

# Infinite variance of U.S. COVID-19 cases & deaths, & Taylor's law of heavy-tailed data

Joel E. Cohen, [cohen@rockefeller.edu](mailto:cohen@rockefeller.edu)

Rockefeller University & Columbia University

ICTP Workshop on Quantitative Human Ecology,  
Trieste 2022-07-28

# Plan

Fluctuation scaling, variance function

$$\textit{variance} = f(\textit{mean})$$

Taylor's law:  $\textit{variance} = a(\textit{mean})^b$

$$\log(\textit{variance}) = \log a + b \log(\textit{mean})$$

Heavy tails & regular variation

$$\Pr(X > x) = L(x)x^{-\alpha}, 0 < \alpha < 2$$

COVID-19 in US

Taylor's law, infinite variance

# Plan

→ Fluctuation scaling, variance function

$$\textit{variance} = f(\textit{mean})$$

Taylor's law:  $\textit{variance} = a(\textit{mean})^b$

$$\log(\textit{variance}) = \log a + b \log(\textit{mean})$$

Heavy tails & regular variation

$$\Pr(X > x) = L(x)x^{-\alpha}, 0 < \alpha < 2$$

COVID-19 in US

Taylor's law, infinite variance

# Variance function

**Population:** Given a non-empty family of random variables  $\{X(s)\}_{s \in S}$ , if each  $X(s)$  has finite mean  $E(X(s))$  & finite variance  $Var(X(s))$ , the **population variance function**  $f$  says:  $Var(X(s)) = f(E[X(s)])$ ,  $\forall s \in S$ .

**Sample:** Sample of size  $n > 1$  is a set  $\{X_1(s), \dots, X_n(s)\}$  of  $n$  iid copies of  $X(s)$ , with sample mean  $\bar{X}_n(s) := (X_1(s) + \dots + X_n(s))/n$ , sample variance  $s_n^2(s)$ . The **sample variance function**  $f_n$  says:  $s_n^2(s) \approx f_n(\bar{X}_n(s))$ .

# Plan

Fluctuation scaling, variance function

$$\text{variance} = f(\text{mean})$$

→ Taylor's law:  $\text{variance} = a(\text{mean})^b$

$$\log(\text{variance}) = \log a + b \log(\text{mean})$$

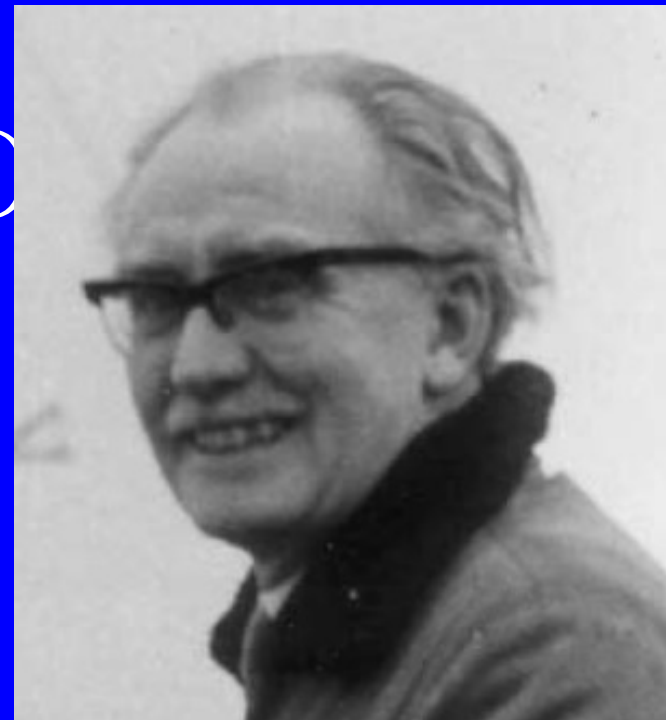
Heavy tails & regular variation

$$\Pr(X > x) = L(x)x^{-\alpha}, \alpha \in (0, 2)$$

COVID-19 in US

Taylor's law, infinite variance

Lionel Roy Taylor  
(1924–2007)



# Taylor's law(s)

**Population TL** holds if each nonnegative random variable  $X(s)$ ,  $\forall s \in S$ , has finite, positive mean & variance &  $\exists a > 0, b$  such that

$$\text{Var}(X(s)) = a\{E(X(s))\}^b.$$

**Sample TL** holds if samples with mean  $\bar{X}_n(s)$ , variance  $s_n^2(s)$  obey,  $\forall s \in S$ , for some  $a > 0, b$ ,  $\log s_n^2(s) \approx \log a + b \log \bar{X}_n(s)$  or

$$\frac{s_n^2(s)}{\{\bar{X}_n(s)\}^b} \approx a > 0.$$

For sample TL, there is no requirement that mean or variance exist or are finite.

# TL data structure: multiple samples, each with multiple observations

Sample number →	s=1	s=2	s=3	s=...
Population size or density in units (quadrats, plots, transects, counties, states, years, days)	$x_{11}$	$x_{12}$	$x_{13}$	$x_{...}$
	$x_{21}$	$x_{22}$	$x_{23}$	...
	$x_{31}$	$x_{32}$	$x_{33}$	...
		$x_{42}$	$x_{43}$	...
		$x_{52}$		...
Mean (weighted)	$m_1$	$m_2$	$m_3$	$m_{...}$
Variance (weighted)	$v_1$	$v_2$	$v_3$	$v_{...}$



# Tornado

USA has more tornadoes than any other country. (Lloyd's)

Category F5 tornado viewed from the southeast as it approached Elie, Manitoba on Friday, June 22nd, 2007. Justin Hobson

2/3/2024

Joel E. Cohen

Creative Commons Attribution-Share Alike 3.0 Unported, 2.5 Generic, 2.0 Generic and 1.0 Generic license. GNU Free Documentation License, Version 1.2 or any later version published by the Free Software Foundation



# TL for tornadoes: size of outbreaks by calendar year

Year→	1954	1955	...	2014
Number of F1+ tornadoes per outbreak ( $\geq 6$ consecutive tornadoes with $\leq 6$ hour gap)	$x_{11}$	$x_{12}$	$x_{13}$	$x_{...}$
	$x_{21}$	$x_{22}$	$x_{23}$	...
	$x_{31}$	$x_{32}$	$x_{33}$	...
		$x_{42}$	$x_{43}$	...
		$x_{52}$		...
Mean (weighted)	$m_1$	$m_2$	$m_3$	$m_{...}$
Variance (weighted)	$v_1$	$v_2$	$v_3$	$v_{...}$

# Outbreaks ( $\geq 6$ tornadoes) cause most damage.

Outbreak is defined as  $\geq 6$  tornadoes starting  $\leq 6$  hours apart.

1972–2010: 79% of tornado fatalities & most economic losses occurred in outbreaks.

No trends in numbers of reliably reported tornadoes or outbreaks in last half century.

Mean & variance of the number of tornadoes per outbreak, & insured losses, increased significantly in last half century.

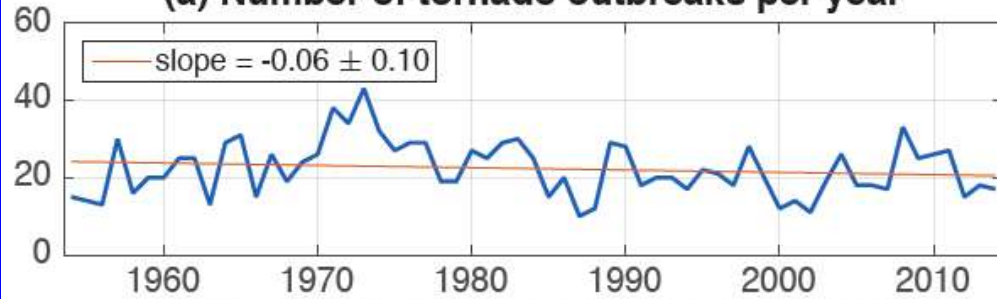
# F1+ tornadoes per outbreak in USA: $\text{variance} \sim (\text{mean})^{4.3}$



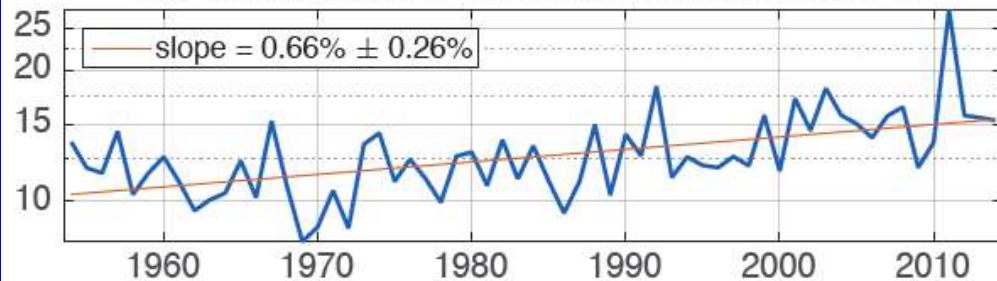
Tippett & Cohen, *Nature Communications* 2016

Michael Tippett

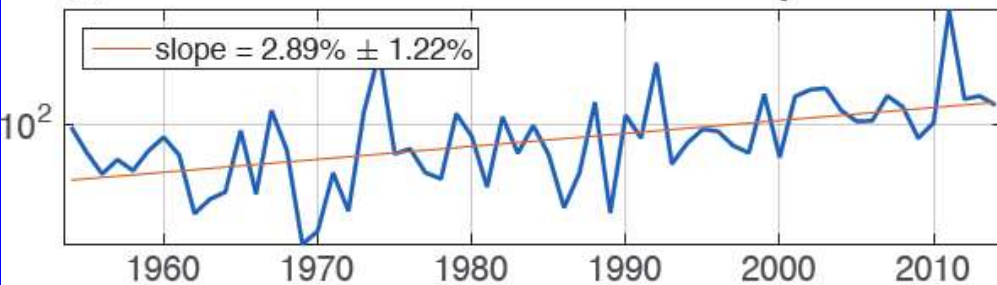
(a) Number of tornado outbreaks per year



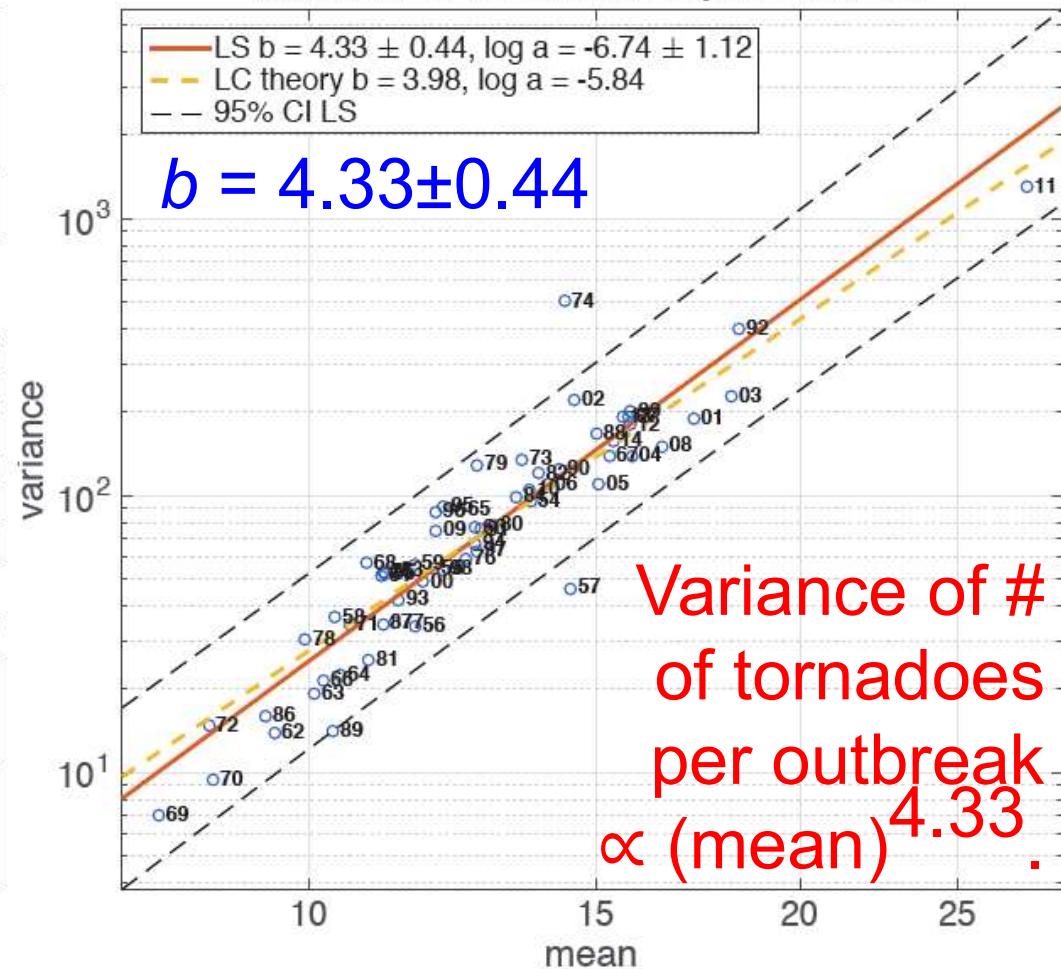
(b) Mean number of tornadoes per outbreak



(c) Variance of the number of tornadoes per outbreak



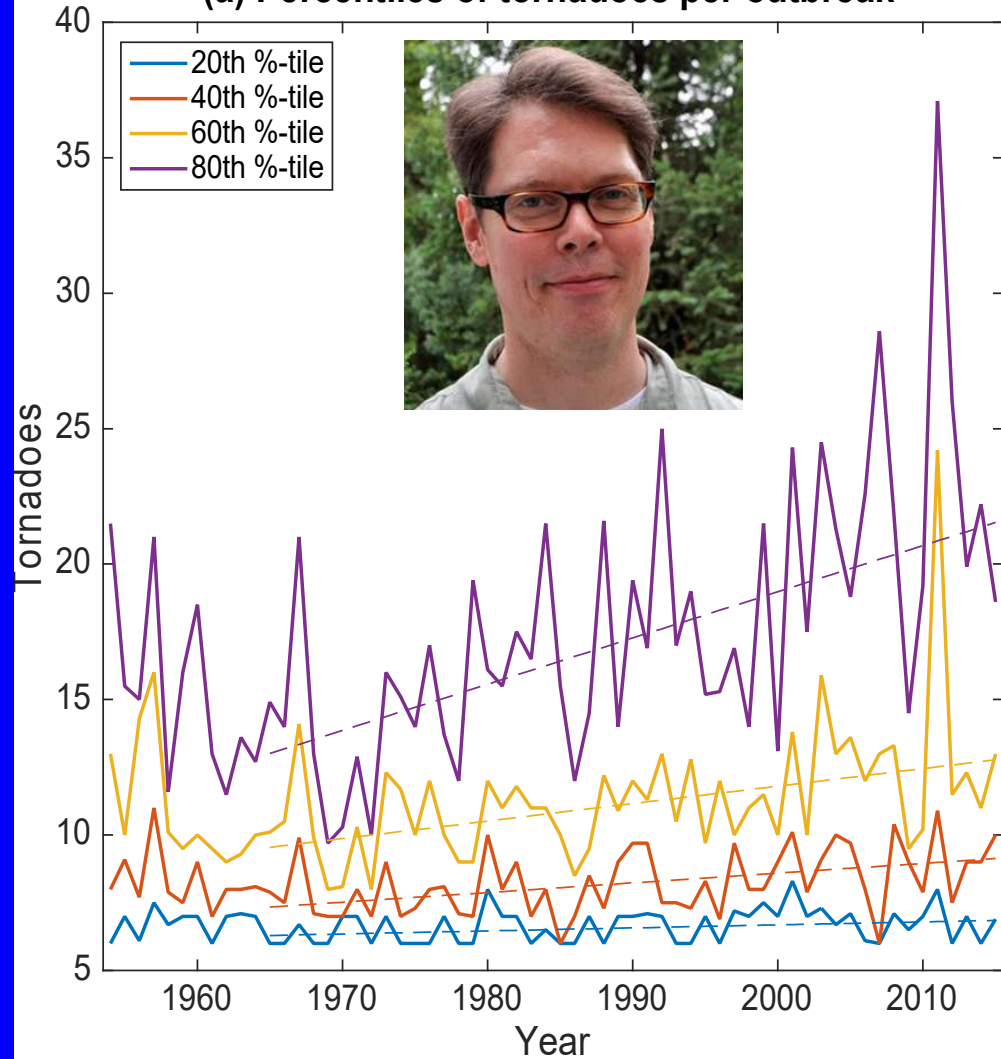
(d) Number of tornadoes per outbreak



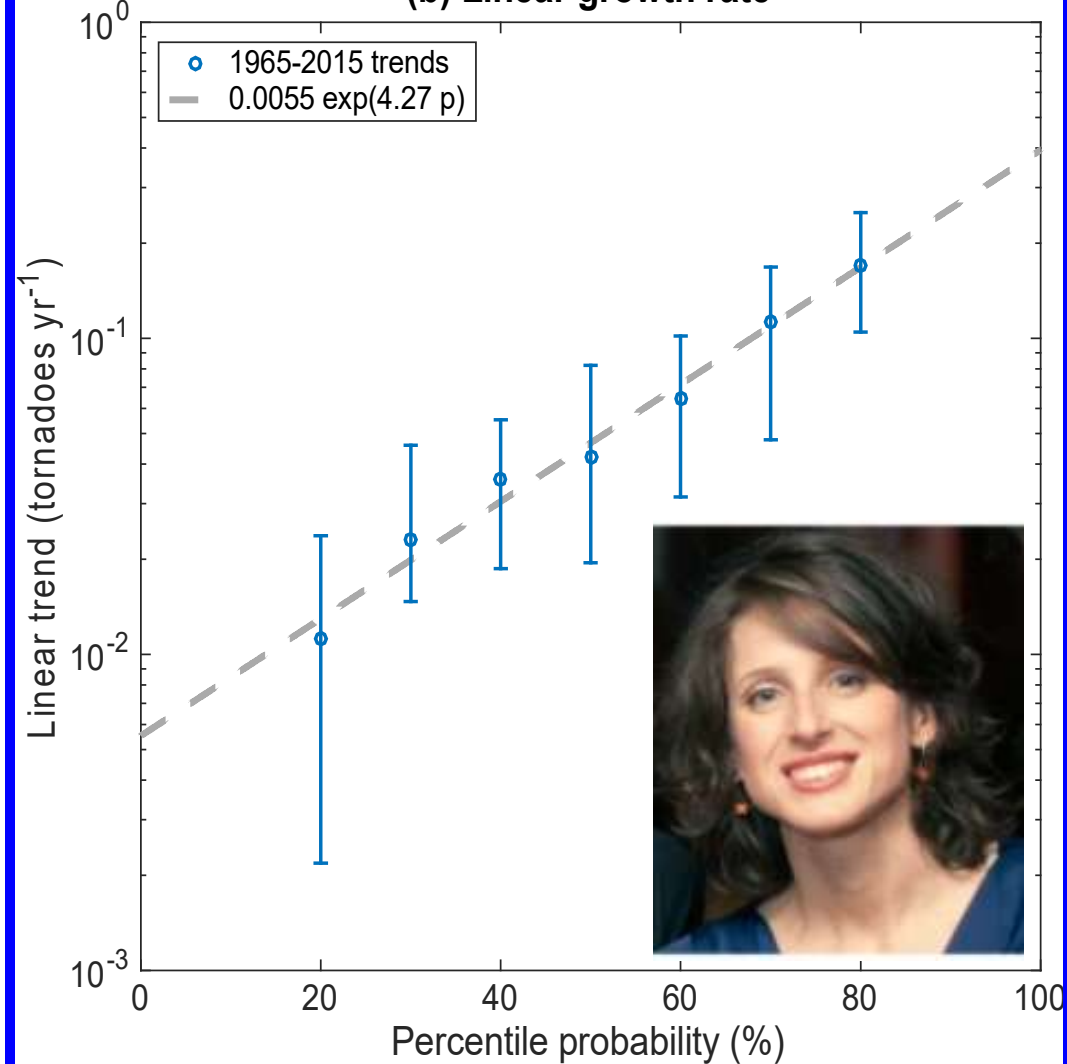
# Higher percentiles increased faster.

“quantile regression”

(a) Percentiles of tornadoes per outbreak



(b) Linear growth rate



# Plan

Fluctuation scaling, variance function

$$\textit{variance} = f(\textit{mean})$$

Taylor's law:  $\textit{variance} = a(\textit{mean})^b$

$$\log(\textit{variance}) = \log a + b \log(\textit{mean})$$

→ Heavy tails & regular variation

$$\Pr(X > x) = L(x)x^{-\alpha}, 0 < \alpha < 2$$

COVID-19 in US

Taylor's law, infinite variance

# Lognormal distribution

A positive-valued random variable  $Y(\mu, \sigma^2)$  with real parameters  $\mu, \sigma^2 \geq 0$  is

**lognormal** if

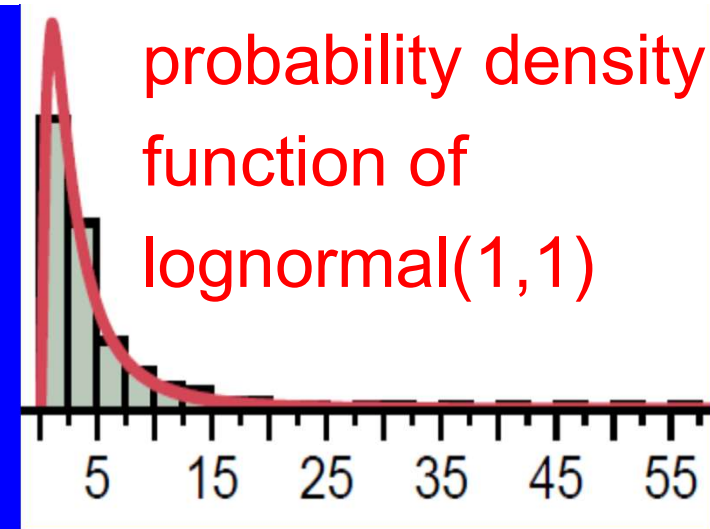
$\log Y(\mu, \sigma^2)$  is normal(mean  $\mu$ , variance  $\sigma^2$ ).

$$E\left(Y(\mu, \sigma^2)\right) = \exp\left(\mu + \frac{\sigma^2}{2}\right),$$

$$\text{Var}\left(Y(\mu, \sigma^2)\right) = [\exp(\sigma^2) - 1] \exp(2\mu + \sigma^2).$$

If  $\sigma^2$  is constant & only  $\mu$  changes, then

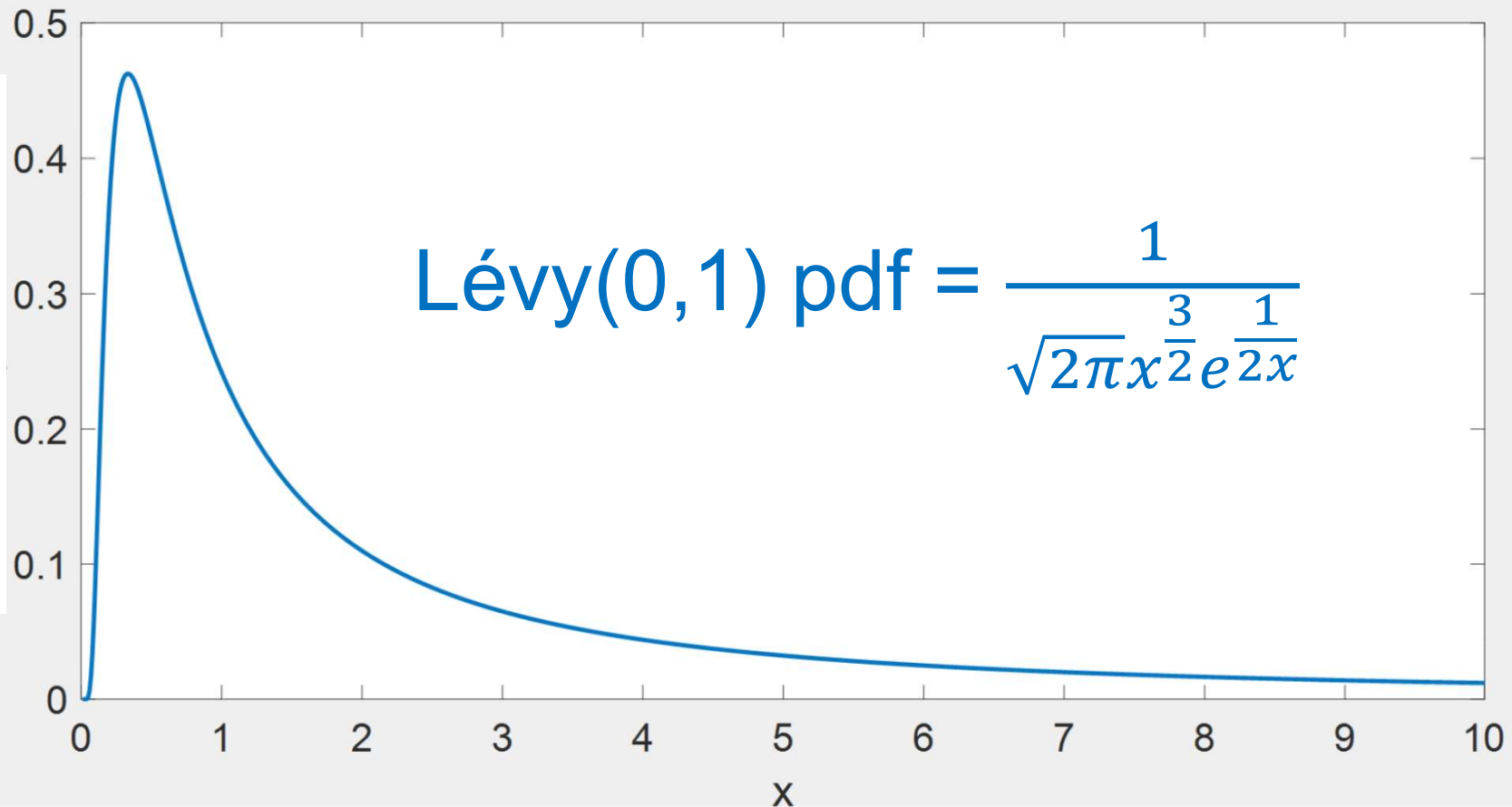
$$\text{Var}(Y(\mu)) = c\{E(Y(\mu))\}^2: \text{TL with exponent 2.}$$



# Lévy distribution $S_{1/2}$

If  $X$  is normal  $\mathcal{N}(0,1)$ , then  $1/X^2$  has Lévy distribution (1924)  $S_{1/2}$  [Helmert 1875, Lüroth 1876] with infinite mean & variance.

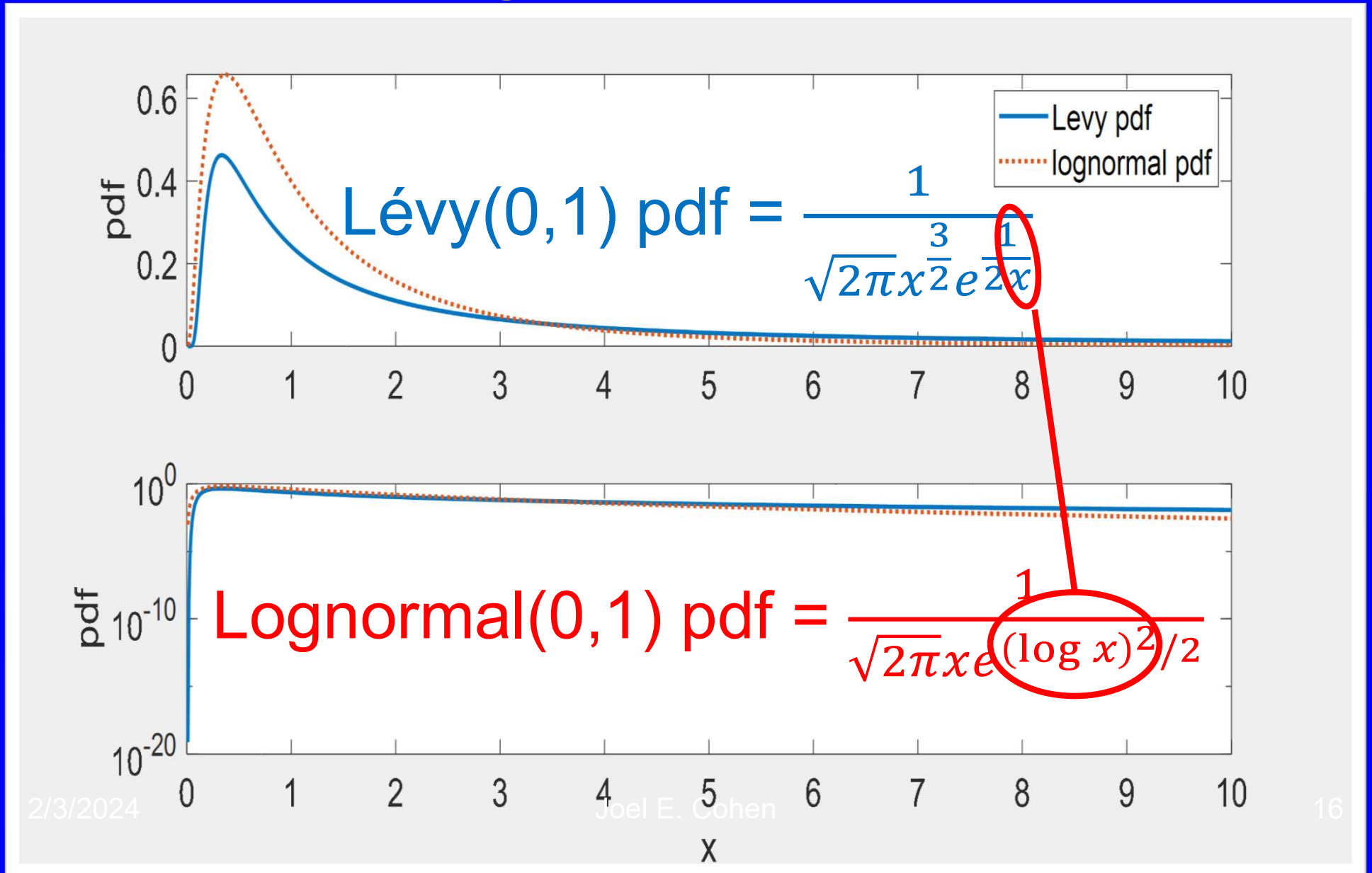
Lévy(0,1) pdf



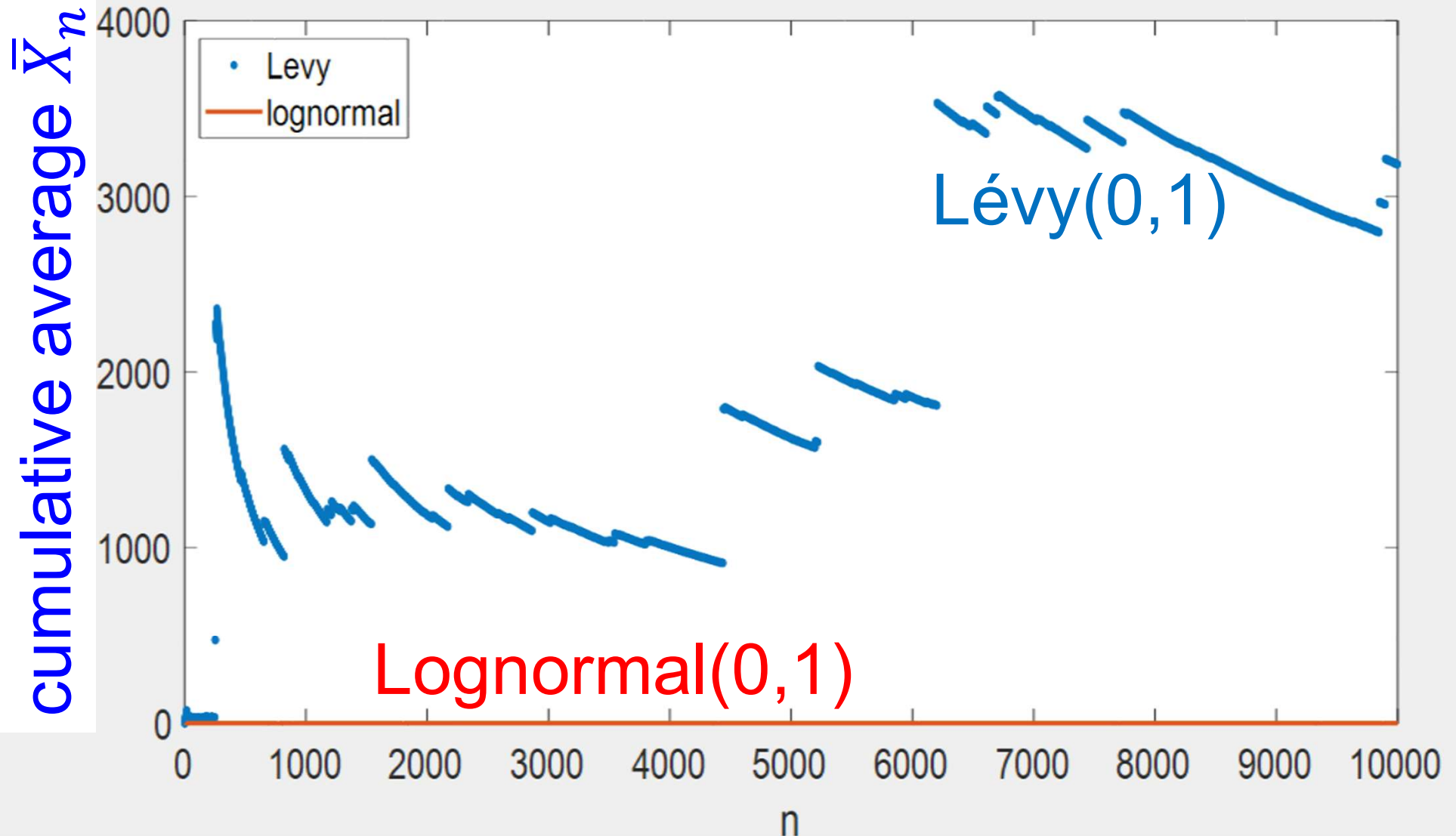
$$\text{Lévy}(0,1) \text{ pdf} = \frac{1}{\sqrt{2\pi}x^{3/2}e^{-\frac{1}{2x}}}$$



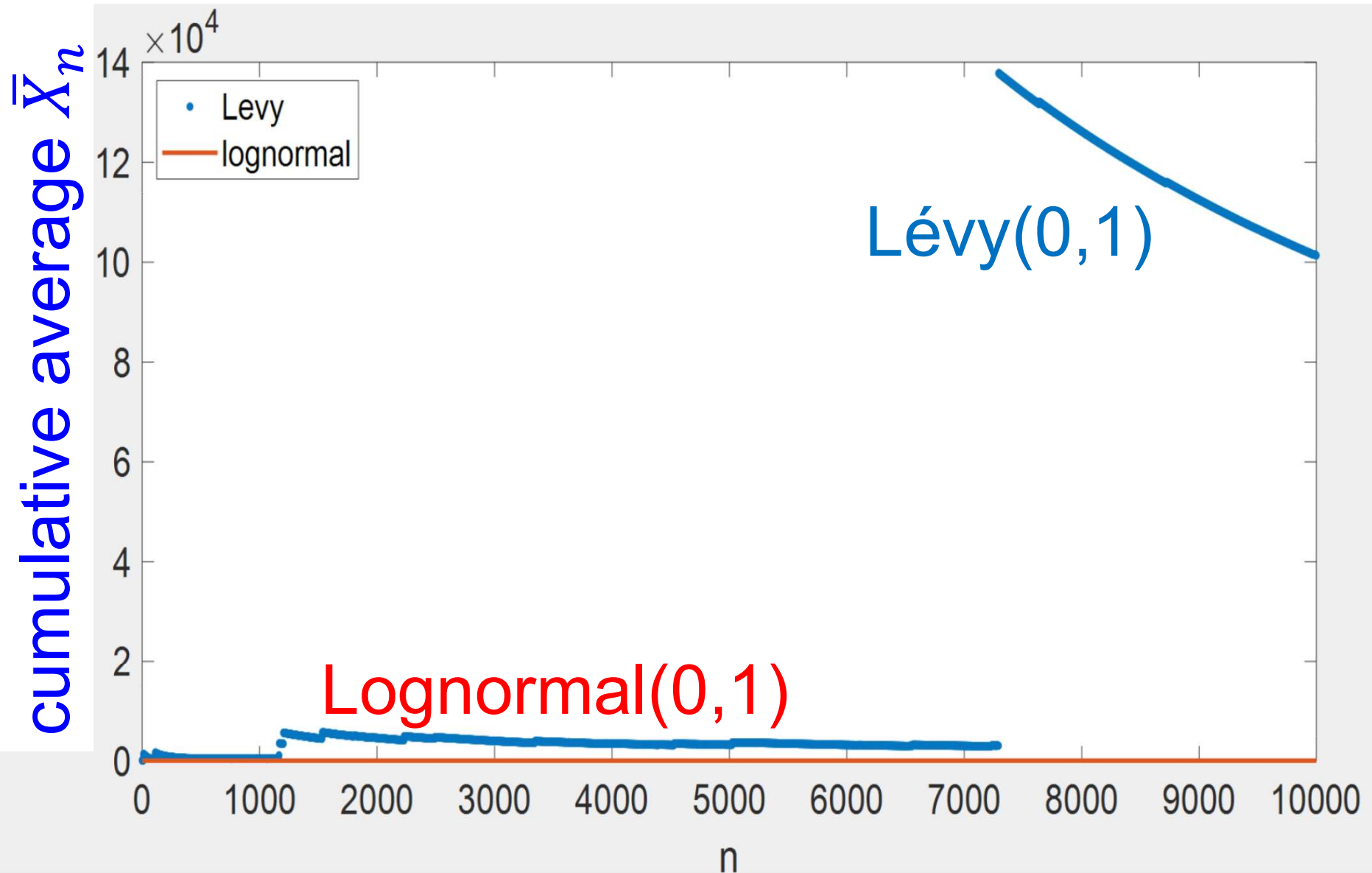
# Lévy distribution has heavier right tail than lognormal distribution.



Lévy cumulative averages grow like  $n \times \text{Lévy}$ . Lognormal averages converge.

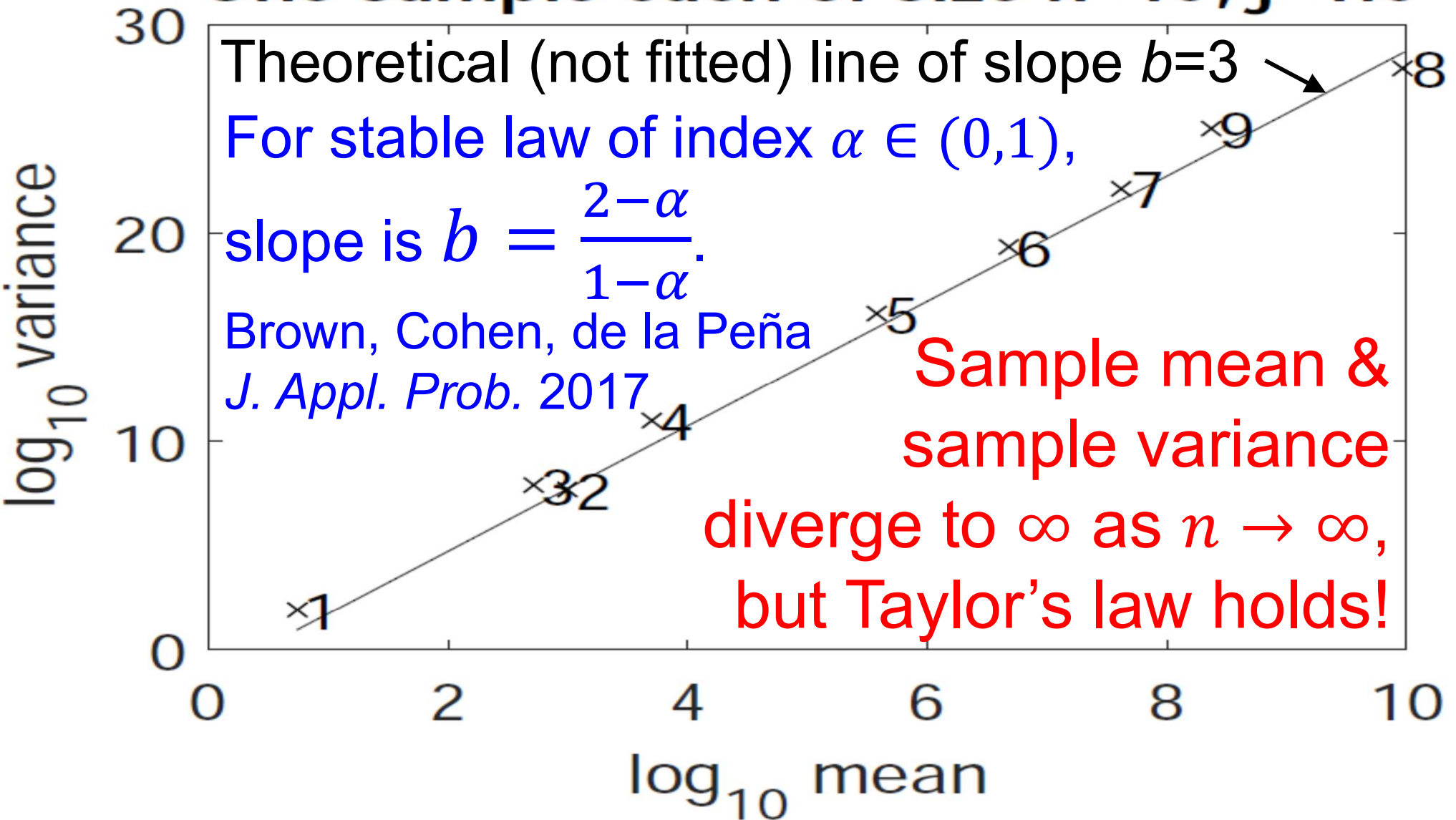


Lévy cumulative averages grow like  $n \times \text{Lévy}$ . Lognormal averages converge.



# Lévy law (stable $\alpha = 1/2$ ) obeys TL with increasing sample sizes.

One sample each of size  $n=10^j, j=1:9$



# Heterogeneous dependent data

Cohen, Davis, Samorodnitsky, *Proc. Roy. Soc. A* 2020

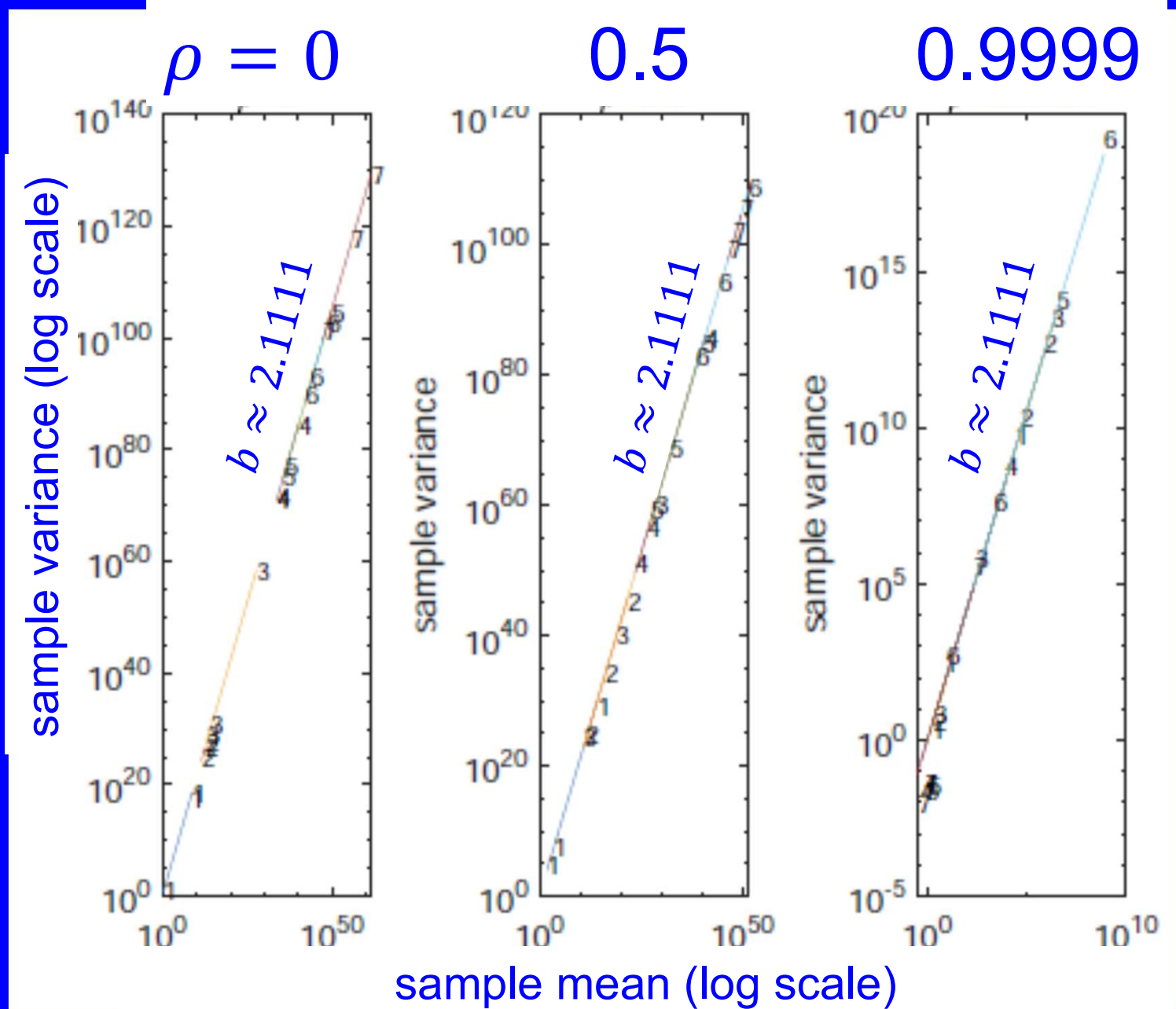
$$X_t := \frac{1}{|Z|^{\frac{1}{\alpha_t}}},$$

$$\Pr\{\alpha_t = 0.1\} = 0.1,$$

$$\Pr\{\alpha_t = 0.9\} = 0.9,$$

$$\text{corr}(Z_s, Z_t) = \rho, \text{ for } s \neq t.$$

$$b = \frac{2 - 0.1}{1 - 0.1} \approx 2.1111.$$

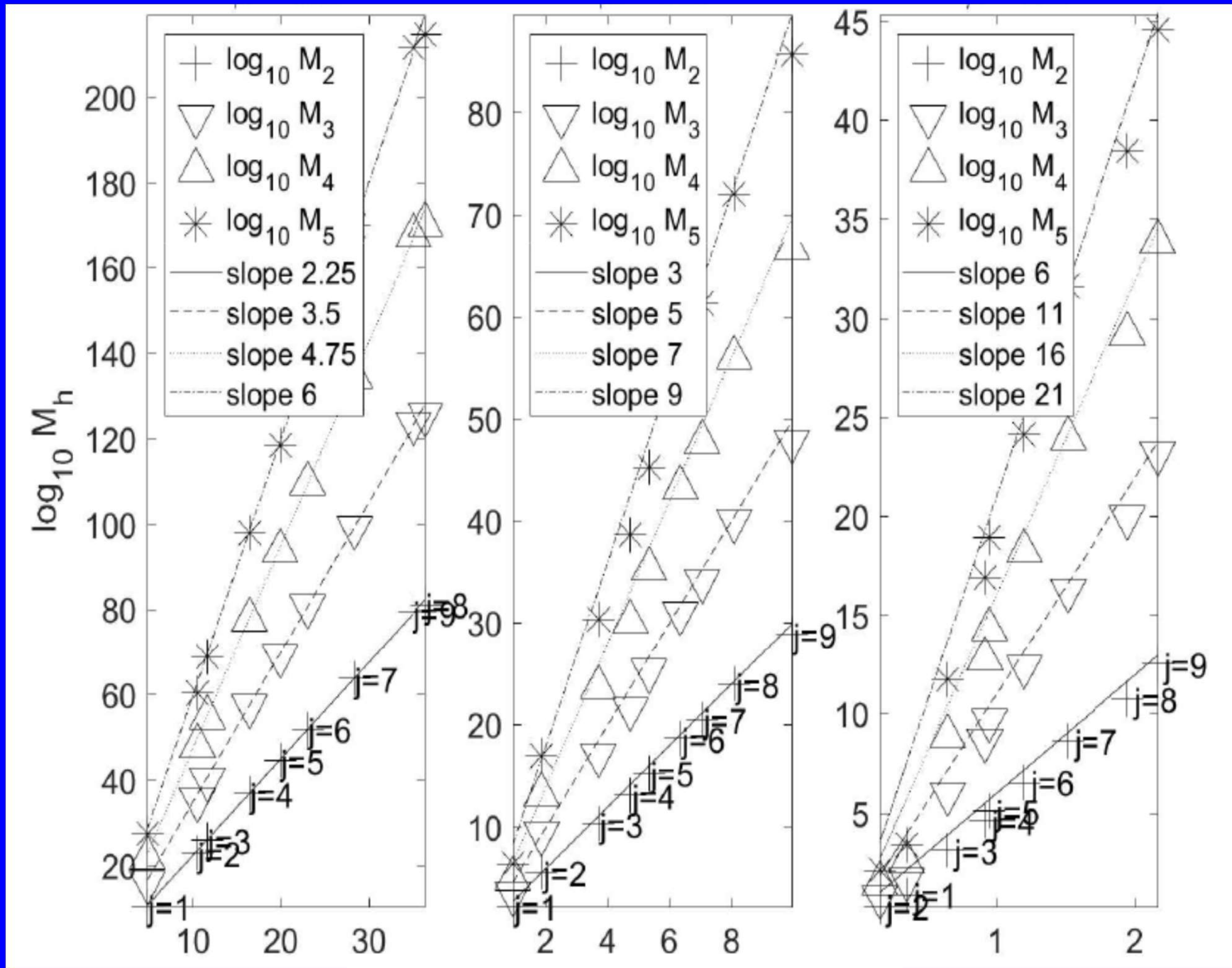


# Log central moments

$\alpha=0.2$

$\alpha=0.5$

$\alpha=0.8$



$$\log(\bar{X}_n), n = 10^j$$

# Many roads lead to TL.

Many models yield TL exactly or asymptotically.

Power-law form & parameter values of TL do not determine underlying mechanisms.

Interpreting the parameters of TL in terms of a specific mechanism requires testing the assumptions against detailed data.



# Plan

Fluctuation scaling, variance function

$$\textit{variance} = f(\textit{mean})$$

Taylor's law:  $\textit{variance} = a(\textit{mean})^b$

$$\log(\textit{variance}) = \log a + b \log(\textit{mean})$$

Heavy tails & regular variation

$$\Pr(X > x) = L(x)x^{-\alpha}, 0 < \alpha < 2$$

→ COVID-19 in US

Taylor's law, infinite variance

# COVID-19 cases & deaths

*New York Times* historical data base has final counts of COVID-19 cumulative cases & cumulative deaths at end of each day, 2020-01-21 to 2021-06-19 by "state" & "county" for days & counties with >0 cases or >0 deaths.

1,436,628 counts by day & county in data downloaded 2021-06-20

# On each date, cumulative cases & deaths within each state by county

State number →	$s=1$	$s=2$	$s=3$	$s=\dots$
County 1	$x_{11}$	$x_{12}$	$x_{13}$	$x_{\dots}$
County 2	$x_{21}$	$x_{22}$	$x_{23}$	$\dots$
County 3	$x_{31}$	$x_{32}$	$x_{33}$	$\dots$
County 4		$x_{42}$	$x_{43}$	$\dots$
$\vdots$		$x_{52}$		$\dots$
Mean = average	$m_1$	$m_2$	$m_3$	$m_{\dots}$
Variance	$v_1$	$v_2$	$v_3$	$v_{\dots}$

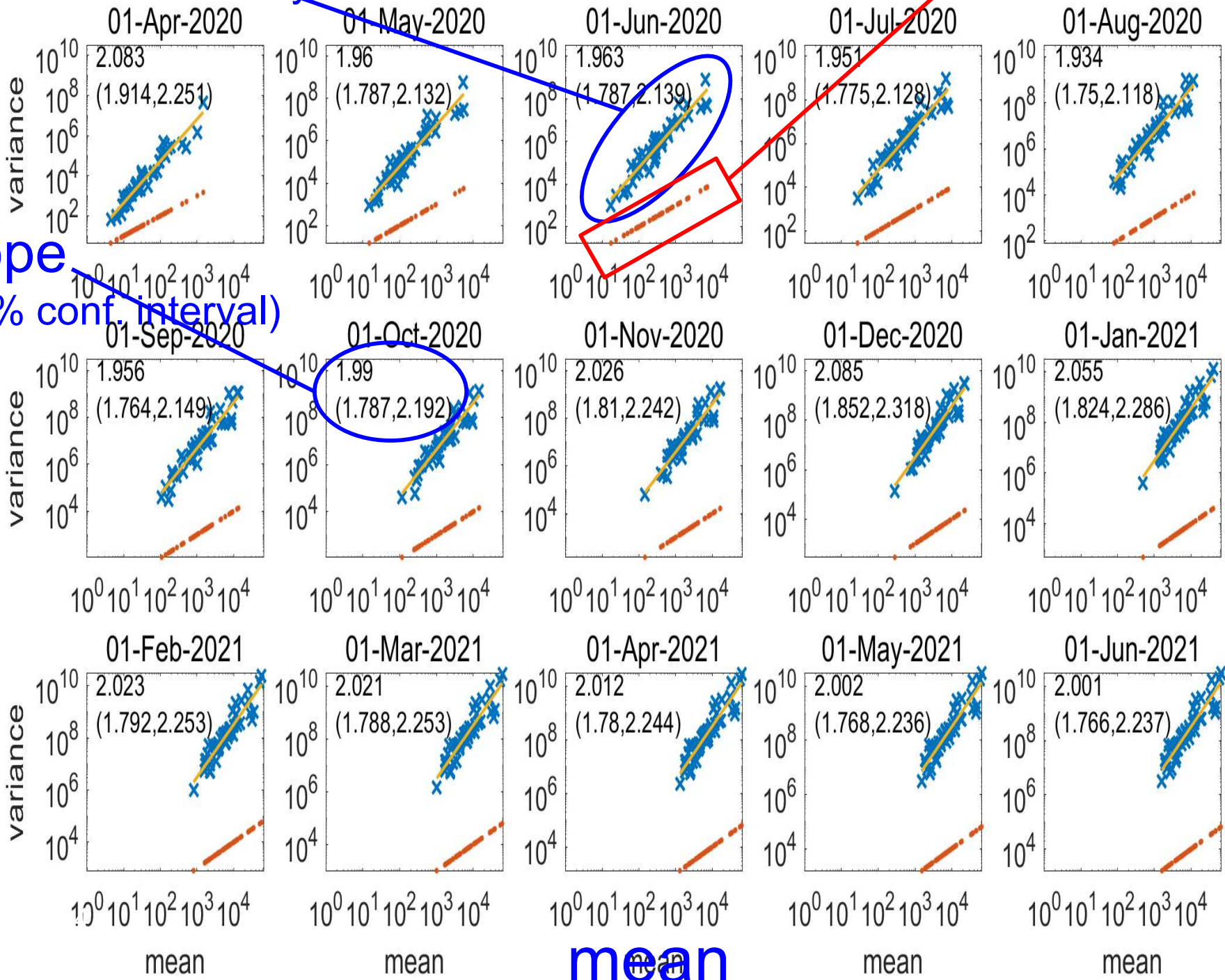
cases

Variance function of cumulative U.S. COVID-19 cases/county by state

only states with  $\geq 7$  counties with  $>0$  counts

Poisson

slope  
(95% conf. interval)  
variance

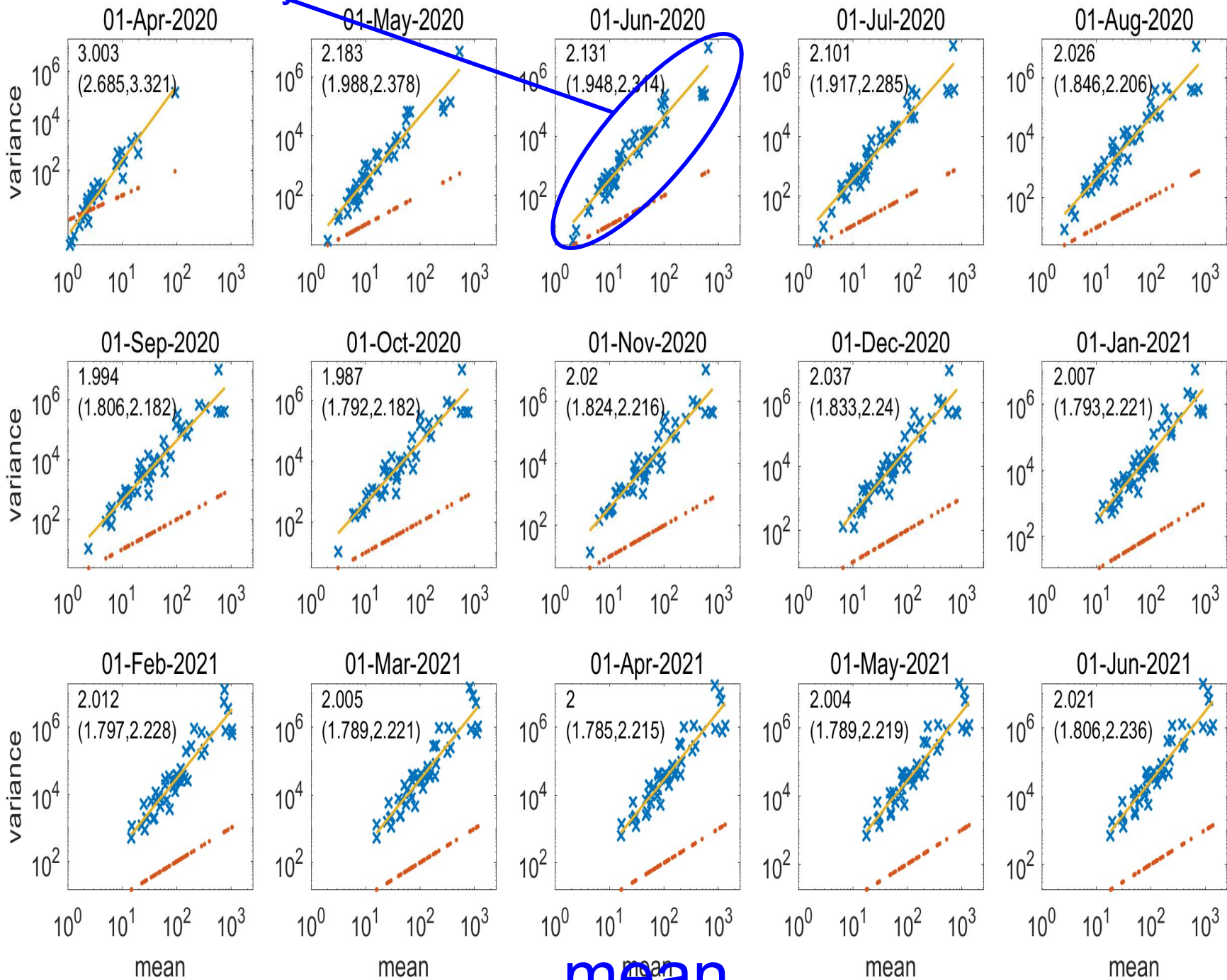


# deaths

## only states with $\geq 7$ counties with $>0$ counts

Variance function of cumulative U.S. COVID-19 deaths/county by state

variance

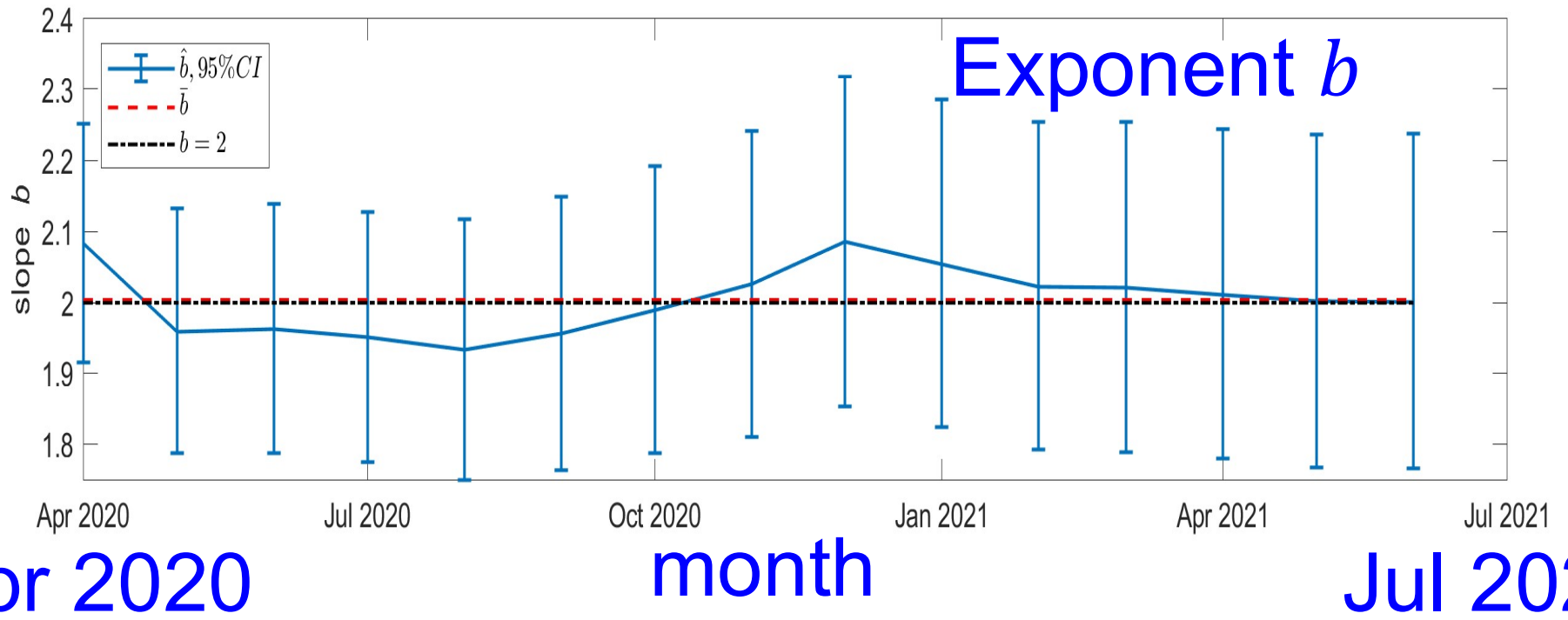
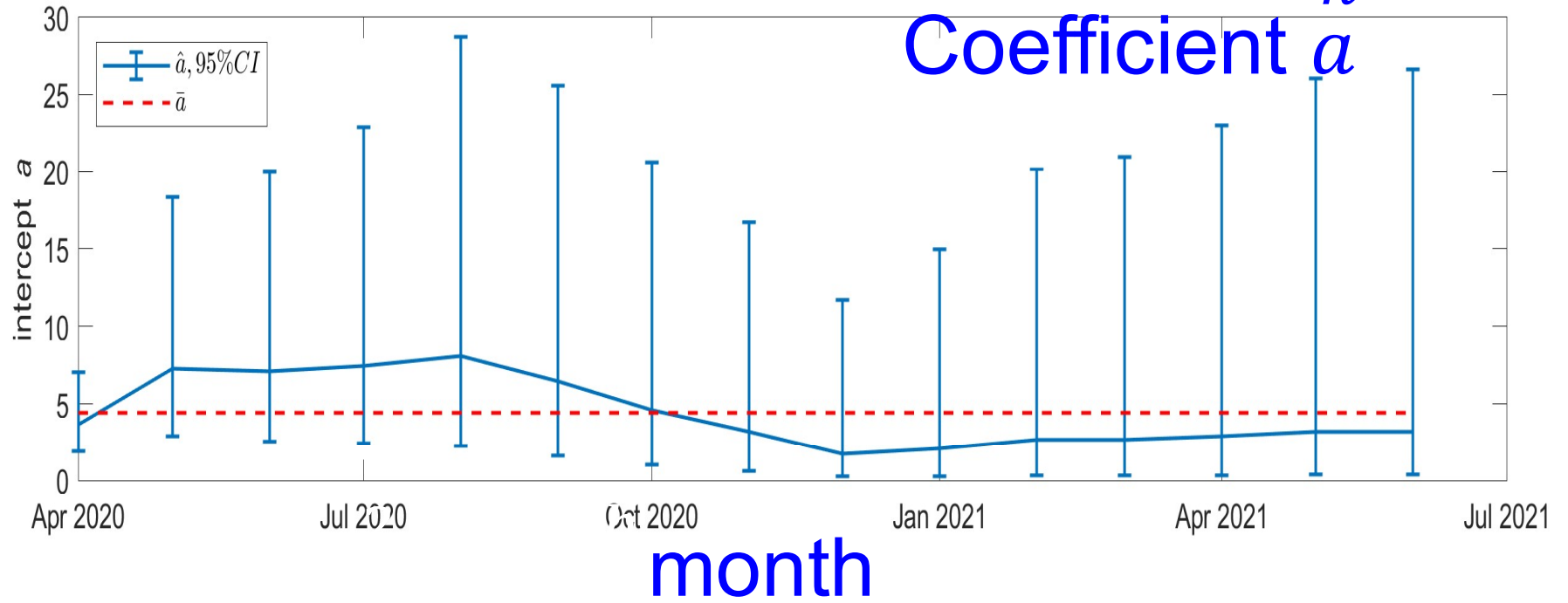


mean



cases

Taylor's law parameters of cumulative U.S. COVID-19 cases/county by state  $s_n^2 = a\bar{X}_n^b$

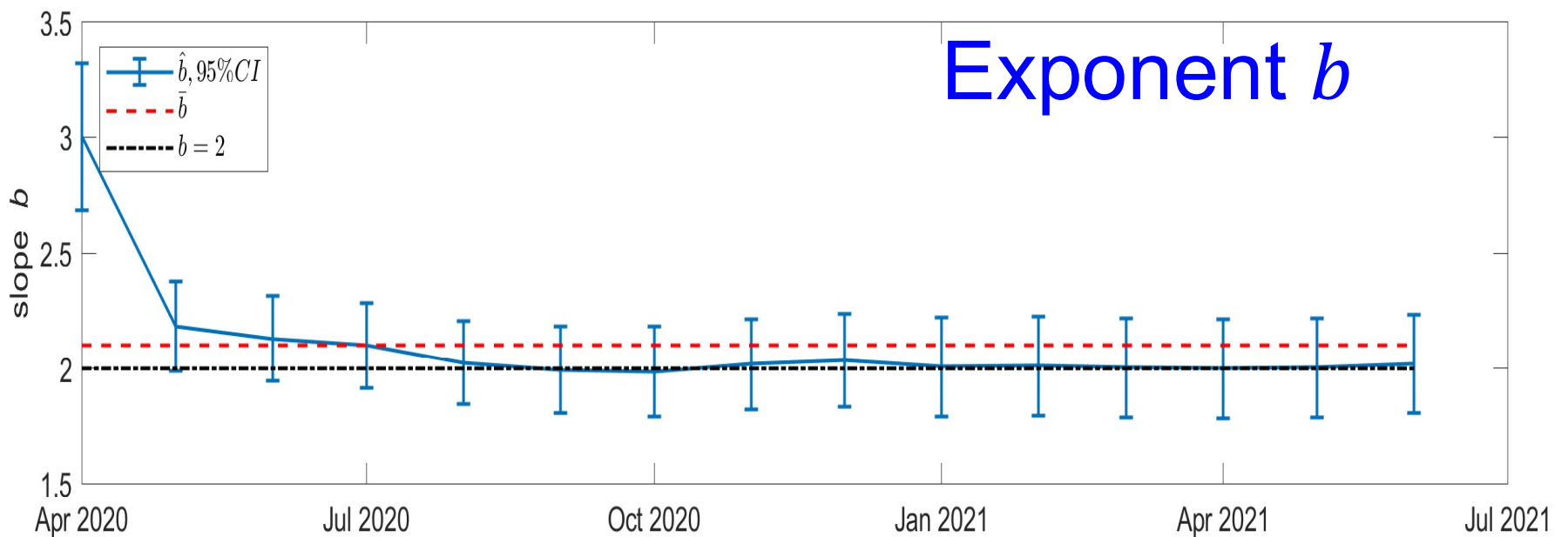
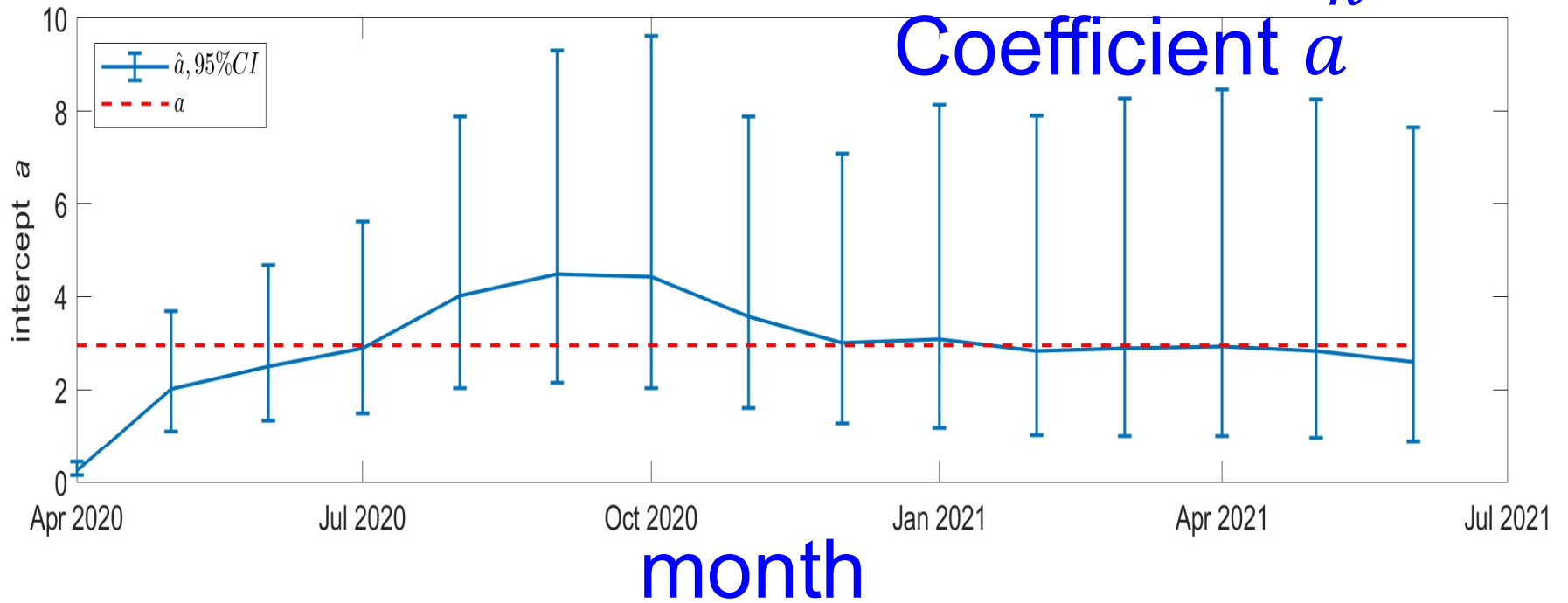


Apr 2020

Jul 2021

# deaths

Taylor's law parameters of cumulative U.S. COVID-19 deaths/county by state  $s_n^2 = a\bar{X}_n^b$



Apr 2020

month

Jul 2021



# Taylor's law describes counties' cumulative cases & deaths.

From April 2020 onward, TL holds:

1. log variance of counts (over counties) increases linearly with the log mean of counts (over counties) from state to state.

2. Slope  $b \approx 2$ .

## Why?

# Survival curve plots probability that counts $> x$ as a function of $x$ .

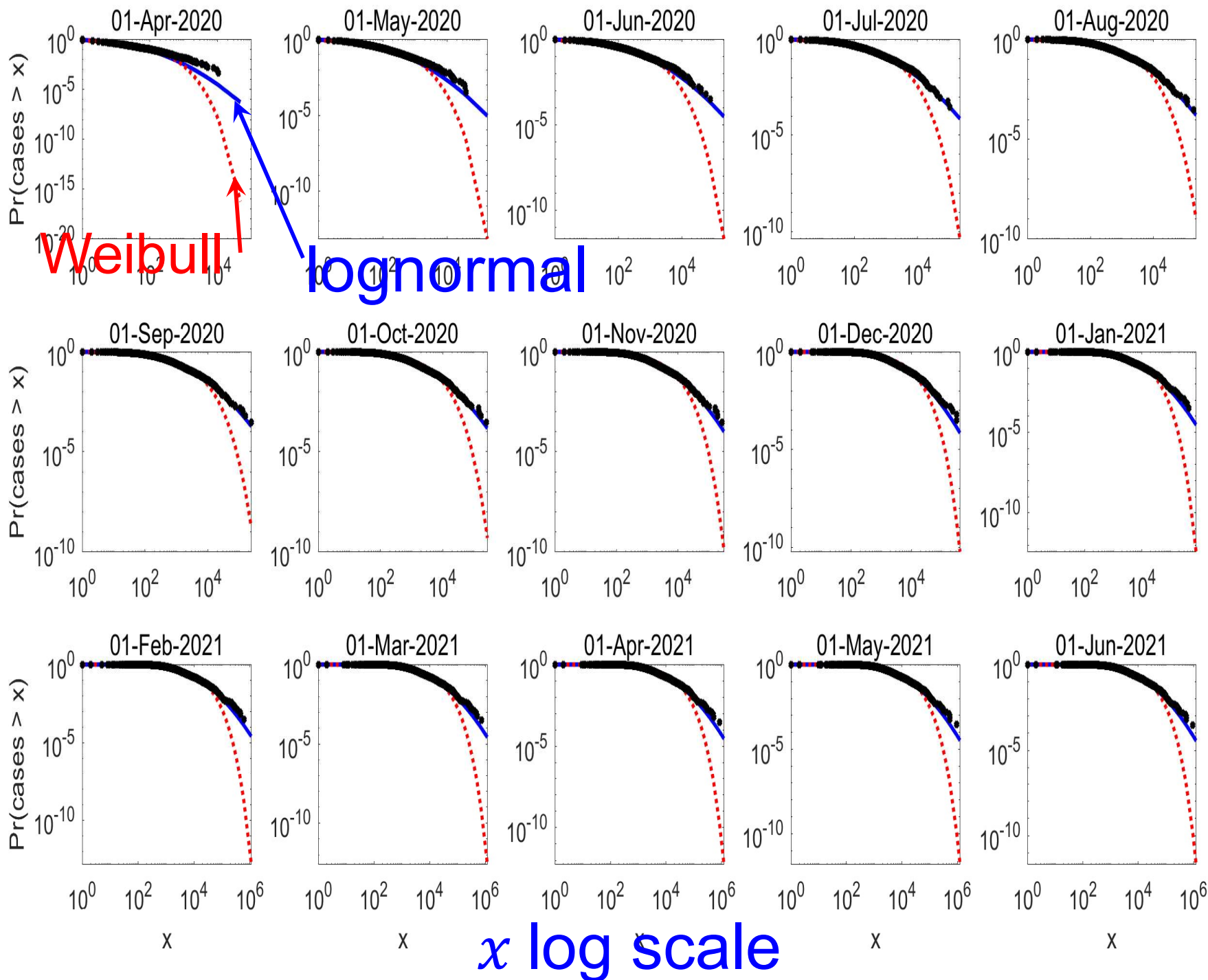
To cover wide ranges of probability  
& of counts, we plot  $\log(\Pr\{X > x\})$   
as a function of  $\log(x)$ .

We also fit lognormal & Weibull  
distributions by maximum  
likelihood to counts of all counties.

cases

Survival curve of cumulative COVID-19 cases/county by date

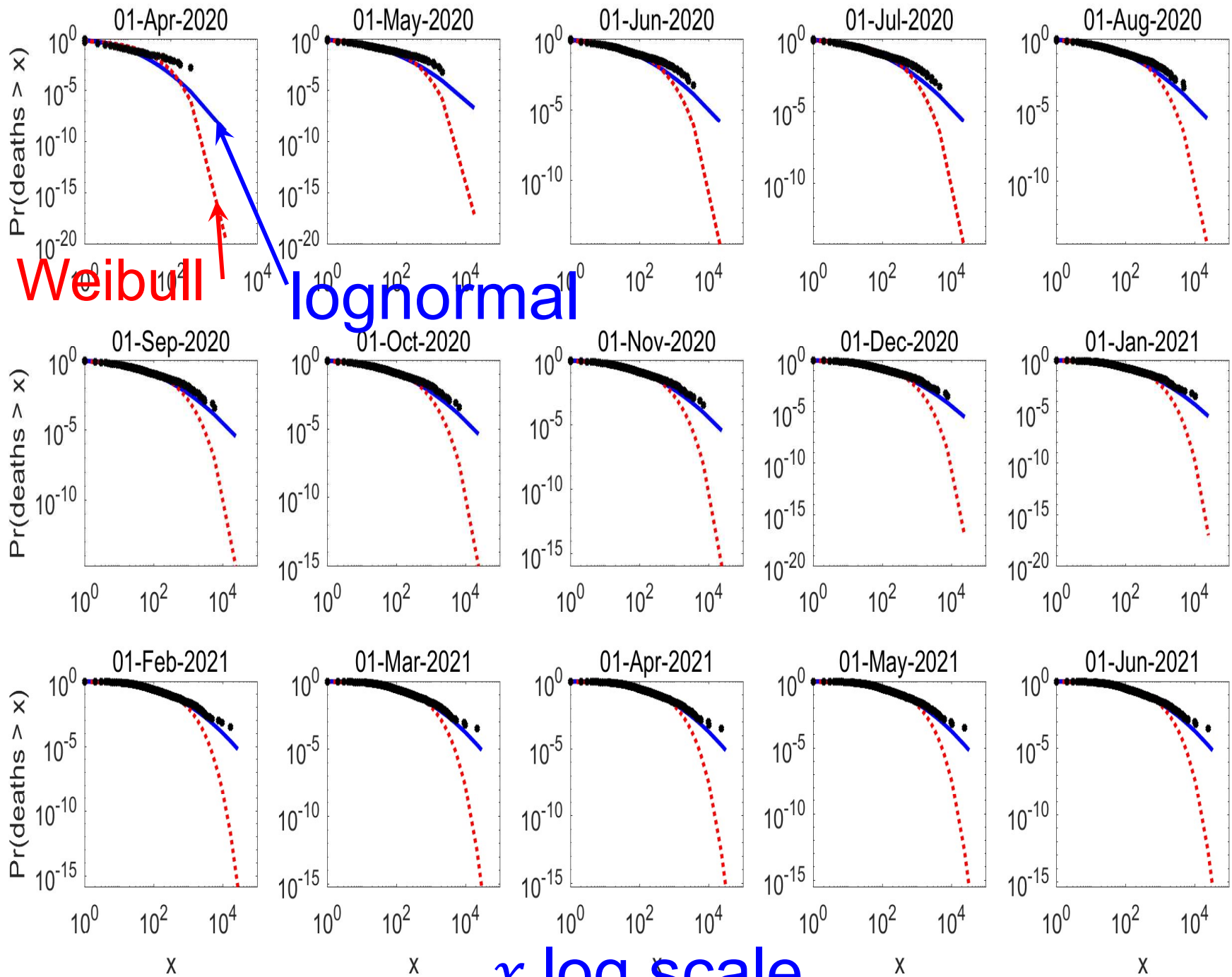
$\Pr(X > x)$  log scale



# deaths

Survival curve of cumulative COVID-19 deaths/county by date

$\Pr(X > x)$  log scale



Weibull

lognormal

$x$  log scale

Lognormal describes >99% of distributions of cases & deaths.

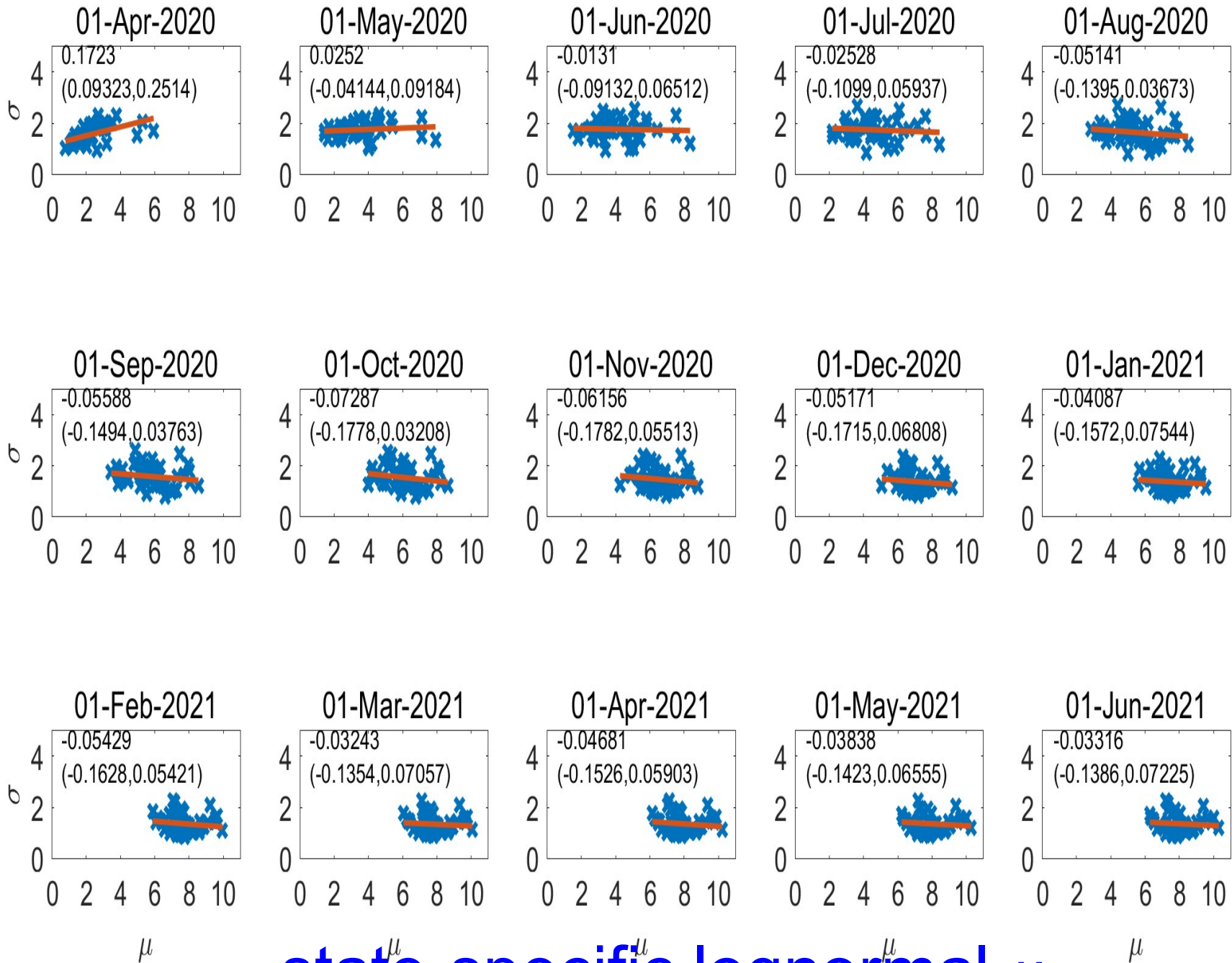
Many models lead to lognormal & Weibull distributions.

Lognormal is much closer than Weibull to data's upper tail, but also falls too fast.

If count has lognormal distribution, if  $\sigma^2$  is constant, & if only  $\mu$  varies from state to state, then Taylor's law holds with slope 2.

# cases Lognormal $\sigma$ as function of $\mu$ for cumulative U.S. COVID-19 cases by state

state-specific lognormal  $\sigma$

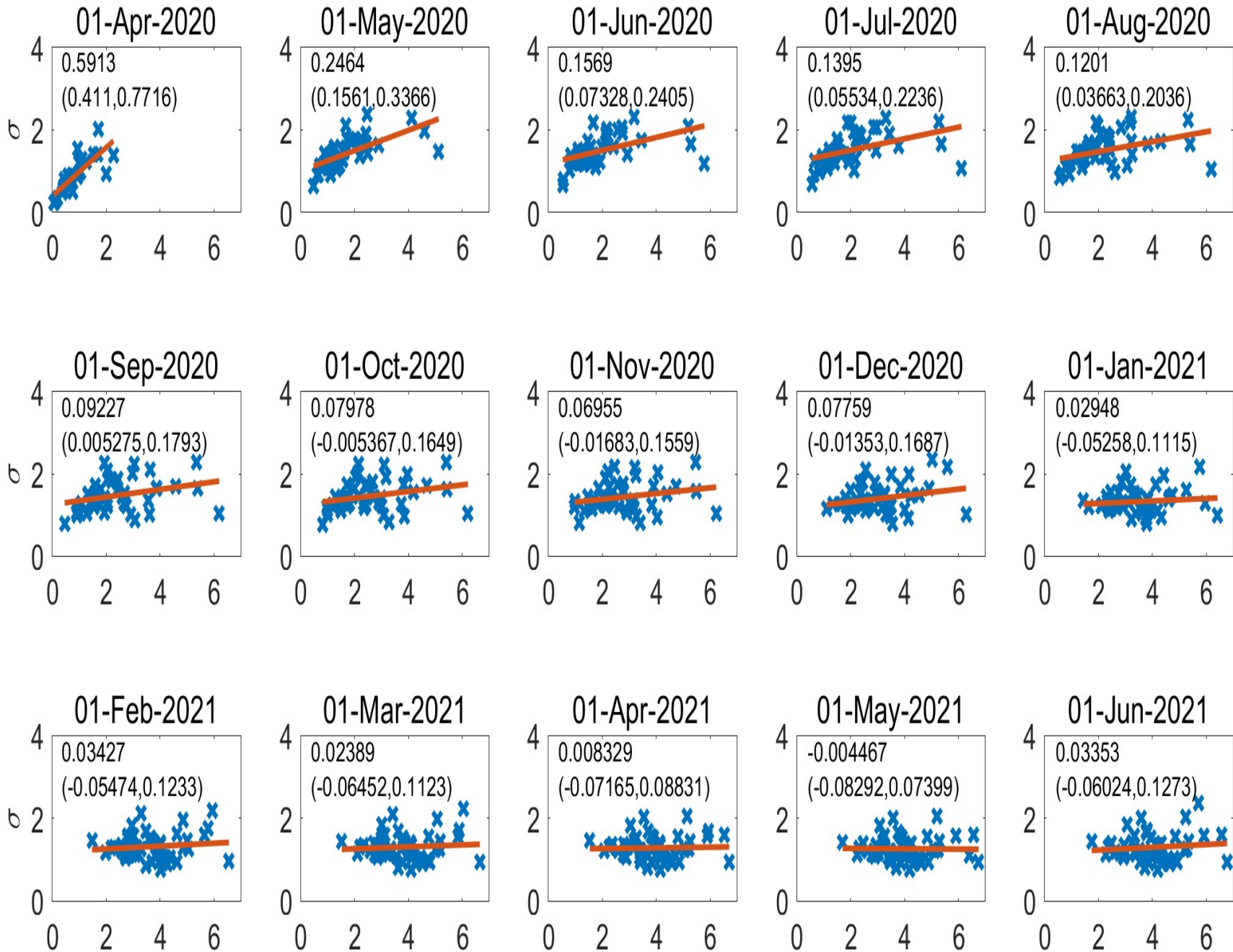


state-specific lognormal  $\mu$



# deaths Lognormal $\sigma$ as function of $\mu$ for cumulative U.S. COVID-19 deaths by state

state-specific lognormal  $\sigma$



state-specific lognormal  $\mu$



Lognormal explains TL with slope 2 for lower 99% of counts but not 1% upper tail.

Lognormal  $\mu$  varies much more widely than lognormal  $\sigma^2$ , generating predicted means & predicted variances that closely approximate Taylor's law with slope 2.

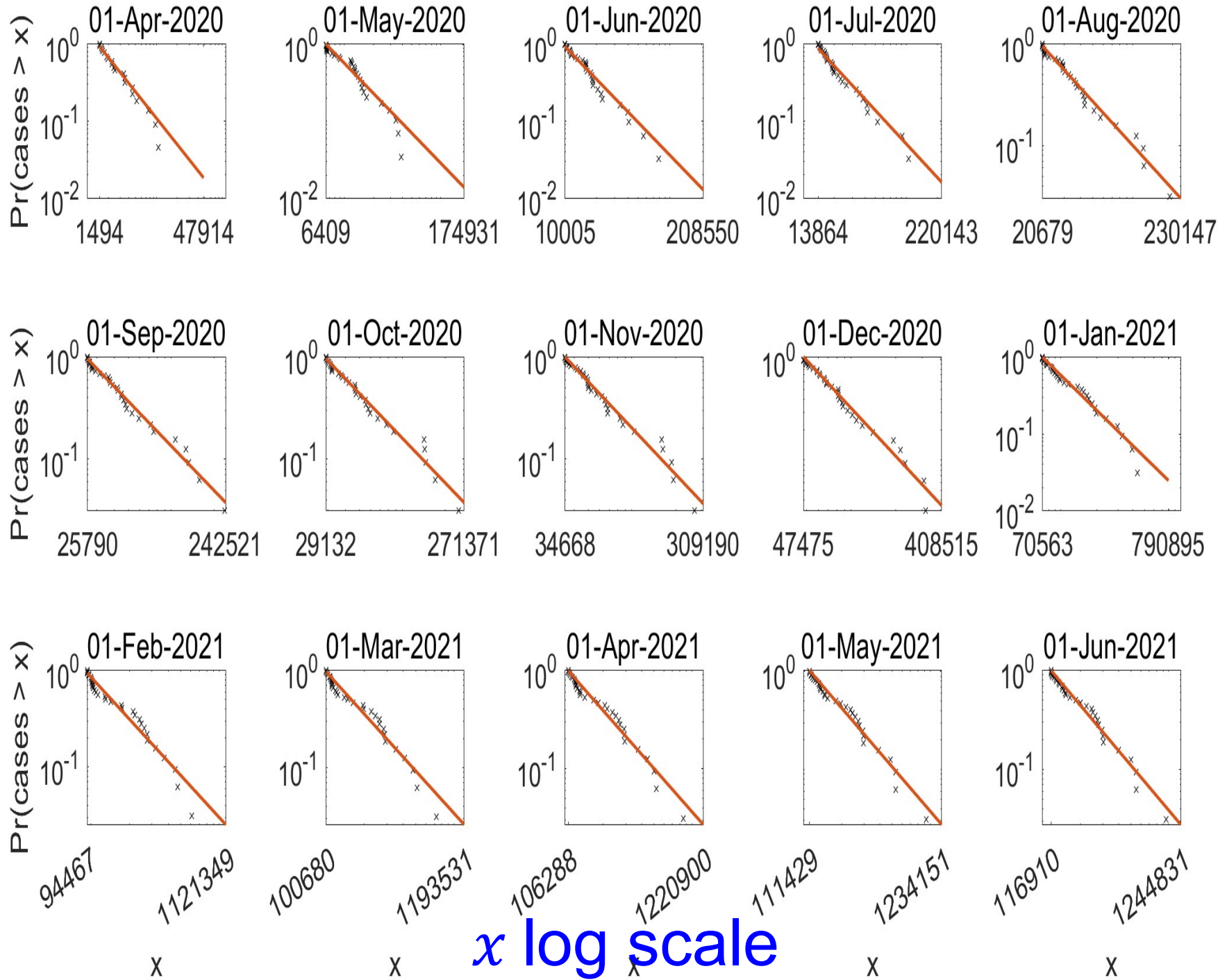
But the very largest counts of cases & deaths are more extreme than lognormal distribution predicts.

We zoom in to the counties with the highest 1% of counts of cases or deaths.

# cases

Survival curve of highest 1% of cumulative COVID-19 cases/county by date

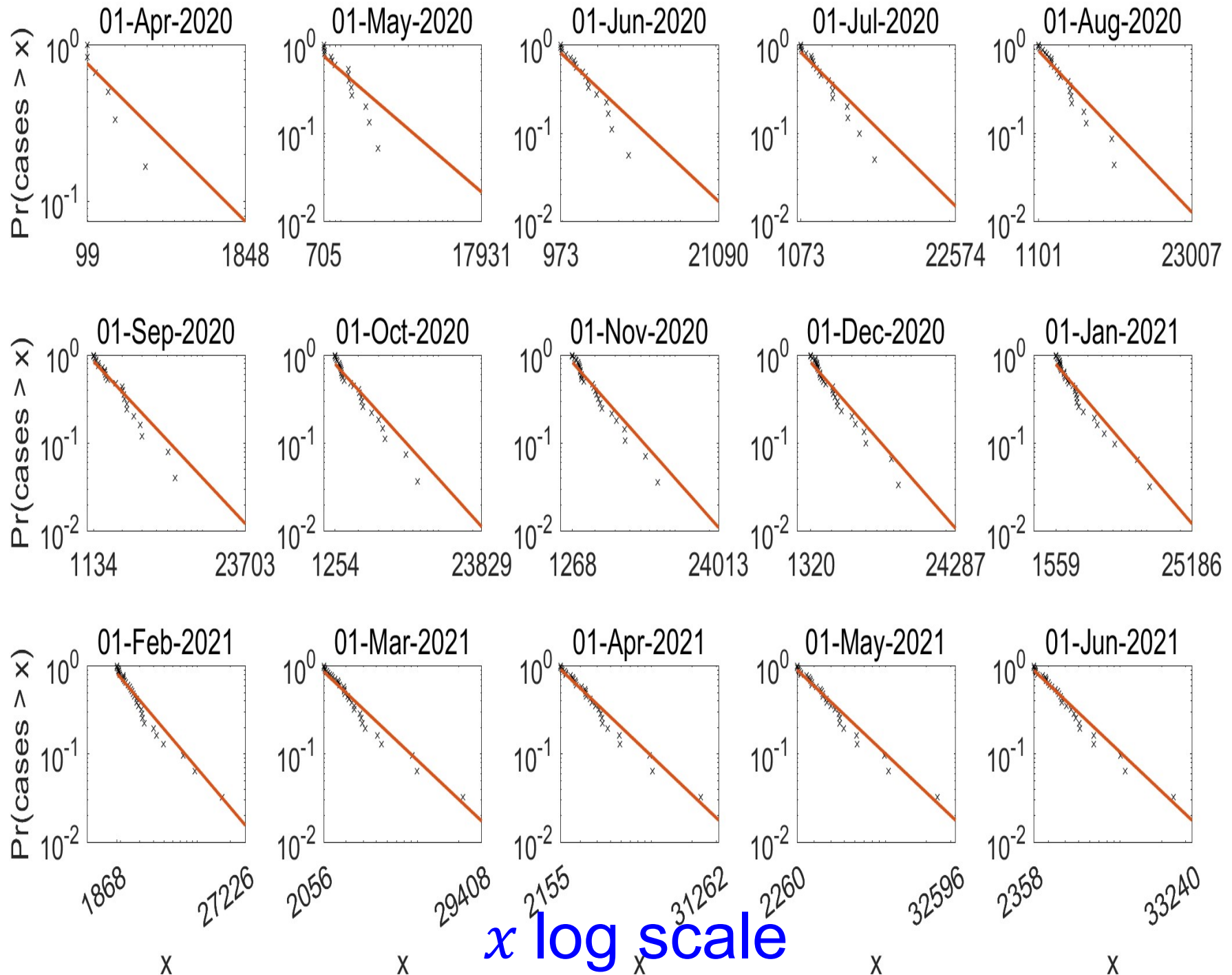
$\Pr(X > x)$  log scale



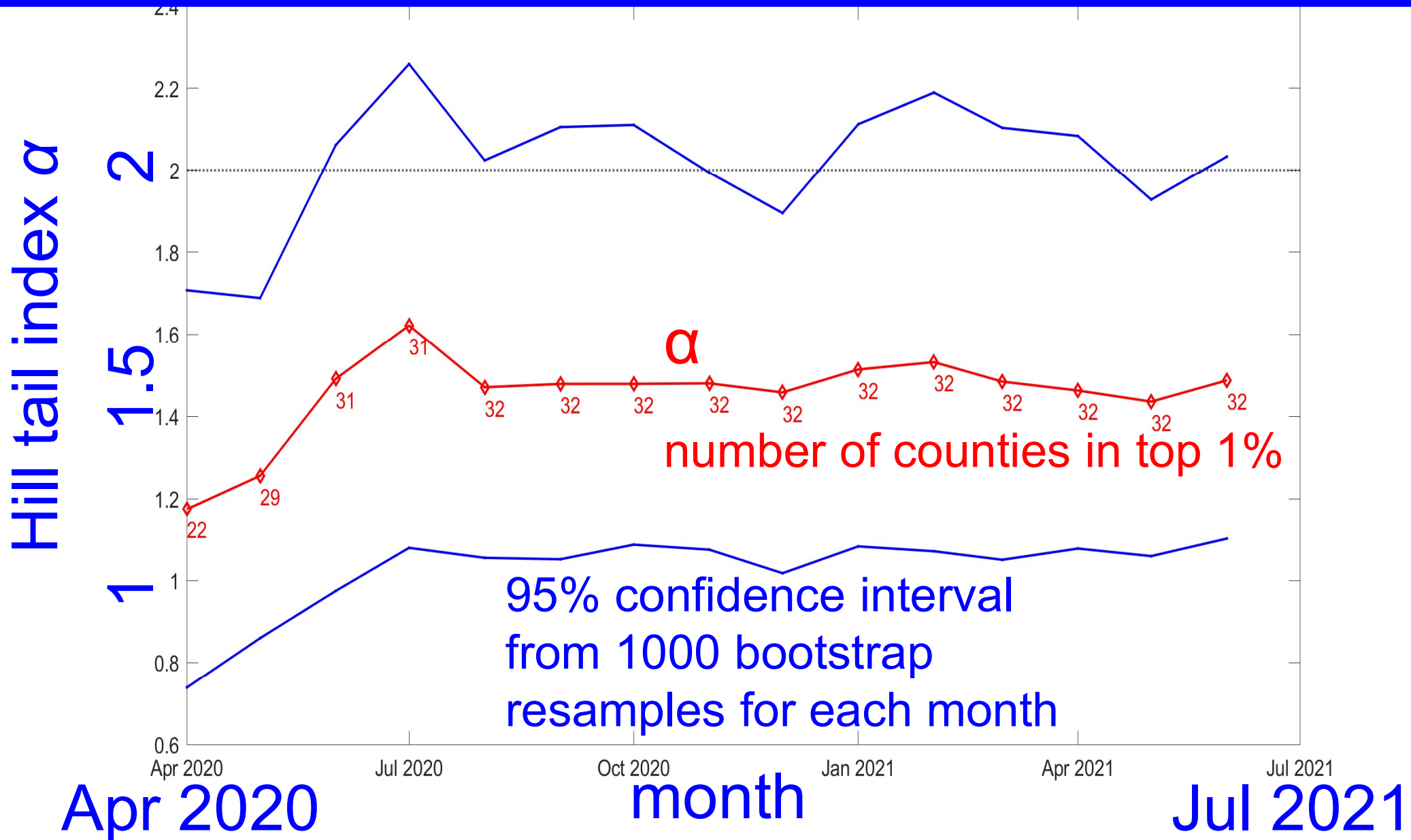
# deaths

Survival curve of highest 1% of cumulative COVID-19 deaths/county by date

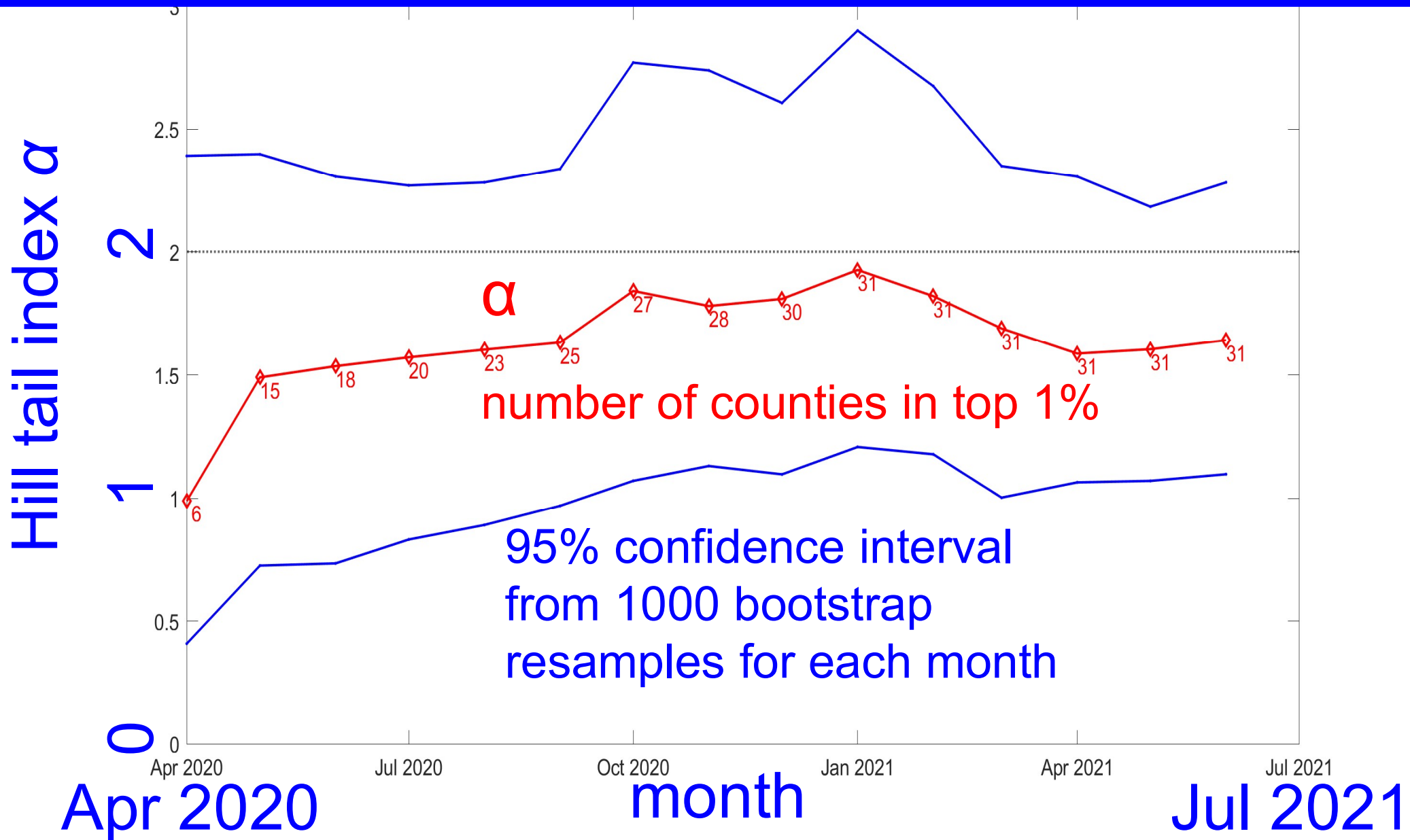
$\Pr(X > x)$  log scale



# For counties with highest 1% of cases, Hill estimates of tail index: $1 < \alpha < 2$



# For counties with highest 1% of deaths, Hill estimates of tail index: $1 < \alpha < 2$





# Empirical survival curves suggest variance is infinite.

The estimated upper tail index is  $1 < \alpha < 2$  in all 15 months for cases & all but first month April 2020 for deaths, so variance is infinite, mean is finite.



Bell 1806 / Wikipedia

"Wonder / Fear / Astonishment"



Regularly varying upper tail with index  $\alpha \in (1,2)$  explains why TL with  $b = 2$  holds even for largest counts where lognormal distribution fails.

Simulations: 100 samples of size 100

4 models in  $RV(\alpha)$ :  $|N(0,1)|^{-\alpha^{-1}}$ ,  $|U|^{-\alpha^{-1}}$ ,  
 $|U_1 U_2|^{-\alpha^{-1}}$ ,  $|U_1 U_2 U_3|^{-\alpha^{-1}}$

$\alpha = 1/2, 1, 3/2$ , with & without dependence.

TL with  $b \rightarrow 2$  is confirmed by mathematics.

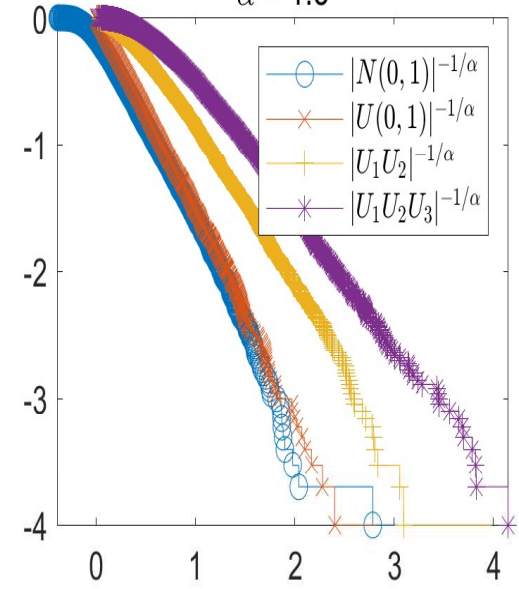
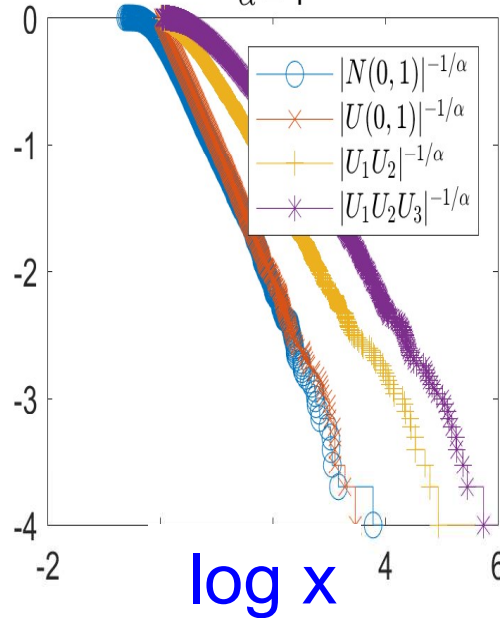
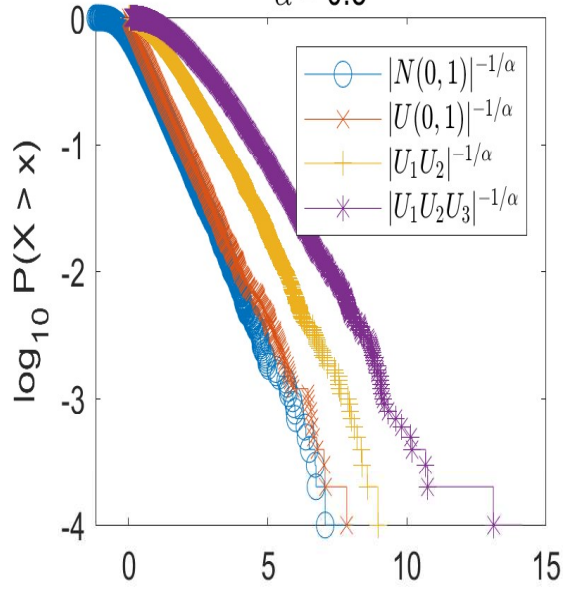
# 4 models in $RV(\alpha)$ , 100 samples of size 100

$\alpha = 1/2$

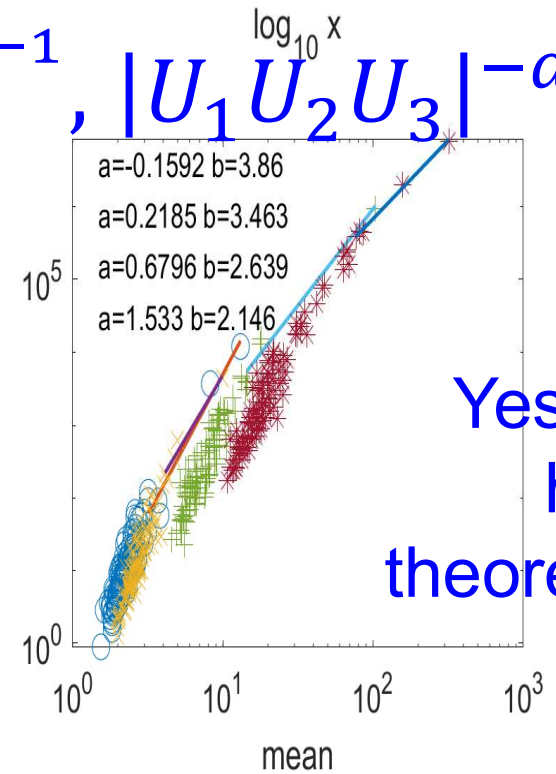
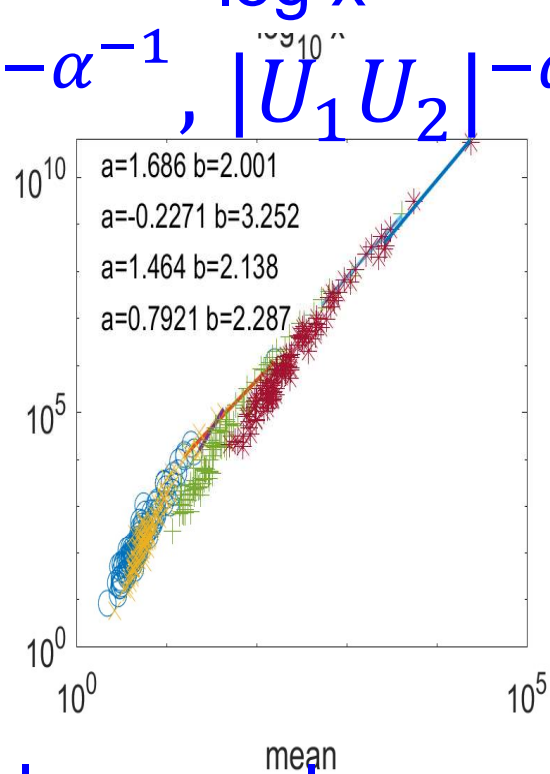
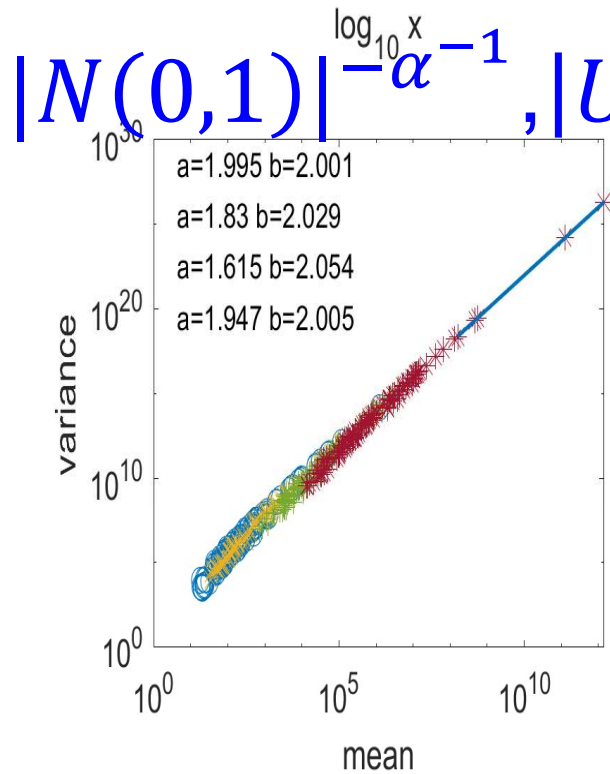
$\alpha = 1$

$\alpha = 1.5$

$\log \Pr(X > x)$



$\log$  sample variance

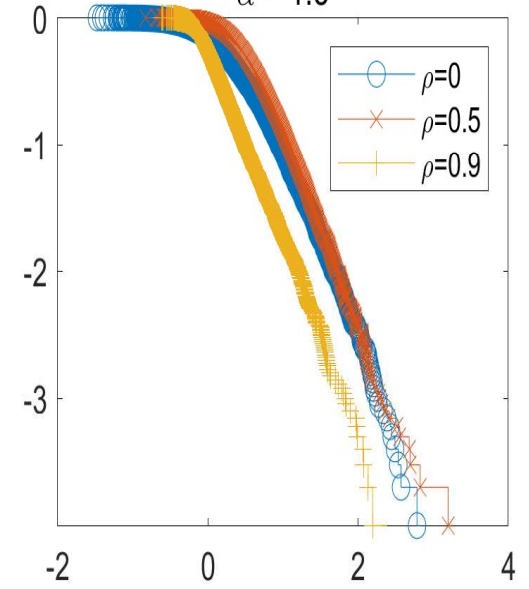
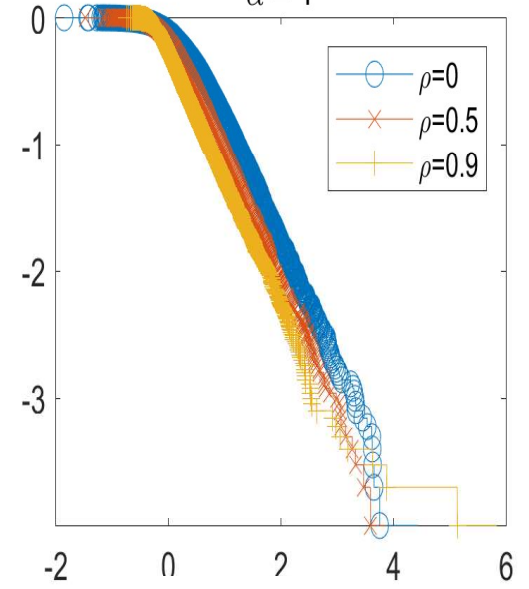
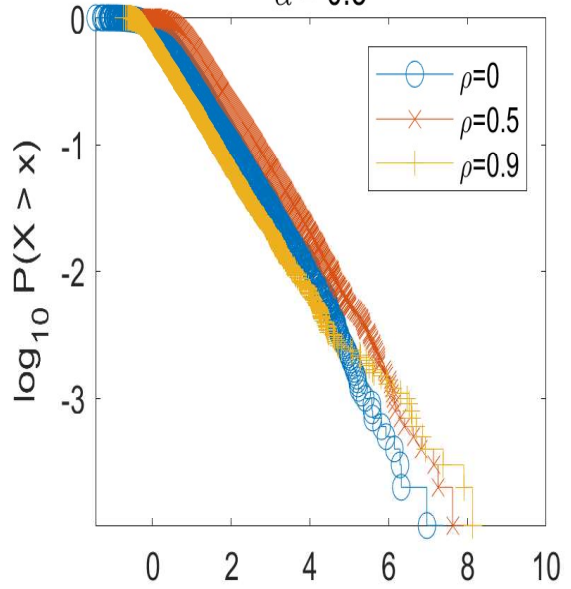


$\log$  sample mean

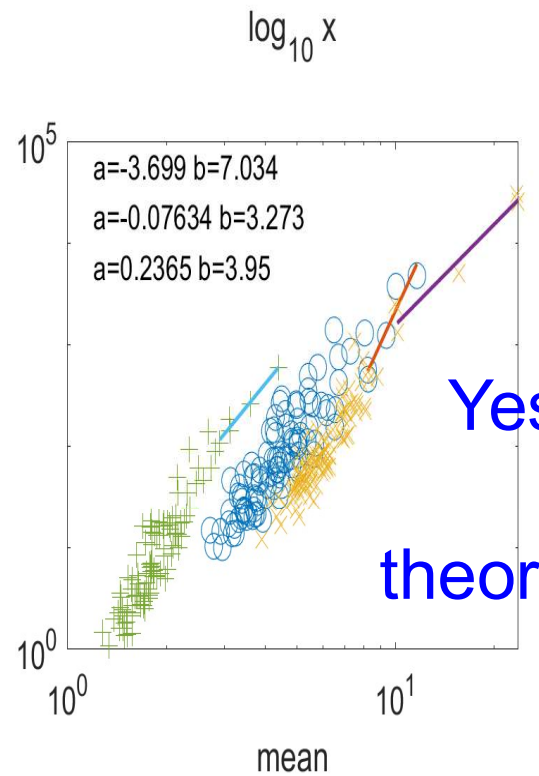
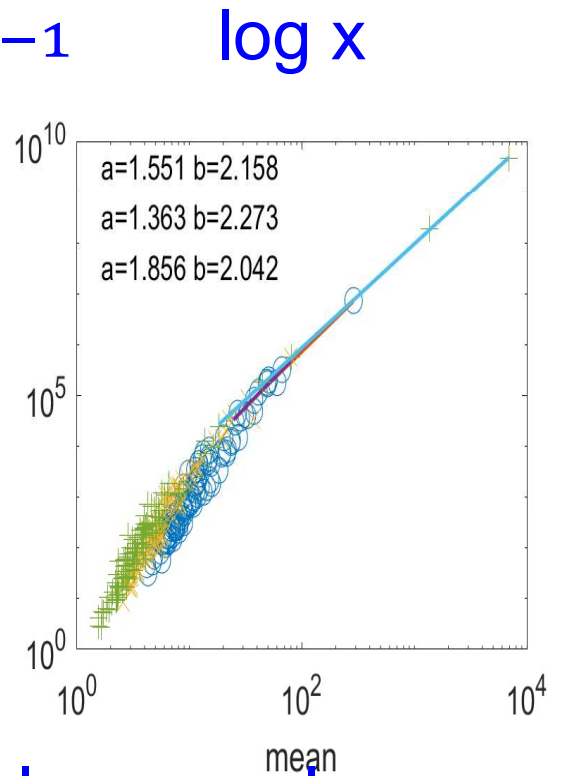
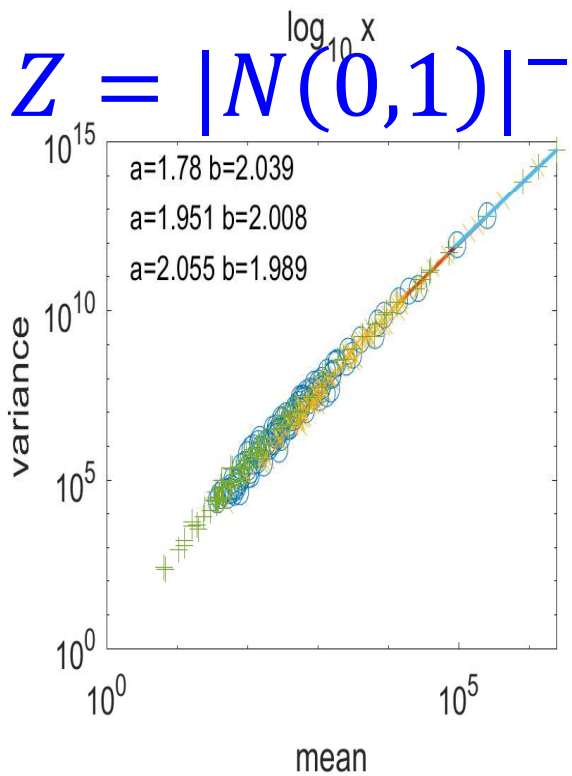
Yes, we have theorems.

$\rho=0, 0.5, 0.9, 100$  samples of size 100

$\log \Pr(X > x)$



$\log$  sample variance



$\log$  sample mean

Yes, we have theorems.

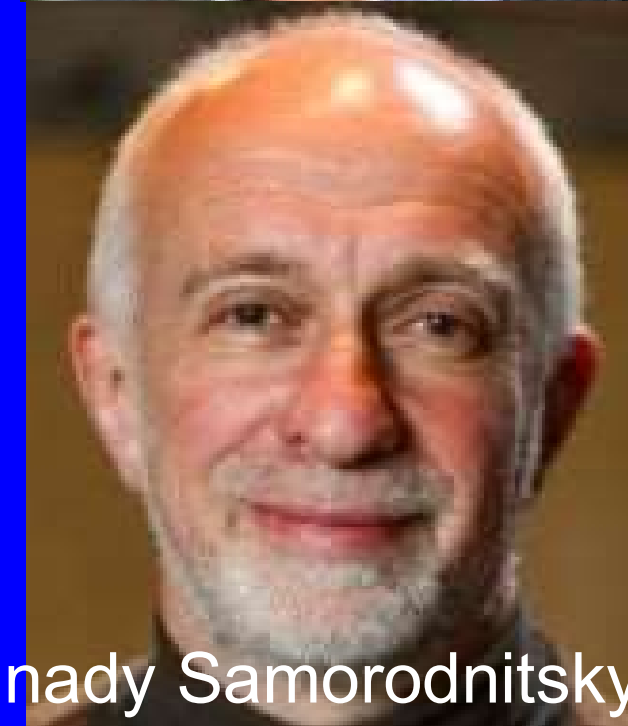
# So what?

If the variances of cases & deaths per county are infinite, facility & resource planning should prepare for unboundedly high counts.

No single county (or state, or other jurisdiction) can prepare for unboundedly high counts.

Cooperative exchanges of support should be planned cooperatively.

# My math collaborators & teachers





Thank you!  
Questions?  
cohen@rockefeller.edu

20190906 La Fage To Florac  
Cévennes "Cham des Bondons  
Chabusse"

2/3/2024

