2.926, 2.601, 2.065, 1.041, 1.800, 3.029, 2.058, 2.326, 2.506,
 204 1.923.

205 **Table 4:** Descriptive statistics for data set 2.

n	Min	Q_1	Median	Mean	Q_3	Max	Skewness	Kurtosis
108	1.041	3.23	5.19	5.75	7.48	16.498	0.98	3.68

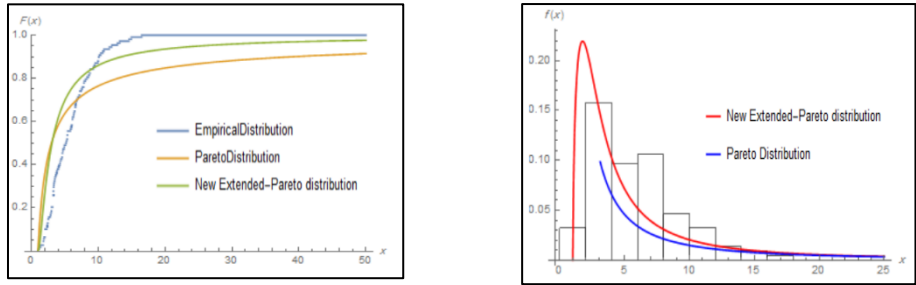
206



214 **Figure 3.** PP plot of the NE-P distribution and the box plot for data set 2.

215 **Table 5.** Parameter estimation for various distributions depending on data set 2.

Model	Parameters		LL	AIC	CAIC	HQIC
	$\hat{\theta}$	$\hat{\beta}$				
NE-P	0.33	2.65	-296.84	597.69	597.81	599.87
PD		0.63	-329.82	661.64	661.68	662.73



222

223

224

225

226

227

228

229

230

231

232

233

234

235

Figure 4. Plots of the fitted CDF (left) and the histogram with fitted PDF (right) of the NE-P model for data set 2.

5.3 Data set 3:

The third data set of a COVID-19 data belonging to the Netherlands of 30 days, that is recorded from 31 March to 30 April 2020. This data formed of rough mortality rate. (see Almongy et al. [3]) The data are as follows:

14.918, 10.656, 12.274, 10.289, 10.832, 7.099, 5.928, 13.211, 7.968, 7.584, 5.555, 6.027, 4.097, 3.611, 4.960, 7.498, 6.940, 5.307, 5.048, 2.857, 2.254, 5.431, 4.462, 3.883, 3.461, 3.647, 1.974, 1.273, 1.416, 4.235.

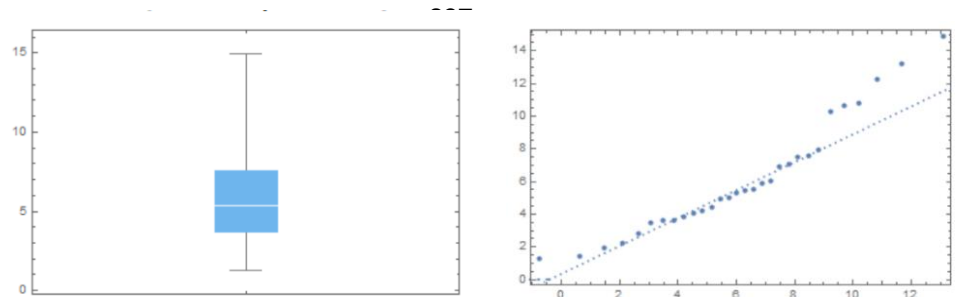
234

235

Table 6: Descriptive statistics for data set 3.

n	Min	Q_1	Median	Mean	Q_3	Max	Skewness	Kurtosis
30	1.273	3.64	5.36	6.15	7.58	14.918	0.83	2.95

236



245

246

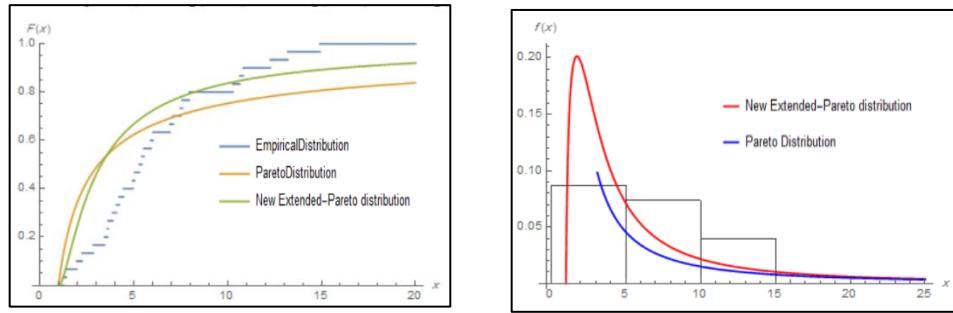
247

Figure 5. PP plot of the NE-P distribution and the box plot for data set 3.

Table 7. Parameter estimation for various distributions depending on data set 3.

Model	Parameters		LL	AIC	CAIC	HQIC
	$\hat{\theta}$	$\hat{\beta}$				
NE-P	0.37	2.32	-85.46	174.93	175.38	175.83
PD		0.60	-94.38	190.76	190.91	191.21

248
249
250
251
252
253
254
255
256



257 **Figure 6.** Plots of the fitted CDF (left) and the histogram with fitted PDF (right) of the NE-P model for
258 data set 3.
259

260
261

5.4 Data set 4:

262
263
264

The fourth data set represents a COVID-19 data belong to Canada of 36 days, from 10 April to 15 May 2020 (see Almetwally et al. [2]). These data formed of mortality rate. The data are as follows:

265
266
267
268

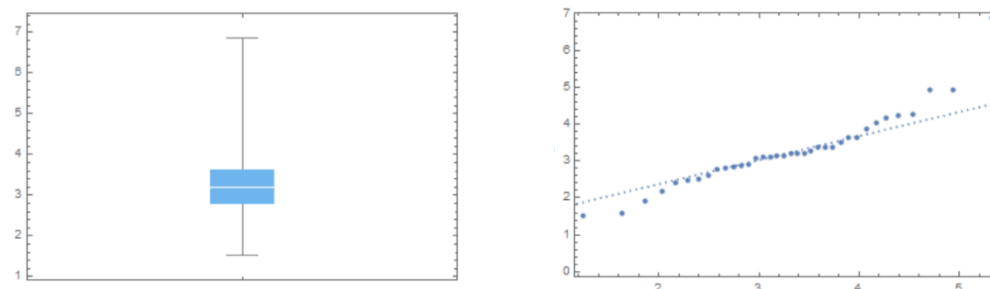
3.1091, 3.3825, 3.1444, 3.2135, 2.4946, 3.5146, 4.9274, 3.3769, 6.8686,
3.0914, 4.9378, 3.1091, 3.2823, 3.8594, 4.0480, 4.1685, 3.6426, 3.2110,
2.8636, 3.2218, 2.9078, 3.6346, 2.7957, 4.2781, 4.2202, 1.5157, 2.6029,
3.3592, 2.8349 ,3.1348, 2.5261, 1.5806, 2.7704, 2.1901, 2.4141, 1.9048.

269
270

Table 8: Descriptive statistics for data set 4.

n	Min	Q_1	Median	Mean	Q_3	Max	Skewness	Kurtosis
36	1.5171	2.77	3.17	3.28	3.63	6.8686	1.21	6.15

271
272



273
274
275
276
277
278
279
280

Figure 7. PP plot of the NE-P distribution and the box plot for data set 4.

281

Table 9. Parameter estimation for various distributions depending on data set 4.

Model	Parameters		LL	AIC	CAIC	HQIC
	$\hat{\theta}$	$\hat{\beta}$				
NE-P	0.32	3.75	-68.20	140.41	140.77	141.51
PD		0.87	-82.13	166.27	166.39	166.82

288

289

290

291

292

293

294

295

296

297

298

299

300

301

302

303

304

305

306

307

308

309

310

311

312

313

314

315

316

317

318

319

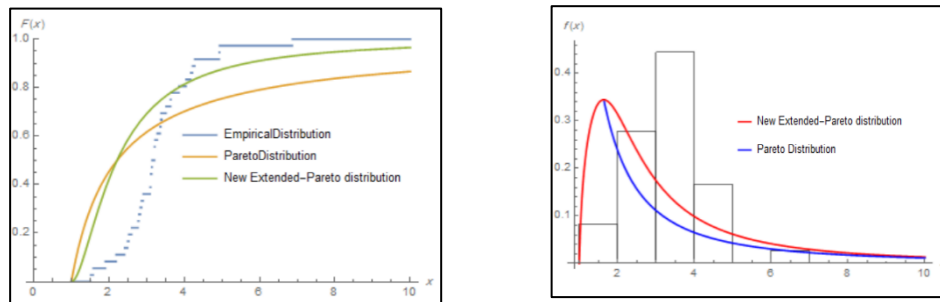


Figure 8. Plots of the fitted CDF (left) and the histogram with fitted PDF (right) of the NE-P model for data set 4.

From Tables 3, 5, 7 and 9, the values of log-likelihood (LL), AIC, CAIC and HQIC are minimum and favorable of NE-P distribution compared with PD distribution, which indicate that our new model is the best comparing with the competing model.

6 Conclusion

COVID-19 data modeling has gained renewed interest among researchers, particularly in the quest to find new models that are more flexible in modeling this data. In this article We found a New Extended -Pareto distribution as a model for these data. Its parameters were estimated by method of maximum likelihood. Performances of MLE were tested through simulation study. Finally, four real data applications of COVID-19 were analyzed in to assess the flexibility of our new model. We encourage researchers to continue exploring new models for modeling this kind of data sets. Future studies can expand the study of NE-P and apply other types of data as well as compare it to other competing distributions.

Disclaimer (Artificial intelligence)

320 Option 1:

321 Author(s) hereby declare that NO generative AI technologies such as Large Language
322 Models (ChatGPT, COPILOT, etc.) and text-to-image generators have been used during the
323 writing or editing of this manuscript.

324 Option 2:

325 Author(s) hereby declare that generative AI technologies such as Large Language Models,
326 etc. have been used during the writing or editing of manuscripts. This explanation will
327 include the name, version, model, and source of the generative AI technology and as well as
328 all input prompts provided to the generative AI technology

329 Details of the AI usage are given below:

- 330 1.
- 331 2.
- 332 3.

333

334 7 References:

335

- 336 1. Aldahlan, M. A., Rabie, A. M., Abdelhamid, M., Ahmed, A. H. N., &
337 Afify, A. Z. (2023). The Marshall–Olkin Pareto Type-I Distribution:
338 Properties, Inference under Complete and Censored Samples with
339 Application to Breast Cancer Data. *Pakistan Journal of Statistics and*
340 *Operation Research*, 603-622.
- 341 2. Almetwally, E. M., Alharbi, R., Alnagar, D., & Hafez, E. H. (2021). A
342 new inverted topp-leone distribution: applications to the COVID-19
343 mortality rate in two different countries. *Axioms*, 10(1), 25.
- 344 3. Almongy, H. M., Almetwally, E. M., Aljohani, H. M., Alghamdi, A. S., &
345 Hafez, E. H. (2021). A new extended Rayleigh distribution with
346 applications of COVID-19 data. *Results in Physics*, 23, 104012.
- 347 4. Bantan, R. A., Jamal, F., Chesneau, C., & Elgarhy, M. (2020). Type II
348 Power Topp-Leone generated family of distributions with statistical
349 inference and applications. *Symmetry*, 12(1), 75.
- 350 5. Boumaraf, B., Seddik-Ameur, N. and Barbu, V. S. (2020). Estimation of
351 Beta-Pareto distribution based on several optimization
352 methods. *Mathematics*, 8(7), 1-12.

353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369

6. Caeiro, F., & Mateus, A. (2023). A new class of generalized probability-weighted moment estimators for the Pareto distribution. *Mathematics*, 11(5), 1076.
7. Jayakumar, K., Krishnan, B. and Hamedani, G. G. (2018). On new generalization of Pareto distribution and its applications. *Communications in Statistics-Simulation and Computation*, 49(5), 1264-1284.
8. Sule, I., Sani, I. D., Audu, I., & Jibril, H. M. (2020). On the Topp Leone exponentiated-G family of distributions: properties and applications.
9. Whittaker, T. A., & Furlow, C. F. (2009). The comparison of model selection criteria when selecting among competing hierarchical linear models. *Journal of Modern Applied Statistical Methods*, 8(1), 15.
10. Zichuan, M., Hussain, S., Iftikhar, A., Ilyas, M., Ahmad, Z., Khan, D. M. and Manzoor, S. (2020). A new extended-family of distributions: properties and applications. *Computational and mathematical methods in medicine*, 2020, 1-13.