

1 The Ongoing COVID-19 Epidemic Curves Indicate Initial Point Spread in China with Log-Normal 2 Distribution of New Cases Per Day with a Predictable Last Date of the Outbreak

3 **Stefan Olsson^{1,2*}**

4 ¹State Key Laboratory of Ecological Pest Control for Fujian and Taiwan Crops, College of Plant Protection

5 ²Plant Immunity Center, Haixia Institute of Science and Technology, College of Life Science

6 Fujian Agriculture and Forestry University, No.15 Shangxiadian Road, Cangshan District, Fuzhou City, Fujian Province, China. P.C. 350002

7

8

9 * E-mail stefan@olssonstefan.com or stefan.olsson@fafu.edu.cn

10

11 **Abstract**

12 During an epidemic outbreak it is useful for planners and responsible authorities to be able to plan
13 ahead to estimate when an outbreak of an epidemic is likely to ease and when the last case can be
14 predicted in their area of responsibility. Theoretically this could be done for a point source epidemic
15 using epidemic curve forecasting. The extensive data now coming out of China makes it possible to test
16 if this can be done using MS Excel a standard spreadsheet program available to most offices. The
17 available data is divided up for whole China and the different provinces. This and the high number of
18 cases makes the analysis possible. Data for new confirmed infections for Hubei, Hubei outside Wuhan,
19 China excluding Hubei as well as Zhejiang and Fujian provinces all follow a log-normal distribution that
20 can be used to make a rough estimate for the date of the last new confirmed cases in respective areas.

21

22 **Introduction**

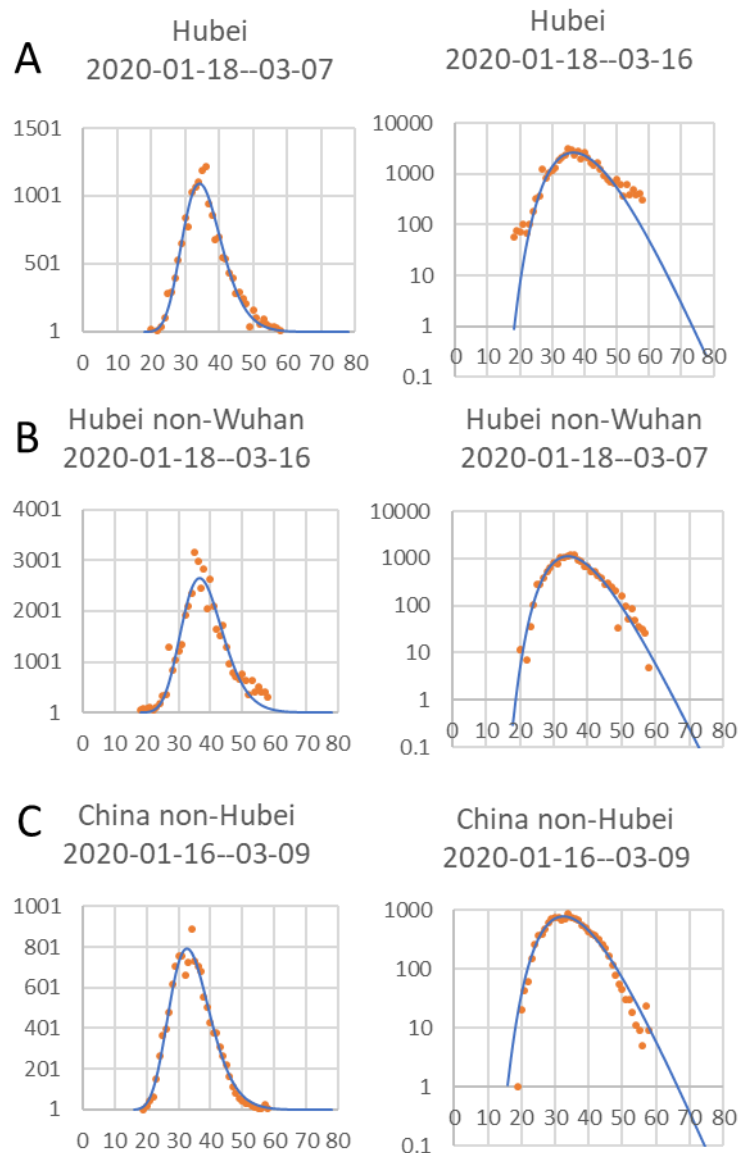
23 In epidemics starting as a point source the number of new cases often follows a log-normal distribution
24 or more precisely a Poisson-Gamma distribution. How this distribution will develop over time can
25 theoretically be determined by fitting a log-normal distribution equation to the data for new cases per
26 day are reported. The estimate will of course be more accurate the further into the outbreak. A literally
27 “breaking point” for the accuracy of the estimate for the end of the outbreak comes after the number of
28 new cases per day have reached its peak. From there on the estimate should be better and better. Here
29 a simple method that could be used without access to special resources for getting such estimates after
30 the peak has been reached is presented using data from the ongoing COVID-19 epidemic in China.

31

32 **Results and discussion**

33 A log normal distribution can be relatively nicely fitted all data sets (Fig 1&2). When using a log scale for
34 the Y-axis it is apparent there are deviations in the early dates especially for Hubei (Fig 1A). This could be
35 caused by a lag in detection of new cases in the beginning of the outbreak. The deviations in the latest
36 dates can have many different causes like changing criteria for new cases, or simply a backlog in cases
37 confirmation due to highly stressed health care system in the worst hit city Wuhan. Both the data from
38 Hubei outside Wuhan (Fig 1B) and China outside Hubei (Fig 1C) on the other hand closely follows a log
39 normal distribution. To see if the same relationships holds also outside Hubei, two provinces with quite
40 different number of cases, Zhejiang with many cases and Fujian with few cases, was also tested (Fig 2).

41



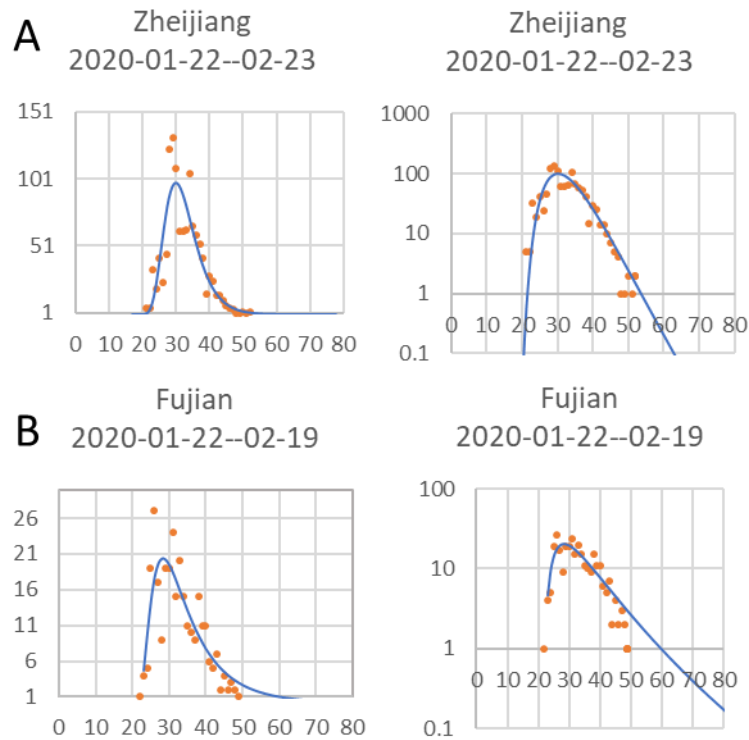
42

43 **Figure 1. Log normal distribution of new confirmed cases for each day since 1 Jan 2020 Hubei, Hubei-nonWuhan**
 44 **and in reest of China.** The Log of day values with start on the first day a case could have been confirmed was used
 45 curve fitting although here in the plot the actual number of days since 1st January was used as X-axis. Number of
 46 new confirmed cases per day and fitted curve (left) and Log number of new cases per day to show start and stop
 47 days (right). Headings shows estimated dates for 1st and last confirmed case. Y axes both to the left and right start
 48 at 1 to highlight the first and last case.

49

50 In Zhejiang the outbreak followed the general pattern very closely (Fig 2A) but for the much smaller
 51 outbreak in Fujian (Fig 2B) the number of cases per day dropped more than the model for the last
 52 days. This is caused by the approximation to log-normal distribution instead of a Poisson distribution that
 53 is more correct for data with few cases (Gonzales-Barron and Butler, 2011) but more difficult to handle

54 using standard Excel curve fitting. This discrepancy mean that the last new infection date will be
 55 overestimated especially for limited outbreaks like the one in Fujian province. From planning point of
 56 view it should however be safer to overestimate the length of the outbreak than underestimate it. A
 57 fairly good estimate of the last data could be done as soon as the number of new confirmed cases per
 58 day started to decrease.



59

60 **Figure 2. Log normal distribution of new confirmed cases for each day since 1 Jan 2020 in two provinces with**
 61 **relatively high numbers of cases, Zhejiang with high numbers and Fujian with low numbers.** The Log of day values
 62 with start on the first day a case could have been confirmed was used curve fitting although here in the plot the
 63 actual number of days since 1st January was used as X-axis. Number of new confirmed cases per day and fitted
 64 curve (left) and Log number of new cases per day to show start and stop days (right). Headings shows estimated
 65 dates for 1st and last confirmed case. Y-axes both to the left and right start at 1 to highlight the first and last case.

66 The estimated start date for when new cases could have been confirmed caused by community spread
 67 was for Hubei and Wuhan the 18th January while outside Hubei the data indicate a 2 day earlier start if
 68 the disease behaved similarly. This is a bit surprising but could indicate that the disease was brought to
 69 Wuhan city and Hubei province from a less populated area and found good conditions for spread in
 70 Wuhan. The estimated start dates for when new cases could be confirmed in the two provinces
 71 Zhejiang and Fujian were both the 22nd January a few days later than in the epicenter for the outbreak.

72

73 Conclusion

74 Plotting new confirmed cases per day against time can be used during a large point source epidemic
 75 outbreak to relatively early after the peak in new cases determine a likely last date for new cases. Such

76 information should be useful to people in charge for planning how to allocate resources. The
77 information will also be available when resources are as most stretched with a large number of active
78 cases just after the peak in number of new cases per day, In addition, if the data continue to fit the curve
79 for a point source outbreak in one area there has most likely been no new introduction of cases or any
80 change to the virus or the likelihood that a person becomes infected within that area. The latter seems
81 to be the case for the COVID-19 outbreak in China 2019-2020 pointing to that the quarantining
82 measures stopping further spread between provinces and cities after the first few days of person-to-
83 person transfer have worked efficiently.

84 **Methods**

85 Official referred to data for the COVID-19 outbreak in China is collected at a Wikipedia page
86 (Anonymous, 2020). Since the kind of analysis here presented is a relatively simple analysis it should be
87 possible to do for anyone using a standard program Microsoft Excel with the standard available Solver
88 plugin for data handling and curve fitting. The logarithm of number of days since the estimated start of
89 the epidemic outbreak were used for fitting a normal distribution equation to the data but in the figures
90 the data was plotted against the non-logged day number with day 1 on the 1st January to ease in
91 determining the actual dates from readings on the X-axis and the values in the spreadsheet files.

92 The MS Excel file used for this analysis is available as Supplementary file and can easily be modified to
93 be used with other data to relatively early after the peak in new confirmed cases be able to predict the
94 end of an epidemic outbreak with a definite starting point having a “first case”.

95

96 **Acknowledgement**

97 When back in my home country Sweden I had to decide when to return to China after the winter break
98 for the Chinese New Year (Spring Festival), I decided to look at the epidemiology data since I have been
99 working with biological control trying to cause epidemics in fungal pathogens attacking plants. I thought
100 of looking for data about the COVID-19 outbreak to be able to determine a time and a route back that
101 limit the chances for me to catch the infection and bring it to my workplace. I found the very good
102 Wikipedia entry I refer to in the methods and would like to thank everyone that has contributed to edit
103 that site. Finally, I want to acknowledge my employer Fujian Agriculture and Forestry University that
104 makes it possible for me to do research in China.

105 **Supplemental file**

106 “Corona model final.xlsx” is a supplemental file containing all data and calculations. In addition, the file
107 also contain instructions for how to use it to fit new data to make predictions.

108 **References**

109 Anonymous (2020). Timeline of the 2019–20 coronavirus outbreak. *Wikipedia*. Available at:
110 https://en.wikipedia.org/wiki/Timeline_of_the_2019%E2%80%9320_coronavirus_outbreak#Case_statistics
111 [e_statistics](https://en.wikipedia.org/wiki/Timeline_of_the_2019%E2%80%9320_coronavirus_outbreak#Case_statistics) [Accessed March 1, 2020].

112 Gonzales-Barron, U., and Butler, F. (2011). A comparison between the discrete Poisson-gamma and
113 Poisson-lognormal distributions to characterise microbial counts in foods. *Food Control* 22,
114 1279–1286. doi:10.1016/j.foodcont.2011.01.029.

115