# Collated Distributions: A New Model for Infectious Disease Spread Using Early COVID-19 Data

Benjamin T. Solomon[1]

[1]Xodus One Management 9808 Amberton Pkwy Dallas, TX 75243

**Abstract: A new modelling technique, Collated Distributions (CD), is presented, to model Infectious Disease Spread (IDS). This modeling shows that there are 3 probability distributions that one must determine to effectively manage public health, infectability, mortality and survival. Thus, leading to better understanding of Infectious Disease Spread (IDS). The Reproduction Model used in COVID-19 disease spread modeling is shown to be doubtful.**

**Collated Distributions have implications in many other fields of study such as dissemination of knowledge, and the rise and fall of civilizations.**

*Keywords: Covid, Mortality, Survival, Disease, Health Policy.*

## I. INTRODUCTION

The purpose of this paper is to provide (i) A means to formulating an informed opinion about future disease spread. And (ii) to understand at which stage in the disease lifecycle is a treatment effective. This is accomplished by characterizing an infectious disease using Collated Distributions [1 & 2] so that rapid development of health & treatment policies to counter the disease spread can be developed early as COVID-19 has shown us.

This paper explains why the authoritative models (reproduction models [3 & 4] such as that of the Washington University's Institute of Health Metrics & Evaluation, IHME, Reproduction model) used to manage health policy are doubtful. Though the mathematics of these models are accepted on a peer-reviewed consensus basis to be valid, but due to implicit assumptions, the statistical implementation is biased, and cannot be corrected. See Discussion section.

The term, context structure, is used as in many cases many models are required to construct a context structure that fits the characteristics of the context structure. Models are best fits for the data and are only as good as the data available. The data is only as good as the test accuracy (whether the tests are accurate or provide substantial false positives or false negatives). It is interaction between the models that gives the context structure its power to provide deep insights into disease spread, i.e. context structure tells much more than the data alone could.

Note, that (i) the Infectious Disease Spread (IDS) model presented here does not belong to any of the three know [5] model types, metapopulation model, cellular automaton model or gravity model, and (ii) The US COVID-19 data [6] used was from February 23rd, 2020 to April 5th, 2020 (early in the pandemic). See Table 1.
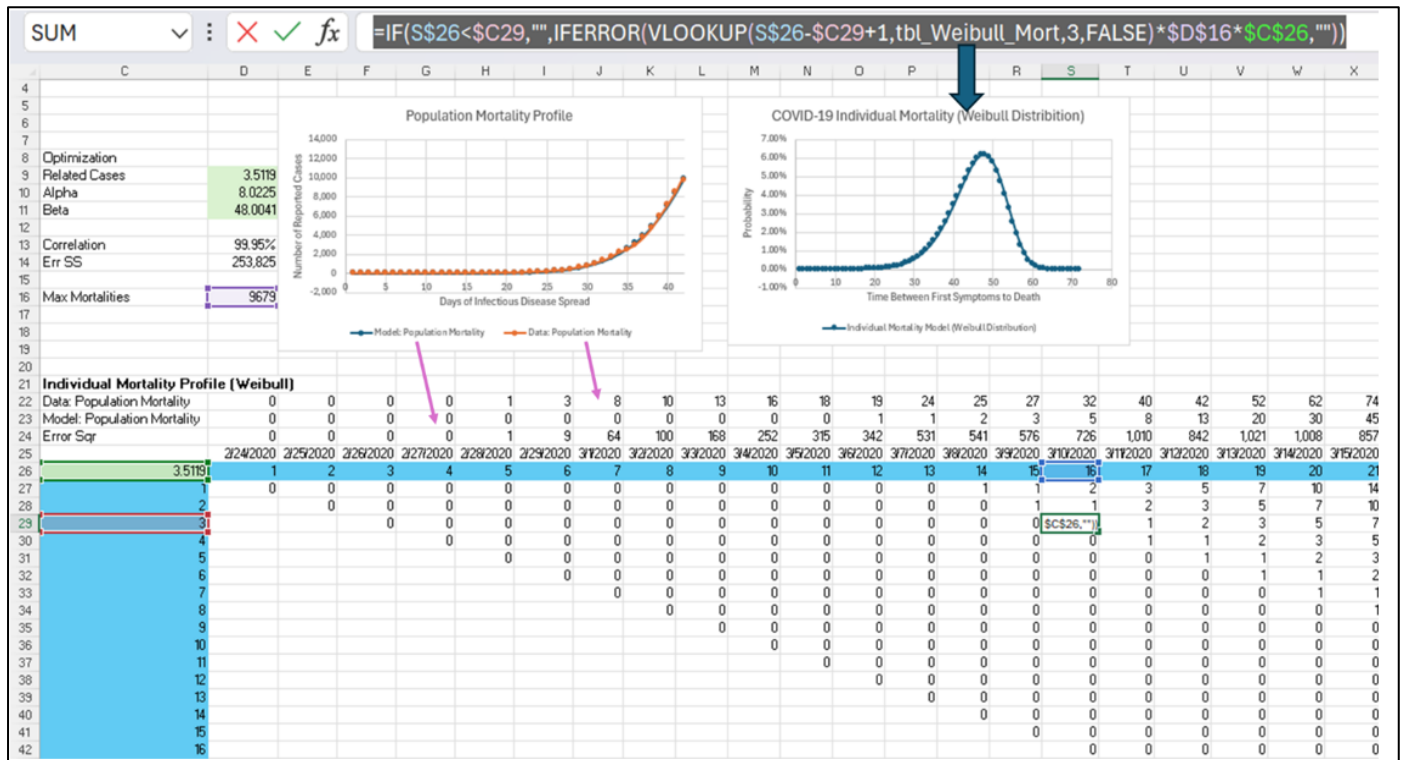
Fig 1 Collated Distribution Model for Mortalities Using MS Excel

Table 1 US COVID-19 Data [6]

| Day | Date | Infected | Cum | Deaths | Cum |
|---|---|---|---|---|---|
| 0 | 2/23/2020 | 35 | 35 | 0 | 0 |
| 1 | 2/24/2020 | 18 | 53 | 0 | 0 |
| 2 | 2/25/2020 | 5 | 58 | 0 | 0 |
| 3 | 2/26/2020 | 2 | 60 | 0 | 0 |
| 4 | 2/27/2020 | 1 | 61 | 0 | 0 |
| 5 | 2/28/2020 | 6 | 67 | 1 | 1 |
| 6 | 2/29/2020 | 5 | 72 | 2 | 3 |
| 7 | 3/1/2020 | 22 | 94 | 5 | 8 |
| 8 | 3/2/2020 | 18 | 112 | 2 | 10 |
| 9 | 3/3/2020 | 22 | 134 | 3 | 13 |
| 10 | 3/4/2020 | 35 | 169 | 3 | 16 |
| 11 | 3/5/2020 | 71 | 240 | 2 | 18 |
| 12 | 3/6/2020 | 104 | 344 | 1 | 19 |
| 13 | 3/7/2020 | 116 | 460 | 5 | 24 |
| 14 | 3/8/2020 | 121 | 581 | 1 | 25 |
| 15 | 3/9/2020 | 176 | 757 | 2 | 27 |
| 16 | 3/10/2020 | 290 | 1047 | 5 | 32 |
| 17 | 3/11/2020 | 245 | 1292 | 8 | 40 |
| 18 | 3/12/2020 | 424 | 1716 | 2 | 42 |
| 19 | 3/13/2020 | 532 | 2248 | 10 | 52 |
| 20 | 3/14/2020 | 724 | 2972 | 10 | 62 |
| 21 | 3/15/2020 | 707 | 3679 | 12 | 74 |
| 22 | 3/16/2020 | 965 | 4644 | 27 | 101 |
| 23 | 3/17/2020 | 1454 | 6098 | 23 | 124 |
| 24 | 3/18/2020 | 2567 | 8665 | 31 | 155 |
| 25 | 3/19/2020 | 5438 | 14103 | 56 | 211 |
| 26 | 3/20/2020 | 5489 | 19592 | 67 | 278 |
| 27 | 3/21/2020 | 7169 | 26761 | 73 | 351 |
| 28 | 3/22/2020 | 8336 | 35097 | 122 | 473 |
| 29 | 3/23/2020 | 9739 | 44836 | 128 | 601 |

| 30 | 3/24/2020 | 10523 | 55359 | 205 | 806 |
| 31 | 3/25/2020 | 13890 | 69249 | 259 | 1065 |
| 32 | 3/26/2020 | 16755 | 86004 | 257 | 1322 |
| 33 | 3/27/2020 | 18747 | 104751 | 404 | 1726 |
| 34 | 3/28/2020 | 19722 | 124473 | 479 | 2205 |
| 35 | 3/29/2020 | 17939 | 142412 | 312 | 2517 |
| 36 | 3/30/2020 | 21801 | 164213 | 489 | 3006 |
| 37 | 3/31/2020 | 25599 | 189812 | 753 | 3759 |
| 38 | 4/1/2020 | 26429 | 216241 | 960 | 4719 |
| 39 | 4/2/2020 | 29420 | 245661 | 1232 | 5951 |
| 40 | 4/3/2020 | 32500 | 278161 | 1257 | 7208 |
| 41 | 4/4/2020 | 34128 | 312289 | 1324 | 8532 |
| 42 | 4/5/2020 | 25827 | 338116 | 1147 | 9679 |

➢ *What are Collated Distributions?*

Fig. 1 depicts a set of individual probability distributions collated into a matrix, where the x-axis or rows are individual probability distribution of events by age, days, or relevant time periods, of events. Each row or cohort starting at the next time-period or day. The y-axis or columns are the next time periods cases per the case probabilistic behavior. At the top is the sum population rows for that time period.

Each row or cohort can be considered as a spontaneous set of cases. For example, as I recall, Google's Android phone COVID-19 tracking did not produce any useful results. Consider a person walking at 4 mph along a city street, with a wind speed of 10 mph, i.e., the relative wind speed ranges from 6 to 14 mph. Let assume an exhale of 2 seconds. In 2 seconds, the exhale could have travelled between 15 and 36 feet. Or if the Google tracking was between phones at most 6 feet apart, the results would definitely be inconclusive, unless people were inside a closed ventilated building. Therefore, it is not possible to realistically model spontaneous infections.

Thus, each cohort are spontaneous infections for the next time periods.

Thus, Collated Distributions are a good reasonable approach to modeling IDS.

➢ *Collated Distribution Infectious Disease Spread (CD-IDS) Basics*

The Collated Distribution Infectious Disease Spread (CD-IDS) model presented in this paper is different other IDS models [5]. See Fig. 2. It is based on the axiom that infectious diseases can be characterize by 3 probability distributions, Individual Infectiousness Profile, Individual Mortality Profile, and Individual Survival Profile (Survival is defined as Infected – Mortalities) derived from a minimum of 2 data sets consisting of new infected, and new deaths. To determine the 3 individual probability distributions (profiles), the starting distributions was obtained after researching [7-10] possible distributions in 2020 publications.
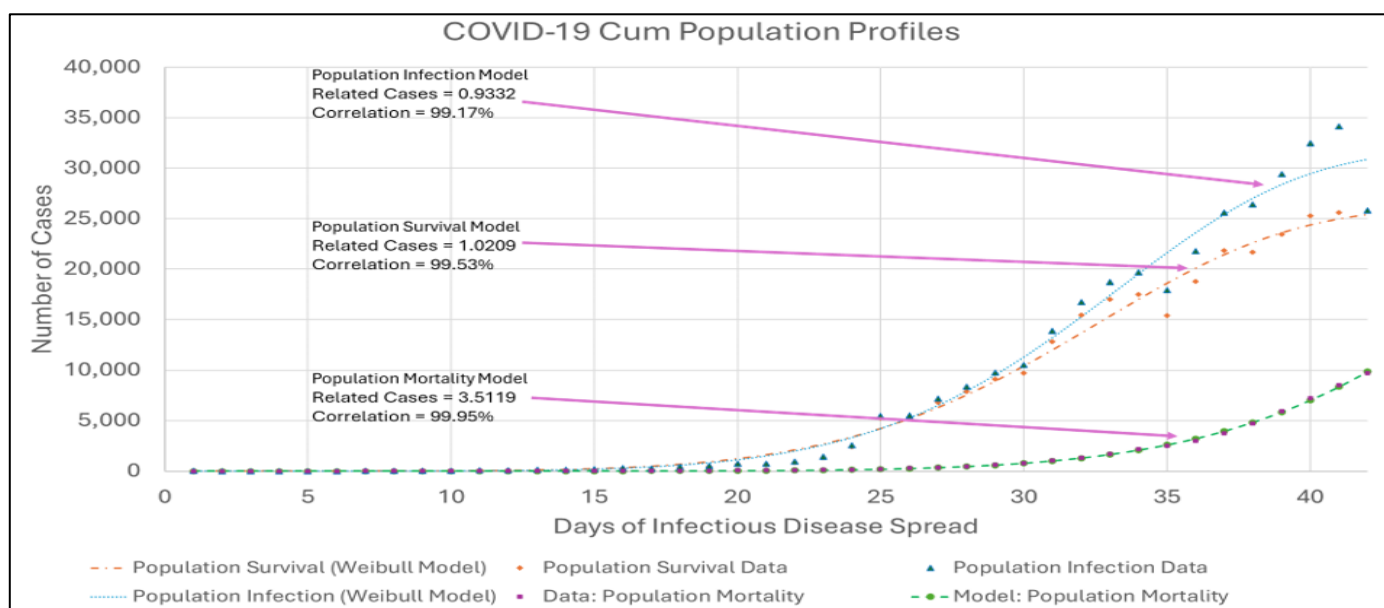


Fig 2 COVID-19 Population Profiles, Infection, Survival & Mortality Between Jan 24, 2020 to Apr 5, 2020

This CD-IDS model uses a new class of solutions, Collated Distributions [1 & 2], that the population's

cumulative cases are derived from the sum of the single individual distribution of cases (1). A rudimentary version

was first pioneered [11 & 12] in the financial services industry, but Collated Distributions are significantly different from these with the use of a row single distribution that

$$\hat{n}_j = N \sum_i (P_{i,j}) \qquad \text{N is cell D16 in Fig. 1} \qquad (1)$$

Where $j$ = day or column of IDS (row 26 in Fig.1), $i$ = cohort or row (column C in Fig. 1), $P_{i,j}$ is the probability of case, and $N$ is the number of cases as of the last day of IDS.

$$\hat{N}_j = \sum_j \hat{n}_j \qquad \text{row 23 in Fig. 1} \qquad (2)$$

$$\hat{n}_{i,j} = R_C N P_{i,j} \qquad \text{cells D21 to X42 in Fig. 1} \qquad (3)$$

The Related Cases, $R_C$, cell C26 or D5 in Fig. 1, is the number of cases associated with each probability and converts this probability matrix $P_{i,j}$ into a population count matrix $\hat{n}_{i,j}$. $R_C$ is the number of related cases modeled to provide a best fit. For Individual Mortality $R_C$ = 3.519 related cases that best fits the data.

## II. METHOD

In this study, Wilcoxon Regression [13, 14] was used, as regression results can lead to technique breakdown, and provide good Goodness-Of-Fit measures but biased modeling if two conditions are not met,

$$Min(SS_E) = Min\left(\sum_j (N_j - \hat{N}_j)^2\right) \qquad \text{cell D14 in Fig. 1} \qquad (4)$$

In this study, MS Excel's GRG Nonlinear was used to find the minimum sum of squared error $SS_E$. Weibull model is defined as a named table, tbl_Weibull_Mort, whose $\alpha$ and $\beta$ (cells D10 & D11, respectively, in Fig.1) are the changing variables for the Weibull Distribution.

"reflects" in the columns. The modeled cases at day j, $\hat{n}_j$ is given by,

See Fig. 1. Thus, the Population for a specific case type (infection, survival, or mortality) up to day $j$, is given by,

Given,

➢ The errors must be Normally distributed.
➢ Heteroscedasticity must not be present, or errors are correlated to some factors (usually x-axis factor).

Wilcoxon Regression [13, 14] based on the Wilcoxon Two-Sample Test, is less sensitive to technique breakdown as it does not assume Normality. It can be summarized as the minimization of the sum of squared error $SS_E$ (cell D14 of Fig. 1) between each data point $j$ of value $N_j$ and model value $\hat{N}_j$, by changing the value of the $R_C$, $\alpha$ and $\beta$ (cells D9, D10 & D11, respectively, in Fig.1),

Note, there are two a necessary check, (i) the Survival distribution cannot start before the Infection distribution, and (ii) the Mortality distribution cannot start before the Infection distribution.
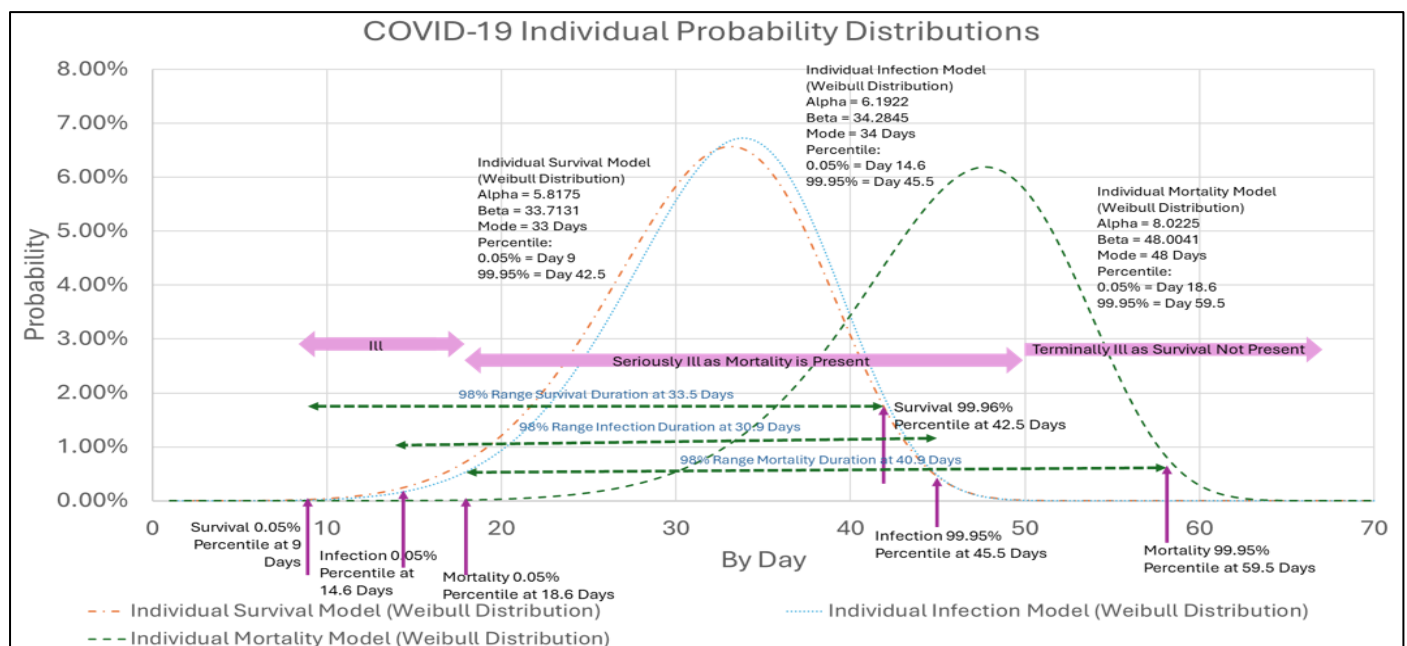


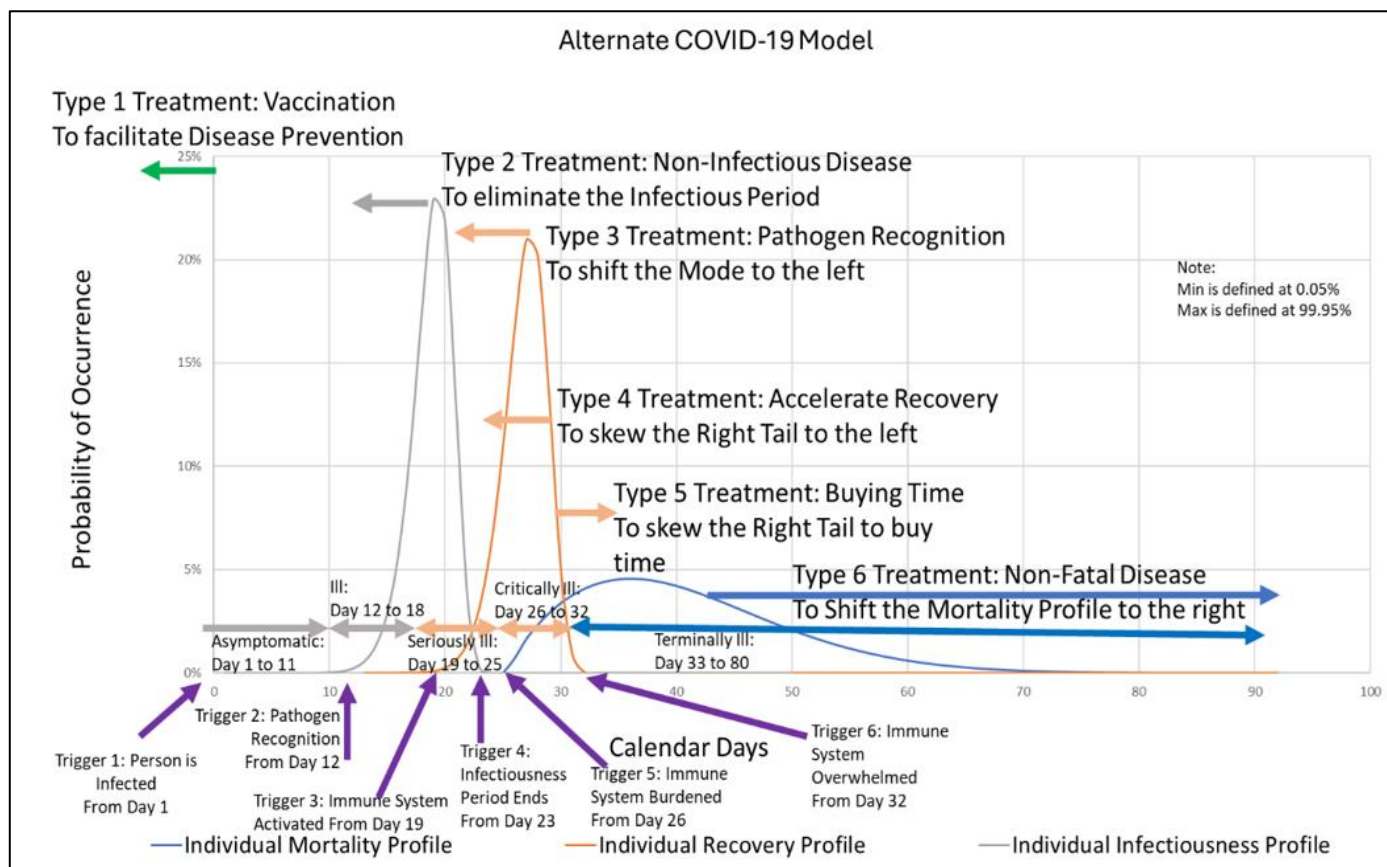Fig 3 COVID-19 Individual Probabilities Prior to the Introduction of US COVID-19 Health Policies

Fig 4 An Alternate COVID-19 Model (2020)

## III. DISCUSSION

➤ *Discussion: Interpretation*

A key part of any analysis is the interpretation of the results. Fig. 3 shows the Individual Infections, Survival & Mortality probability density functions, and depicts the 0.05 & 99.95 percentiles of cum probabilities.

- Mortality ranges between day 18.6 to day 59.5 or a range of 40.9 days.
- Infections ranges between day 14.6 to day 45.5 or a range of 30.9 days.
- Survival ranges between day 9 and day 42.5 or a range of 33.5 days.

Both the Individual Infection & Survival overlap each other. That is a patient is responding to the infection almost immediately. And that it takes about 9 to 15 days for the viral load to become infectious to others. With infections being active for 31 days, would suggest that, if necessary, an effective quarantine period should be 31 days.

Fig. 3 shows a proposed disease stages,

- Ill, defined as Mortality not being present.
- Seriously Ill, defined as Mortality being present.
- Terminally Ill, defined as Survival not being present.

These disease stages' durations are based on the 0.05 & 99.95 percentiles because human lives are at stake but could

be changed to other percentiles based on experience in the field.

Fig. 4 depicts an alternate model based on a more complex modeling (but the author lost this model after his PC crashed). However, it is a good example of the interpretation of results. From a statistical perspective, it shows (i) the three individual probability distributions (infection, survival / recovery & mortality) (ii) the disease stages, (iii) the 6 triggers present at each disease stage, and (iv) how to determine a cure strategy i.e. what is the intention of the research to find a cure?

➤ *Discussion: Critique*

The IHME model is essentially a variation of current epidemiological theory [15] (5),

$$I_{final} = N - (N - I_0)e^{\left(\frac{-R_0 I_{final}}{N}\right)} \tag{5}$$

Where $N$ is the size of the susceptible population, $I_0$ the number of primary cases infected, and $I_{final}$ the expected final size of an outbreak of an infectious disease. If $R_0 > 1$, each primary infection will, on average, generate more than one secondary infection. If $R_0 < 1$, each primary case will, on average, fail to replace itself (although short chains of transmission are still possible) and each single introduction will lead to no more than a minor outbreak.

The IHME model [15] (6) at its essence is the ln() function that has been found to be produce unacceptable forecast results [16]. It suffers from two major problems.

The cumulative mortalities $y^t_j$ at time $t$ in location $j$ in the IHME mortality model (6) is given by,

$$log\left(y^t_j\right) = \frac{p_j}{2}\left(1 + \frac{2}{\sqrt{\pi}}\int_0^{\alpha_j(t-\beta_j)} e^{(-\tau^2)}d\tau\right) + \varepsilon_{t,j} \qquad (6)$$

$$y^t_j = e^{\frac{p_j}{2}\left(1+\frac{2}{\sqrt{\pi}}\int_0^{\alpha_j(t-\beta_j)} e^{(-\tau^2)}d\tau\right)+\varepsilon_{t,j}} + \delta_{t,j} = \left(\Delta_{t,j}\right)e^{\frac{p_j}{2}\left(1+\frac{2}{\sqrt{\pi}}\int_0^{\alpha_j(t-\beta_j)} e^{(-\tau^2)}d\tau\right)} + \delta_{t,j} \qquad (8)$$

That is,

$$\Delta_{t,j} = e^{\varepsilon_{t,j}} \neq N(0,\sigma) \qquad (9)$$

Where $\sigma$ is some standard deviation value and $\delta_{t,j}$ is the true model errors of the true model of population statistic. (8) shows that $\Delta_{t,j}$ acts as a multiplier distorting $y^t_j$, the statistic one is interested in. For example, for $n = 42$ days, mean $\bar{y}_t$ of the known infected (dependent variable) is 8,050 but mean of this dependent variable as a ln() function $\sum[\ln(y_t)]/n$ is 6.6059 which is 739.4, or,

$$\bar{y}_t \neq e^{\sum \ln(y_t)/n} \qquad (10)$$

There is a very big difference between 8,050 and 739.4. Thus, the minimization of the error sum of squares of the transformed data is not the same as minimization of error sum of squares of the data.

$$\Delta_{cum(t,j)} = N(0,\sigma) \qquad \text{where standard deviation } \sigma\rightarrow0 \text{ as } t\rightarrow\infty \qquad (11)$$

Or as $t\rightarrow\infty$ the standard deviation of the model cum errors $\sigma\rightarrow0$ as the errors cancel out. This is another form of heteroscedasticity. Therefore, the residual sum of squares $\rightarrow0$ and thus, the F-Ratio skyrockets but the model usually gives misleading results.

## IV. CONCLUSION

This paper has shown that it is possible to statistically determine statistical disease properties and treatment strategies. The Collated Distributions models should lead to a better understanding and parametrization of generic population spreads and therefore, make informed opinions on how public health management should be conducted and possibly which drug treatment is likely to succeed.

$$\varepsilon_{t,j} \sim N(0, V_t) \qquad (7)$$

Where level $p$ controls the ultimate level, slope $\alpha$ controls speed of infection, and inflection $\beta$ is the time at which the rate of change is maximal.

The first model problem. Even though the log model errors $\varepsilon_{t,j}$ are Normally distributed $N(0,V_t)$, this is the log of the data errors $\Delta_{t,j}$ with respect to the exponential function (8),

Heteroscedasticity [16] is present when the model (dependent) errors exhibit a monotonic behavior with respect to the dependent variable i.e. errors get larger or smaller, even after excluding outliers or is some function of the dependent variable. This is due to (i) missing factors (variables) in the context structure, (ii) incorrectly modeled factors/independent variables and/or (iii) when an independent variable represents more than one underlying factor. Heteroscedasticity is a major problem [17] when the range of values of the dependent variable is very large, as is the case with modeling infectious disease spread and therefore the dire need for a different approach.

The second problem with the IHME's models is the use of cums. From the author's 40-years of working with data, cum models give extremely good fit because cums cancel, not minimize, the noise in the data. However, cums can substantially bias the non-cum model results, the statistic one is interested in. That is, the cum model errors $\Delta_{cum(t,j)}$ is,

## REFERENCES

[1]. Benjamin Solomon, Solomon's Method for Collated Distributions Used in Mortgage-Backed Securities, SeekingAlpha.com, April 20, 2020. https://seekingalpha.com/article/4338509-solomons-method-for-collated-distributions-used-in-mortgage-backed-securities

[2]. Benjamin Solomon, A Critique of Dodd-Frank: Forecasting Securitized Mortgage Credit & Default Risk, Scholar's Press, 2021, https://www.morebooks.shop/store/gb/book/a-critique-of-dodd-frank/isbn/978-613-8-95350-0

[3]. IHME Staff, CurveFit, https://ihmeuw-msca.github.io/CurveFit/methods/#statistical-model, accessed 04/23/2020.

[4]. Mark E.J. Woolhouse, Daniel T. Haydon, Rustom Antia, Emerging pathogens: the epidemiology and evolution of species jumps, TRENDS in Ecology and Evolution Vol.20 No.5 May 2005.https://www.cell.com/action/showPdf?pii=S0169-5347%2805%2900038-8

[5]. Caroline E. Waltersa1, Margaux M. I. Meslébc, Ian M. Halla, Modelling the global spread of diseases: A review of current practice and capability, Epidemics, Volume 25, December 2018.https://www.sciencedirect.com/science/article/pii/S1755436517301135

[6]. 1Point3Acres.com (also used by the CDC and John Hopkins in 2020), accessed between 2/23/2020 and 4/5/2020.

[7]. Weier Wang, Jianming Tang, Fangqiang Wei, Updated understanding of the outbreak of 2019 novel coronavirus (2019-nCoV) in Wuhan, China, Journal of Medical Virology, 29 January 2020. https://onlinelibrary.wiley.com/doi/full/10.1002/jmv.25689?af=R

[8]. Robert Verity, Lucy C Okell, Ilaria Dorigatti, Peter Winskill, Charles Whittaker, Natsuko Imai, Gina Cuomo-Dannenburg, Hayley Thompson, Patrick G T Walker, Han Fu, Amy Dighe, Jamie T Griffin, Marc Baguelin, Sangeeta Bhatia, Adhiratha Boonyasiri, Anne Cori, Zulma Cucunubá, Rich FitzJohn, Katy Gaythorpe, Will Green, Arran Hamlet, Wes Hinsley, Daniel Laydon, Gemma Nedjati-Gilani, Prof Steven Riley, Sabine van Elsland, Erik Volz, Haowei Wang, Yuanrong Wang, Xiaoyue Xi, Prof Christl A Donnelly, Prof Azra C Ghani, Prof Neil M Ferguson, Estimates of the severity of coronavirus disease 2019: a model-based analysis, The Lancet, Infectious Disease, March 30, 2020. https://www.thelancet.com/pdfs/journals/laninf/PIIS1473-3099(20)30243-7.pdf

[9]. Qun Li, M.Med., Xuhua Guan, Ph.D., Peng Wu, Ph.D., Xiaoye Wang, M.P.H., Lei Zhou, M.Med., Yeqing Tong, Ph.D., Ruiqi Ren, M.Med., Kathy S.M. Leung, Ph.D., Eric H.Y. Lau, Ph.D., Jessica Y. Wong, Ph.D., Xuesen Xing, Ph.D., Nijuan Xiang, M.Med., Yang Wu, M.Sc., Chao Li, M.P.H., Qi Chen, M.Sc., Dan Li, M.P.H., Tian Liu, B.Med., Jing Zhao, M.Sc., Man Liu, M.Sc., Wenxiao Tu, M.Med., Chuding Chen, M.Sc., Lianmei Jin, M.Med., Rui Yang, M.Med., Qi Wang, M.P.H., Suhua Zhou, M.Med., Rui Wang, M.D., Hui Liu, M.Med., Yinbo Luo, M.Sc., Yuan Liu, M.Med., Ge Shao, B.Med., Huan Li, M.P.H., Zhongfa Tao, M.P.H., Yang Yang, M.Med., Zhiqiang Deng, M.Med., Boxi Liu, M.P.H., Zhitao Ma, M.Med., Yanping Zhang, M.Med., Guoqing Shi, M.P.H., Tommy T.Y. Lam, Ph.D., Joseph T. Wu, Ph.D., George F. Gao, D.Phil., Benjamin J. Cowling, Ph.D., Bo Yang, M.Sc., Gabriel M. Leung, M.D., and Zijian Feng, M.Med., Early Transmission Dynamics in Wuhan, China, of Novel Coronavirus–Infected Pneumonia, The New England Journal of Medicine, January 29, 2020. https://www.nejm.org/doi/full/10.1056/NEJMoa2001316

[10]. A. C. Ghani, C. A. Donnelly, D. R. Cox, J. T. Griffin, C. Fraser, T. H. Lam, L. M. Ho, W. S. Chan, R. M. Anderson, A. J. Hedley, G. M. Leung, Methods for Estimating the Case Fatality Ratio for a Novel, Emerging Infectious Disease, American Journal of Epidemiology, Volume 162, Issue 5, 1 September 2005. https://academic.oup.com/aje/article/162/5/479/82647

[11]. Von Furstenberg, George M., 1969, "Default risk on FHA-insured home mortgages as a function of the terms of financing: a quantitative analysis", The Journal of Finance 24-3: 459-477.https://www.jstor.org/stable/2325346?seq=1

[12]. Esaki, L'Heureux & Snyderman, Commercial Mortgage Update, Real Estate Finance, The Quarterly Review of Commercial Finance Techniques, Spring 1999, Vol. 16, No. 1.

[13]. Solomon BT, Real World Data Modeling: Applications in Statistics, Physics & Medicine, Scholar's Press, 2021. https://www.morebooks.shop/store/gb/book/real-world-data-modeling/isbn/978-613-8-95346-3

[14]. Benjamin T. Solomon; Dagmar Horvath. "Comparing a Type 2 Diabetic to Non-Diabetics' Blood Glucose Levels." Volume. 11 Issue.1, January 2026 International Journal of Innovative Science and Research Technology (IJISRT) 134-145 https://doi.org/10.38124/ijisrt/26jan107

[15]. IHME Staff, CurveFit, https://ihmeuw-msca.github.io/CurveFit/methods/#statistical-model, accessed 04/23/2020.

[16]. Bucevska V. (2011) Heteroscedasticity. In: Lovric M. (eds) International Encyclopedia of Statistical Science. Springer, Berlin, Heidelberg https://link.springer.com/referenceworkentry/10.1007%2F978-3-642-04898-2_628#toc

[17]. Andrew C. McCarthy, COVID-19 Projection Models Are Proving to Be Unreliable, National Review Institute, 9 April 2020. https://www.nationalreview.com/corner/coronavirus-pandemic-projection-models-proving-unreliable/