

COVID-19 cases and deaths in the United States, Taylor's law of fluctuation scaling, and heavy tails

Joint work with

Richard A. Davis & Gennady Samorodnitsky



Joel E. Cohen, cohen@rockefeller.edu

Rockefeller University & Columbia University

2021-09-17, Ecodep Conference (Paul Doukhan),

Paris Seine University, Cergy-Pontoise, France

Ashokan Reservoir, Ulster County, NY, JEC 20151015

Outline

1. COVID-19 U.S. data analysis:
Taylor's law & heavy upper tails
2. Simulations of an idealization
3. Mathematics

A global pandemic of COVID-19

By 12 September 2021, coronavirus SARS-CoV-2 caused >225 million reported cases of COVID-19 disease and >4.6 million deaths.

U.S. reported more cases (41.8 million) and more deaths (678,000) than any other country.

<https://www.worldometers.info/coronavirus/countries-where-coronavirus-has-spread/>

Administrative structure of U.S.

~56 first-level divisions := "states": states, territories, possessions, Washington DC

~3100 second-level subdivisions := "counties":

"... 3,006 counties; 14 boroughs and 11 census areas in Alaska; ... 64 parishes in Louisiana; Baltimore city, Maryland; St. Louis city, Missouri; that part of Yellowstone National Park in Montana; Carson City, Nevada; and 41 independent cities in Virginia."

U.S. Census, *Geographic Areas Reference Manual*, 1994

~55 counties/state \approx ~number of states

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 \rightarrow	$j=1$	$j=2$	$j=3$	$j=\dots$
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}

Cumulative U.S. cases/county by "state" as of June 1, 2021 (from lowest to highest mean)

state	number of cases	mean cases per county	variance cases per county
Virgin Islands	3465	1155	6.7595e+05
Vermont	24224	1614.9	3.1139e+06
South Dakota	1.2419e+05	1881.7	1.8494e+07
⋮	⋮	⋮	⋮
Massachusetts	7.0713e+05	47142	1.7181e+09
Arizona	8.8145e+05	58764	1.9367e+10
California	3.791e+06	65363	3.03e+10

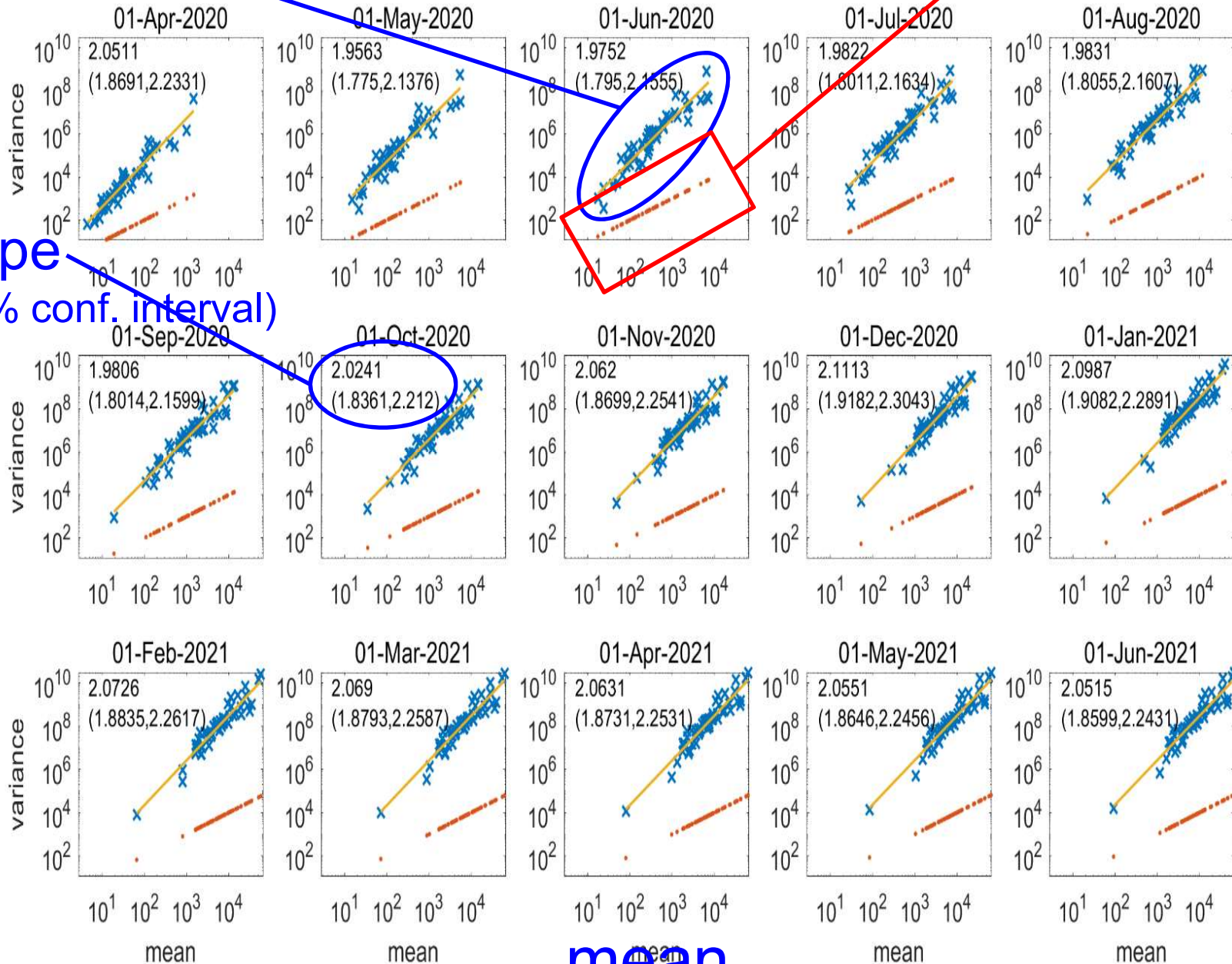
cases

Cumulative U.S. COVID-19 cases/county by state

Poisson

slope
(95% conf. interval)

variance

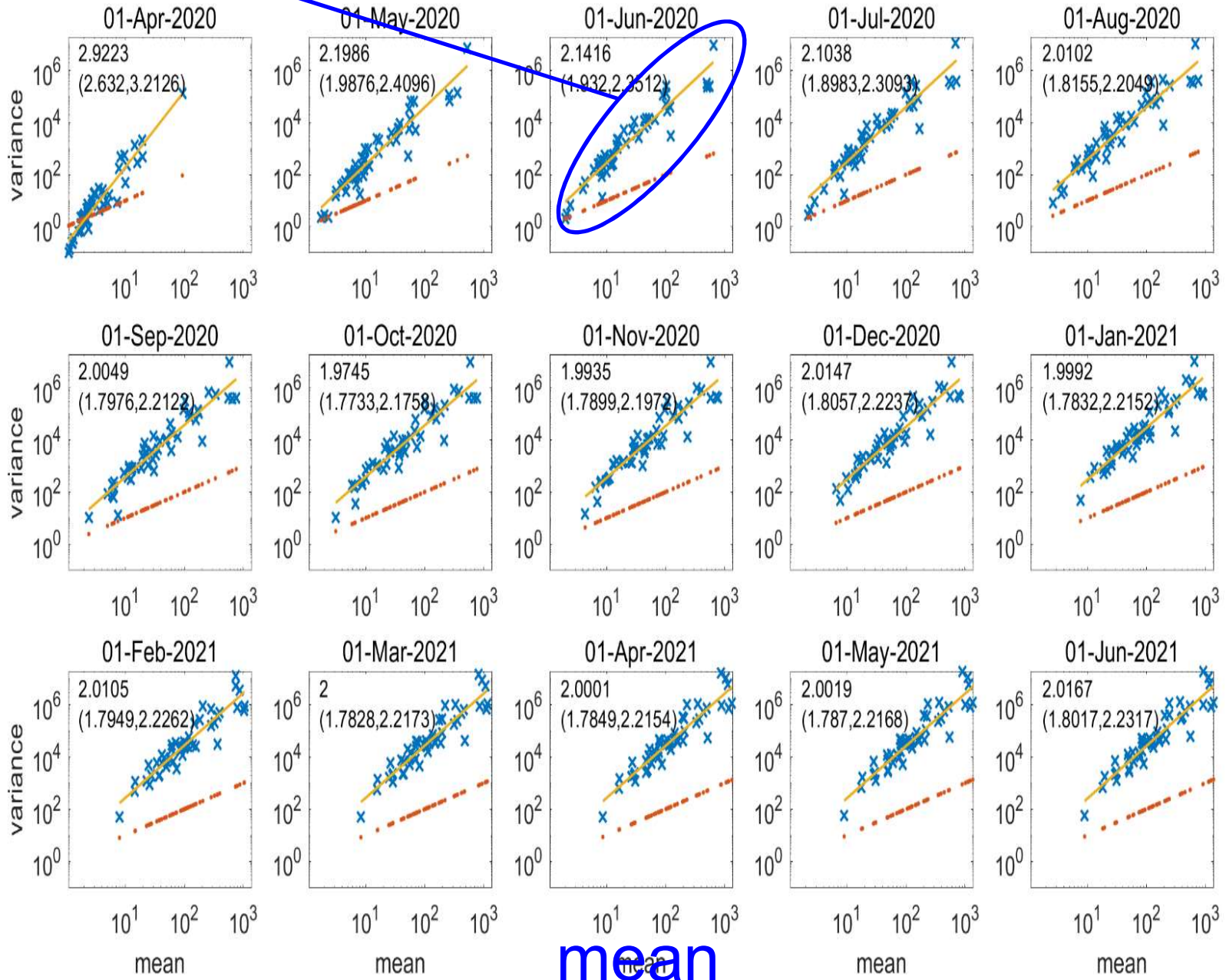


mean

deaths

Cumulative U.S. COVID-19 deaths/county by state

variance



mean

Major findings: counties' cumulative counts of cases or deaths

After first few months,
log variance (over counties) increases, state
by state, linearly as a function of log mean
(over counties), and
slope is close to (not significantly different
from) 2.

Why?

TL data structure: multiple samples, each with multiple observations

Sample number →	$j=1$	$j=2$	$j=3$	$j=\dots$
Counts or nonnegative quantities	x_{11}	x_{12}	x_{13}	x_{\dots}
	x_{21}	x_{22}	x_{23}	\dots
	x_{31}	x_{32}	x_{33}	\dots
		x_{42}	x_{43}	\dots
		x_{52}		\dots
Mean = average	m_1	m_2	m_3	m_{\dots}
Variance	v_1	v_2	v_3	v_{\dots}

Taylor's law *Nature* 1961

In multiple samples, the variance is proportional to a power of the mean.

$$\text{variance} \approx a(\text{mean})^b, a > 0.$$

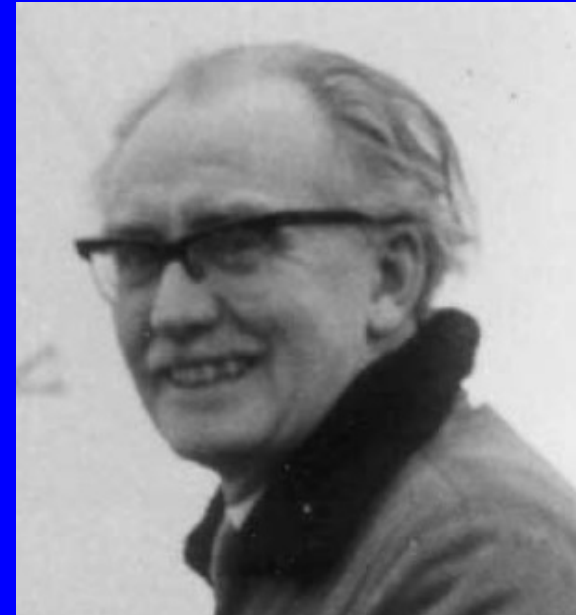
$$\log(\text{variance}) \approx \log(a) + b \cdot \log(\text{mean}).$$

$$\text{variance}/(\text{mean})^b \approx a, \quad a > 0.$$

Taylor measured population density. The pattern applies more widely, but not universally.

Taylor stated no model of deviations from exact equality.

Lionel Roy Taylor
(1924–2007)



Heavy tails

Sample mean, sample variance are always finite. **But** some distributions have infinite mean or infinite variance.

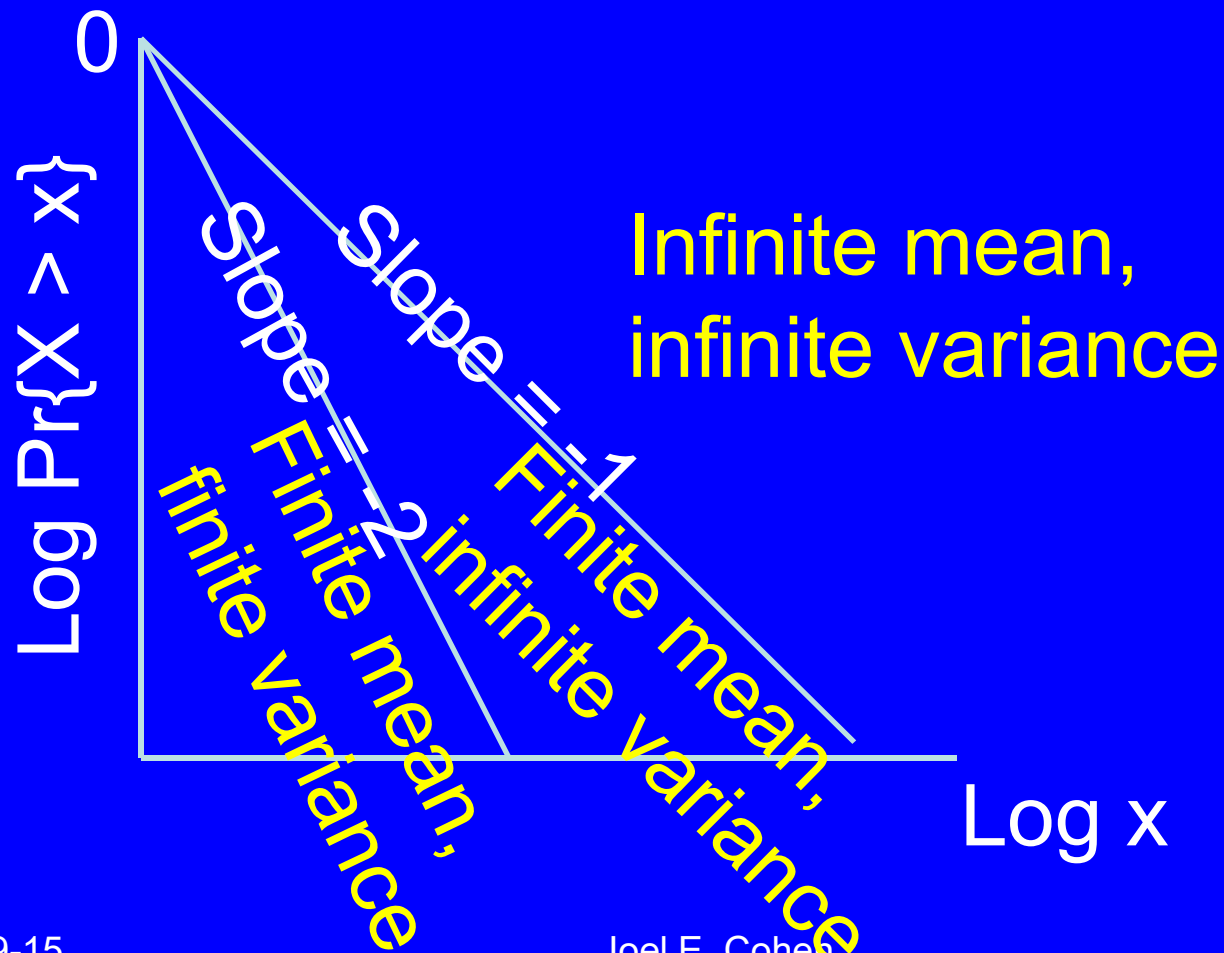
Then sample mean (or sample variance) does not converge to a finite value, but instead moves to infinity, with increasing sample size.



Bell 1806 / Wikipedia

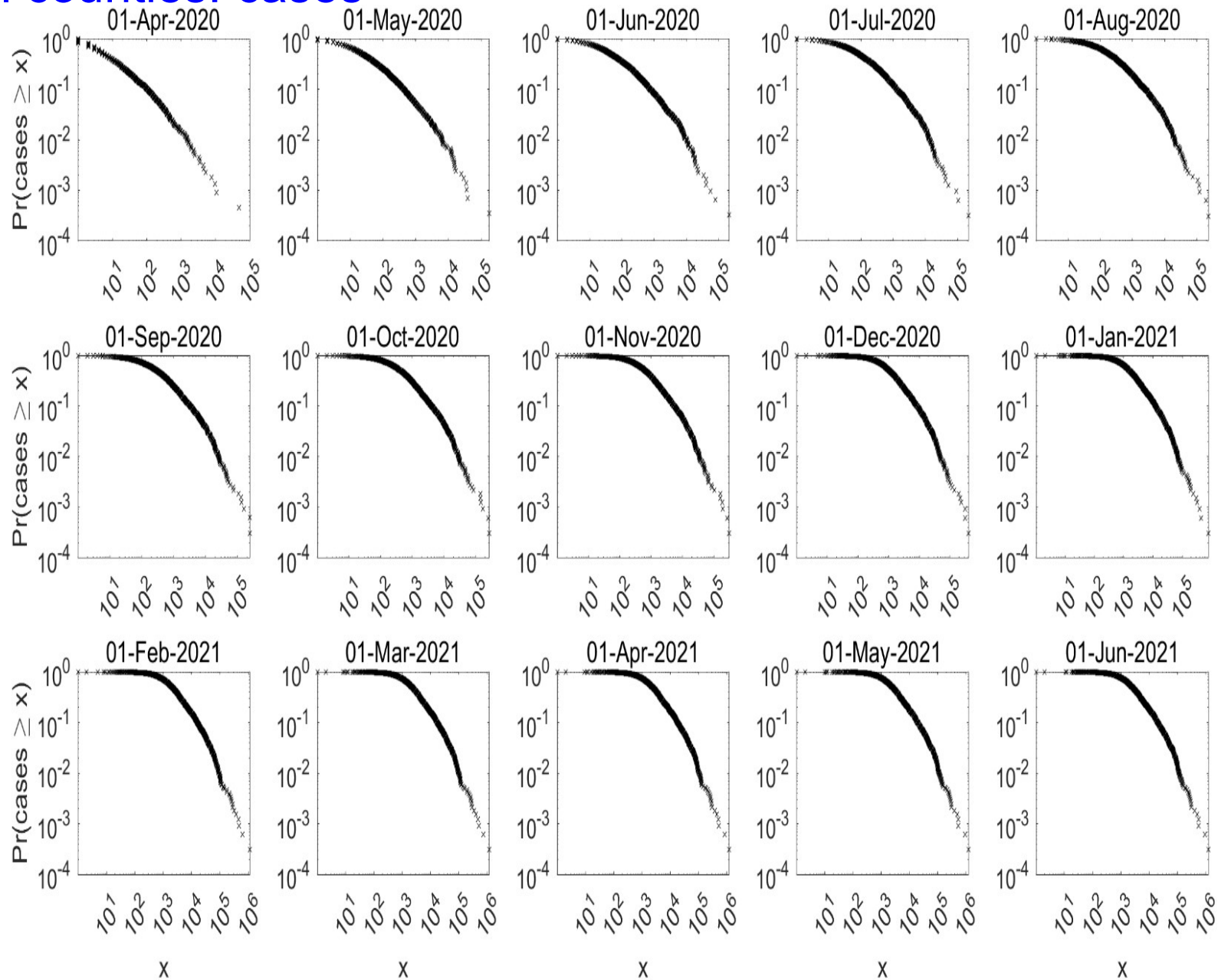
"Wonder / Fear / Astonishment"

Upper tail of survival curve reveals which moments (if any) are finite.



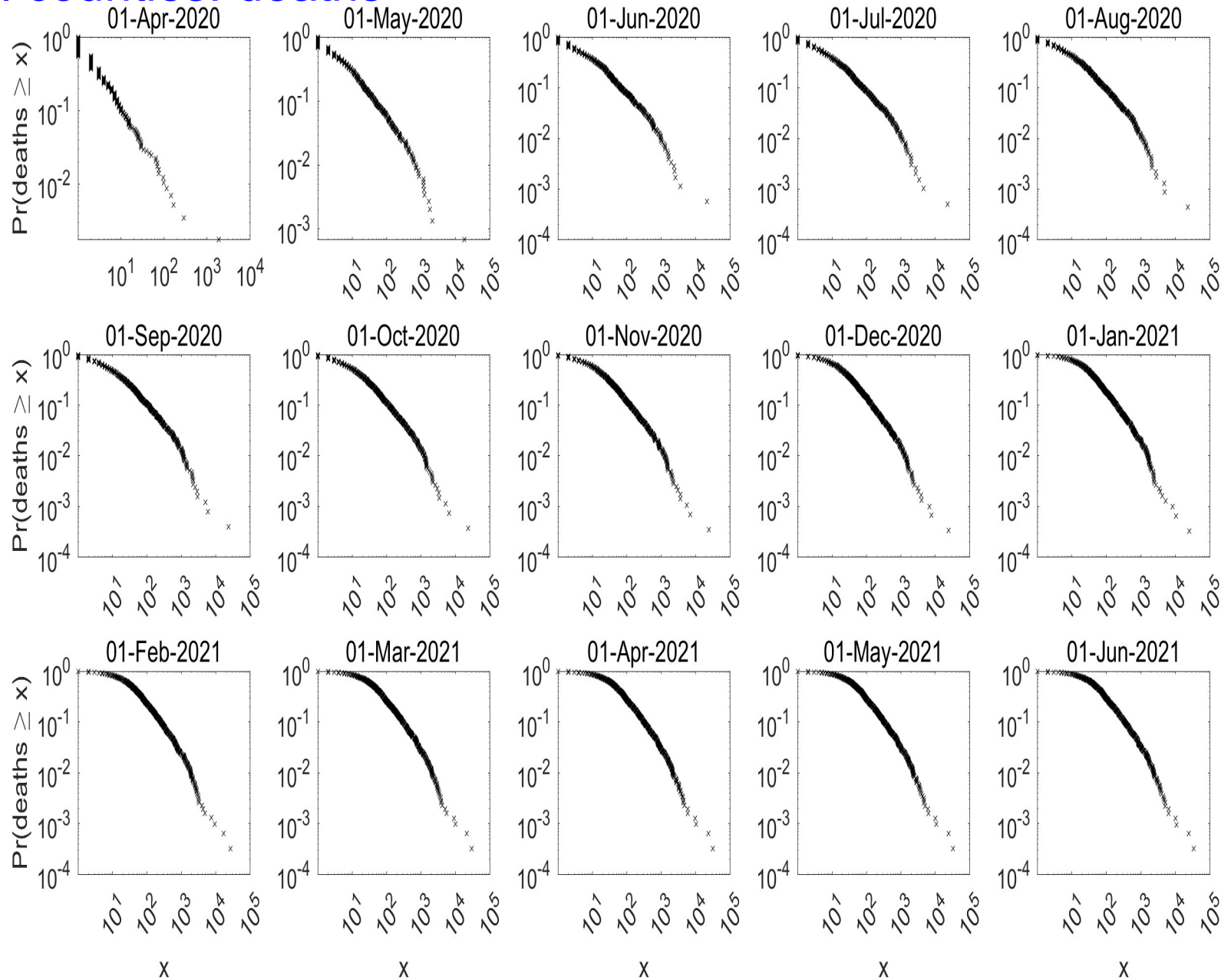
Survival curve of all counties: cases

Cumulative COVID-19 cases/county by date



Survival curve of all counties: deaths

Cumulative COVID-19 deaths/county by date

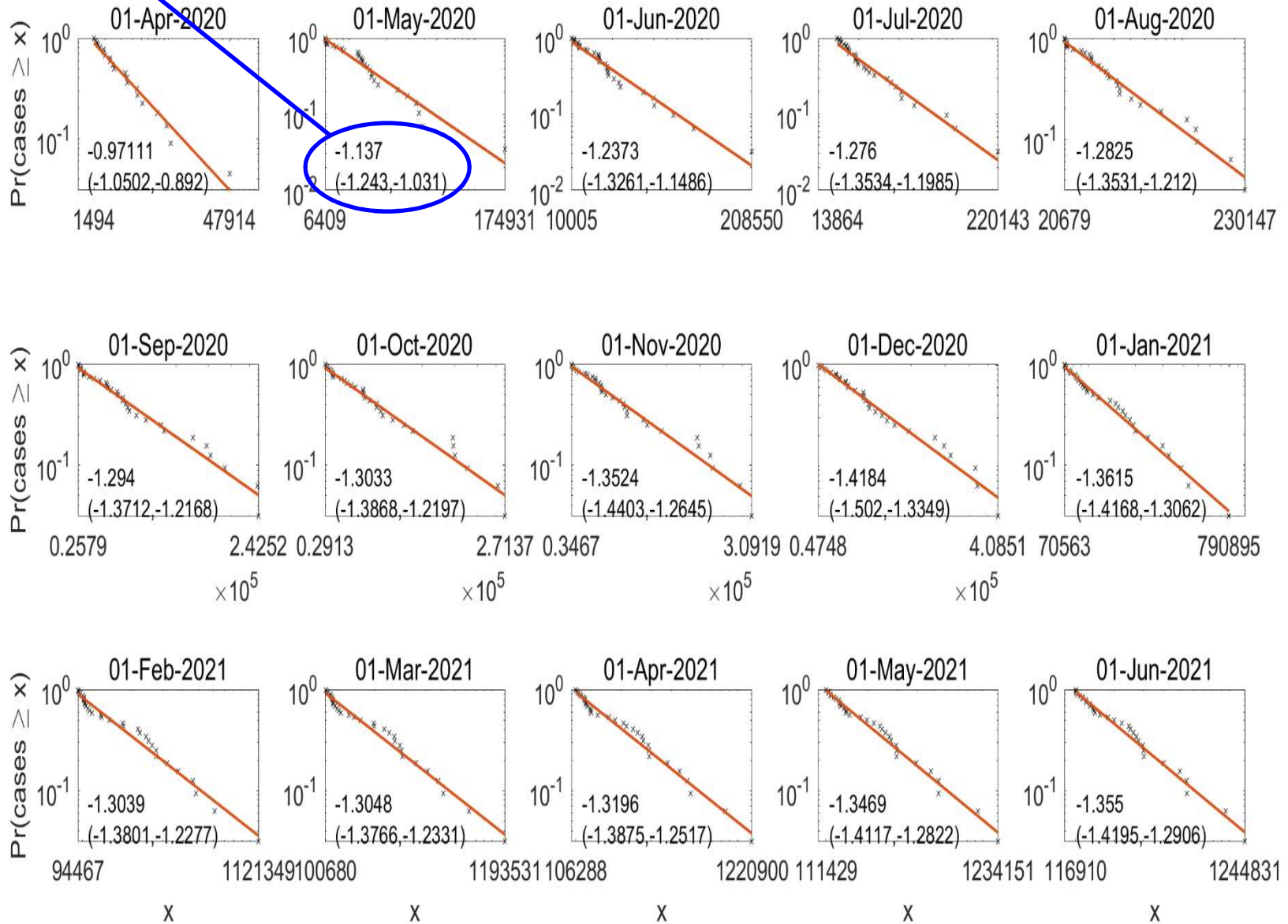


Cases and deaths by county are not Pareto-distributed over their entire range.

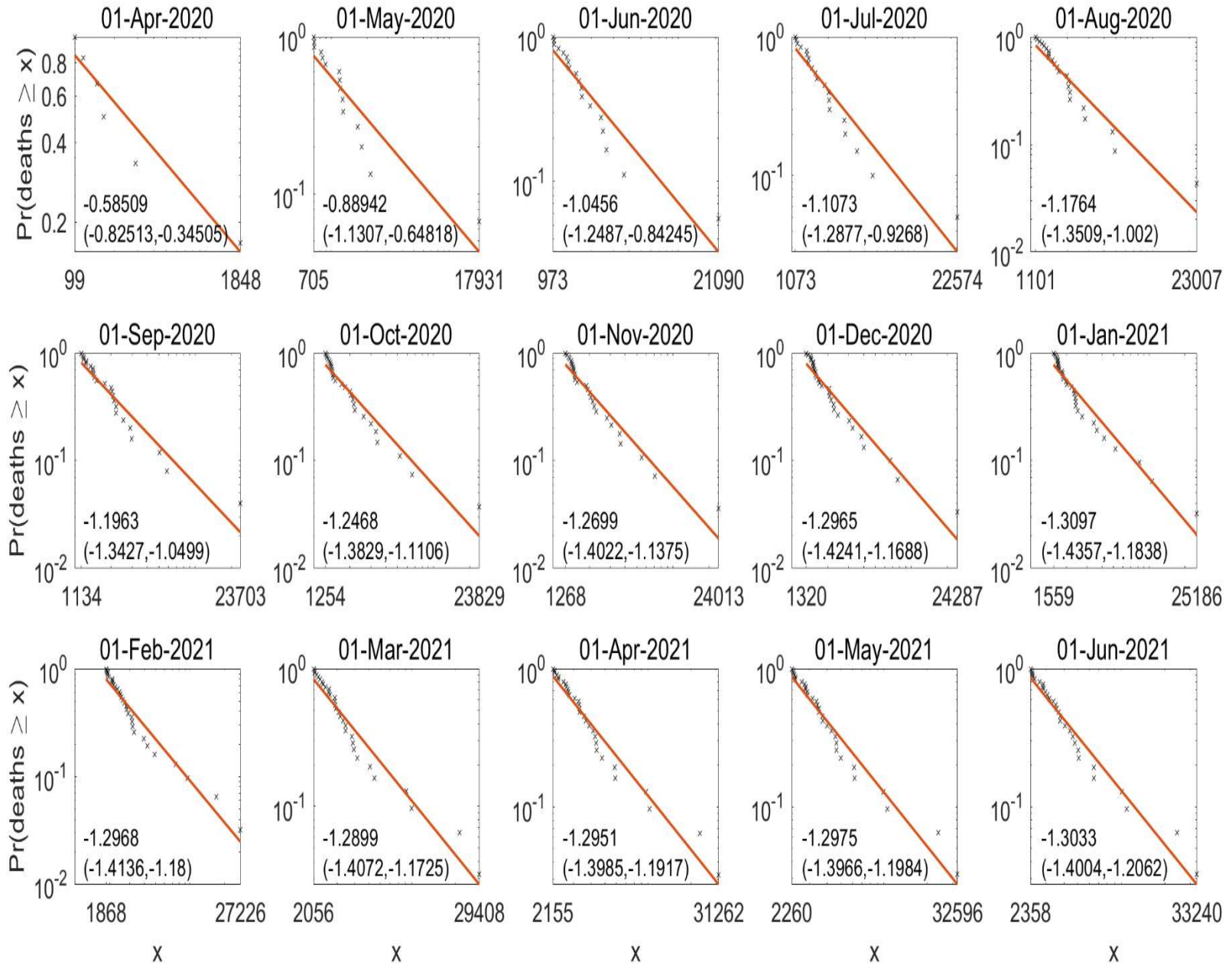
We zoom in to the counties with the highest 1% of numbers of cases or deaths.

slope
(95% conf. interval)

Highest 1% of cumulative COVID-19 cases/county by date



Highest 1% of cumulative COVID-19 deaths/county by date



Major findings: highest 1% of counties' cumulative counts of cases or deaths

After first few months, in the highest 1% of counts of cases or deaths by county, $\log \Pr\{count > x\}$ decreases linearly as $\log x$ increases, and $-1 > \text{slope} > -2$.

Empirical survival curves suggest variance is infinite.

The slope of the upper tail on log-log coordinates is statistically significantly greater than -2 in all 15 months for cases & deaths. Variance is infinite in all cases.

The slope is statistically significantly less than -1 in 14 most recent months for cases & 11 most recent months for deaths. Except at beginning of epidemic, mean is finite.

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.

Wanted: theory of TL for multiple samples of heavy-tailed data

How do sample mean & sample variance behave in multiple samples with fixed or similar sample size, when: data come from distributions with infinite variance or infinite mean, and number of samples \approx each sample size?

Outline

1. COVID-19 U.S. data analysis:
Taylor's law & heavy upper tails
2. Simulations of an idealization
3. Mathematics

Probability distribution with regularly varying tail $\Pr(X > x)$

X is a regularly varying (RV) rv with index $\alpha > 0$ iff

$$\forall t > 0, \lim_{x \rightarrow \infty} \frac{\Pr(X > tx)}{\Pr(X > x)} = \frac{1}{t^\alpha}.$$

We write: $X \in RV(\alpha)$.

Then

$$\Pr(X > x) = \frac{L(x)}{x^\alpha}, \lim_{t \rightarrow \infty} \frac{L(tx)}{L(t)} = 1, \forall x > 0.$$

Probability distribution with regularly varying tail $\Pr(X > x)$

Suppose $X \in RV(\alpha)$.

If $\alpha \in (0,1)$, then mean & all higher moments of X are infinite.

If $\alpha \in (1,2)$, then mean is finite, but variance & all higher moments of X are infinite.

If $\alpha \geq 2$, then X has finite mean & variance.

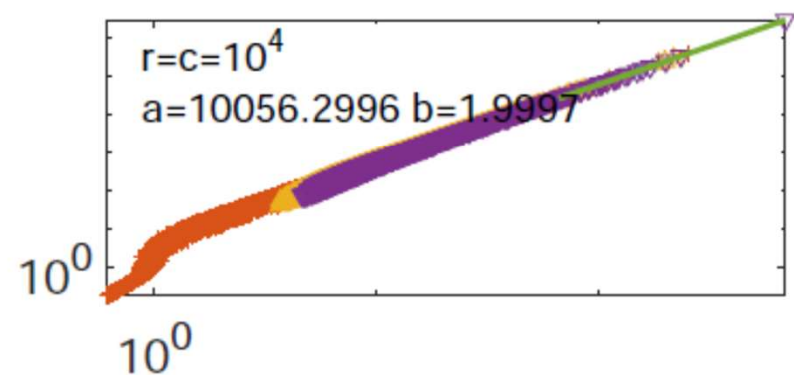
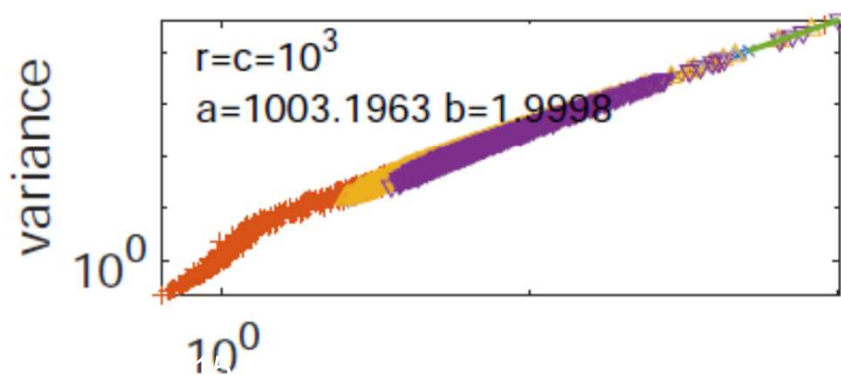
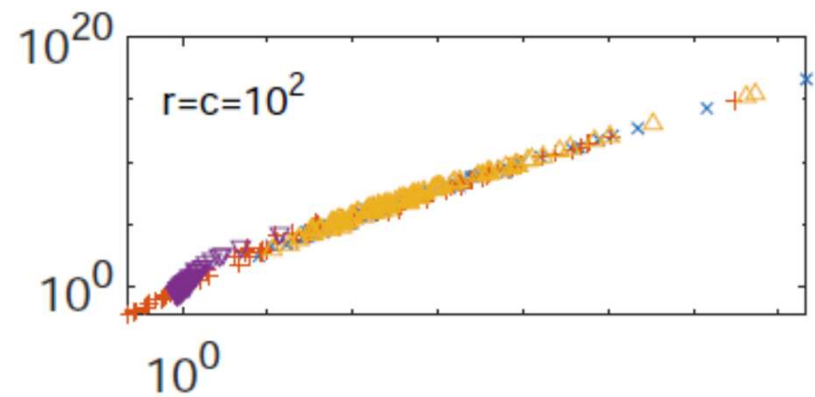
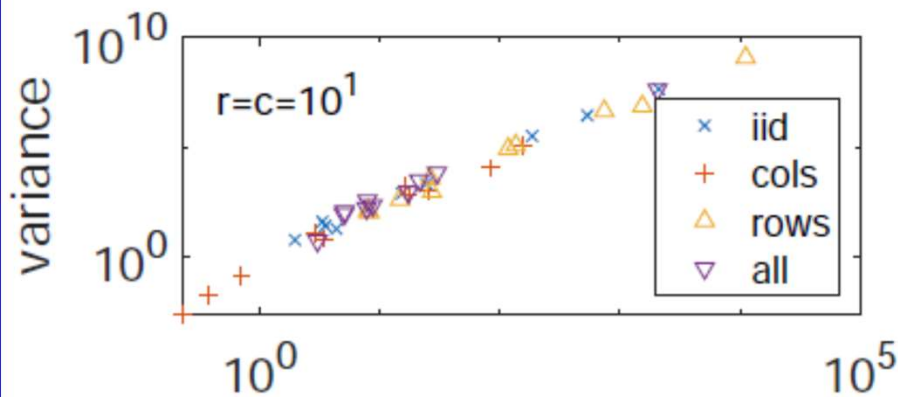
Simulations of idealized cases

We create $r \times c$ matrices, $r = c = 10, 100, 1000, 10^4, 10^5, 10^6$, with $RV(\alpha)$ elements (iid or asymptotically independent within columns, rows, or both), $\alpha \in (0,1) \cup (1,2)$.

For each matrix, we plot log variance of each column as a function of log mean of the same column.

For larger means of large samples, slopes are close to 2.

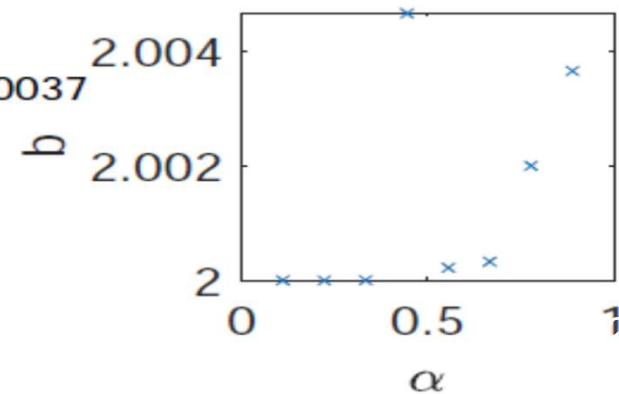
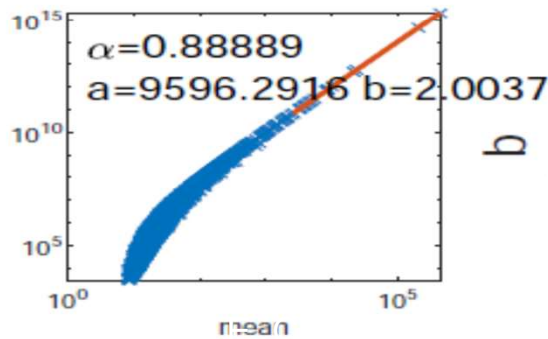
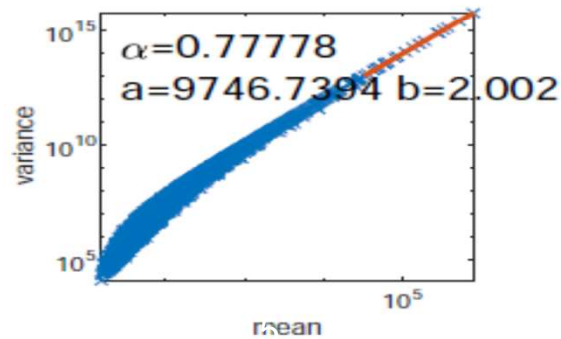
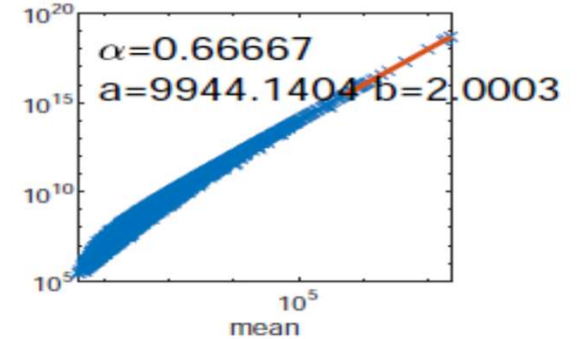
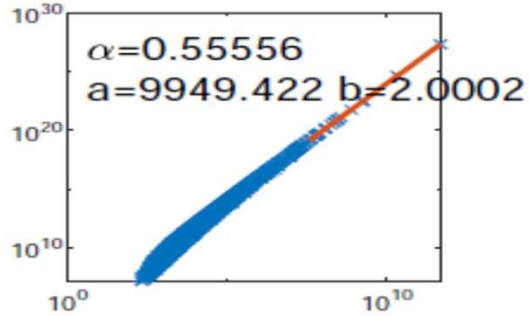
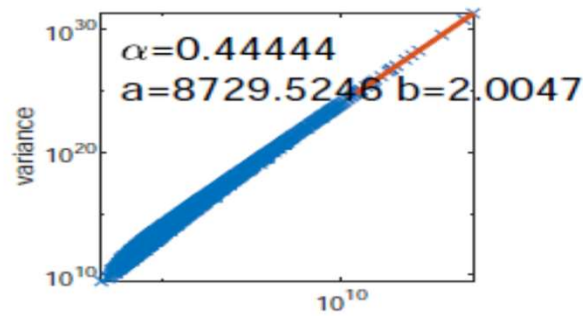
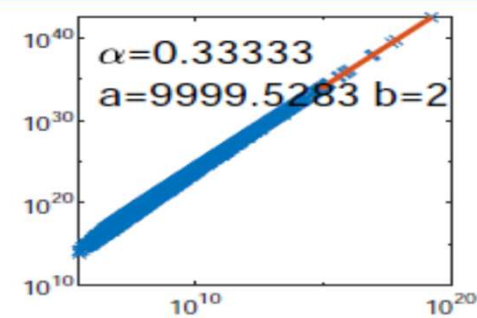
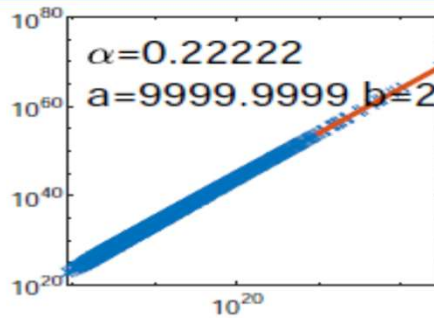
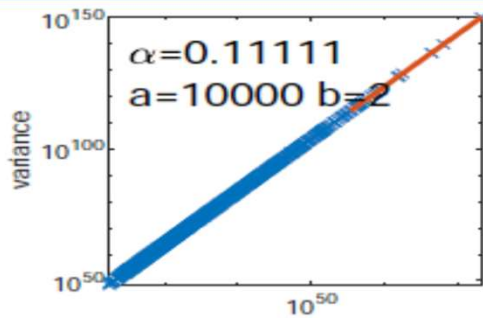
Lévy law: $\alpha = \frac{1}{2}$, $X = |N(0,1)|^{-2}$:
straight line fitted to upper 0.25% of
sampled means has slope ≈ 2 .



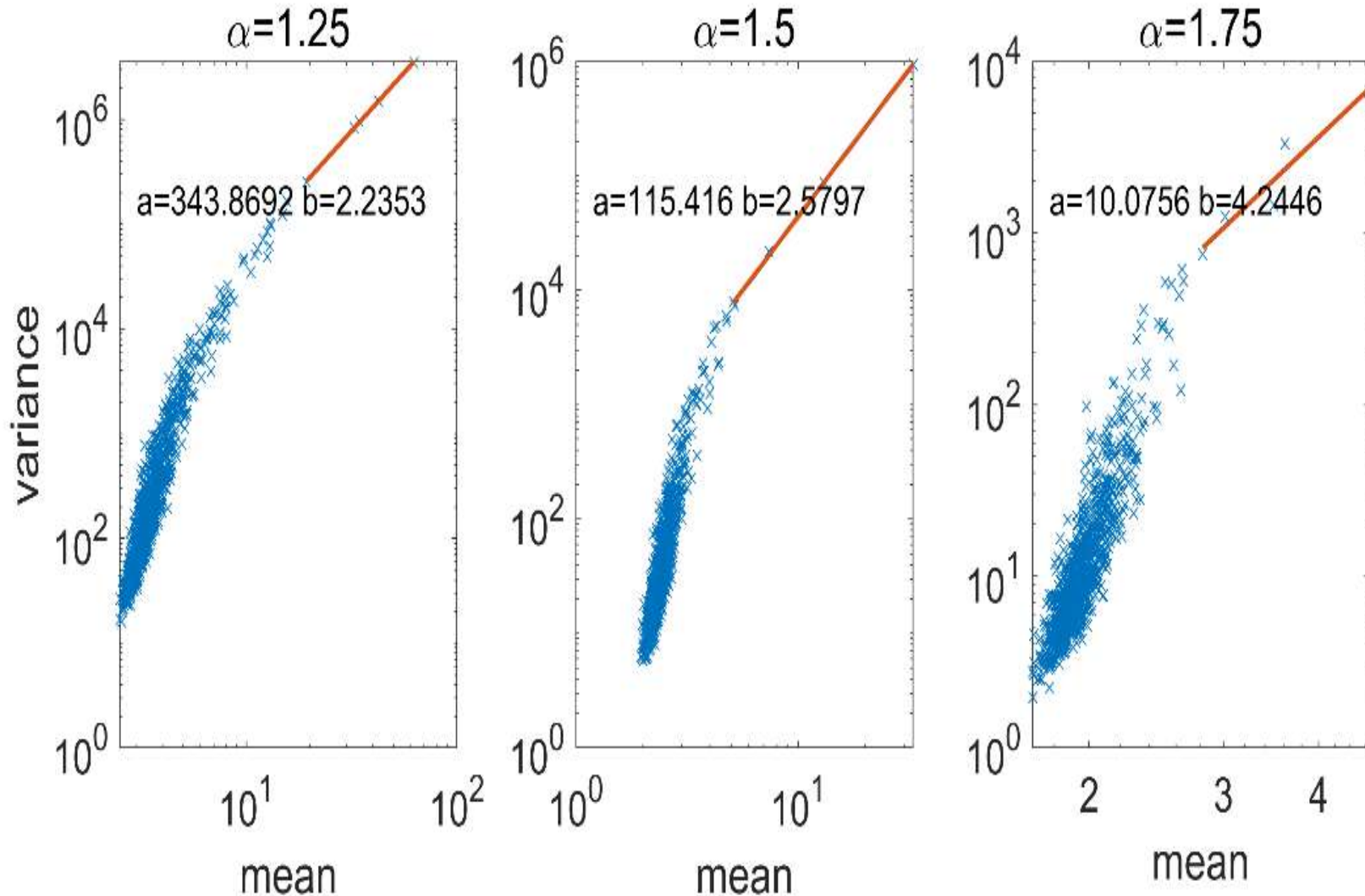
mean

mean

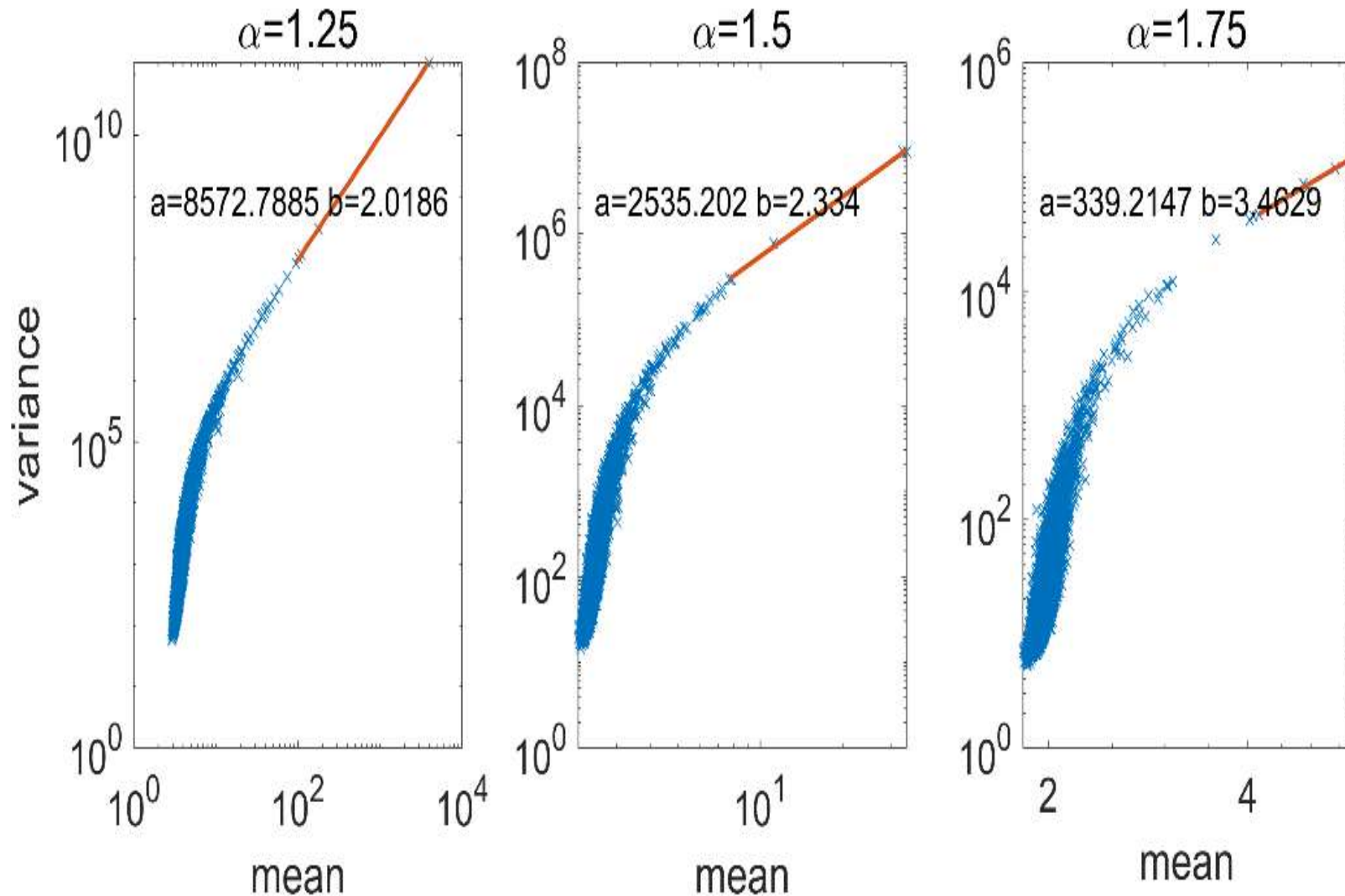
10⁴ samples, each sample size 10⁴, for $\alpha \in (0,1)$



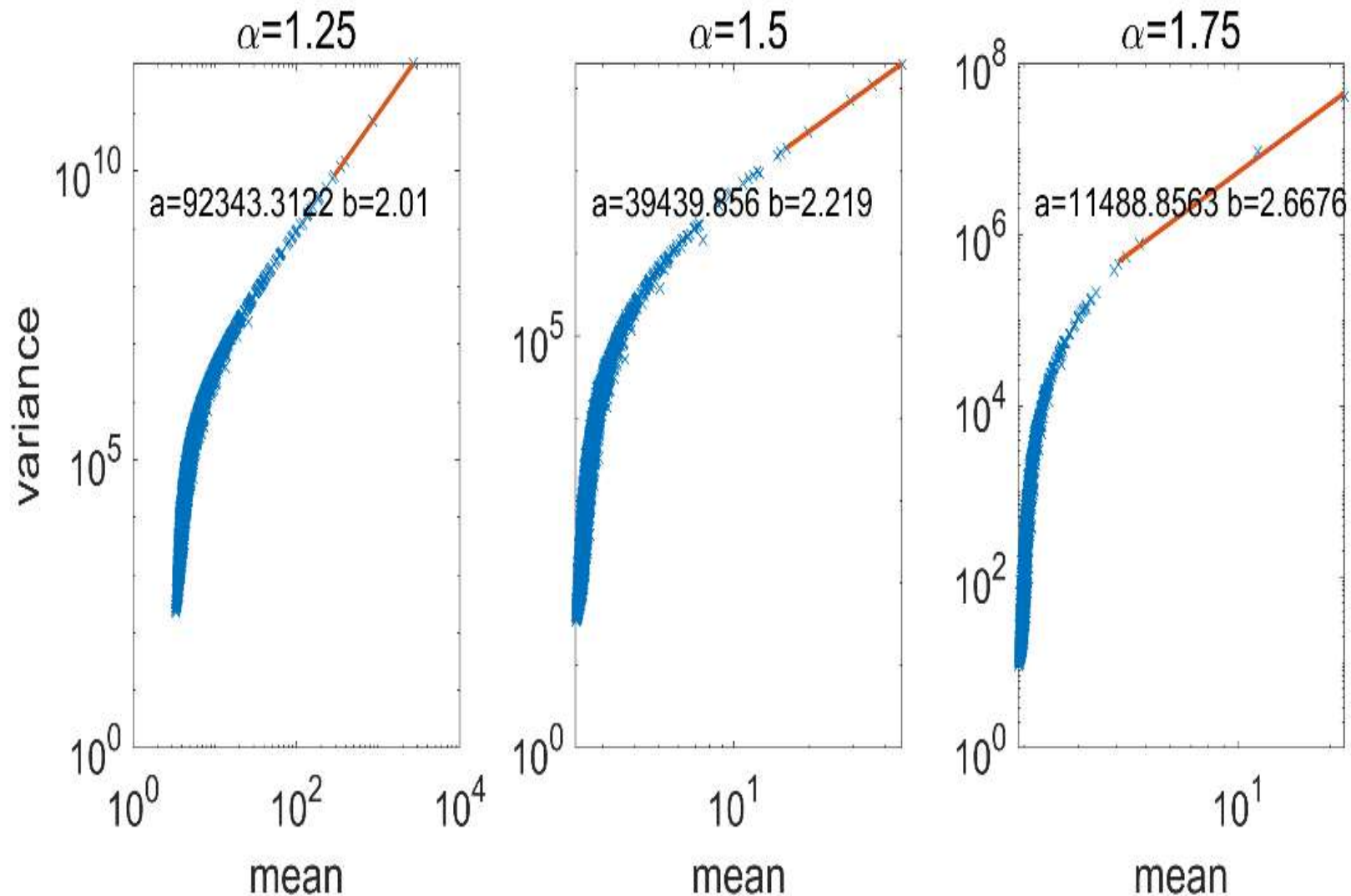
10^3 samples, each sample size 10^3 , for $\alpha \in (1,2)$



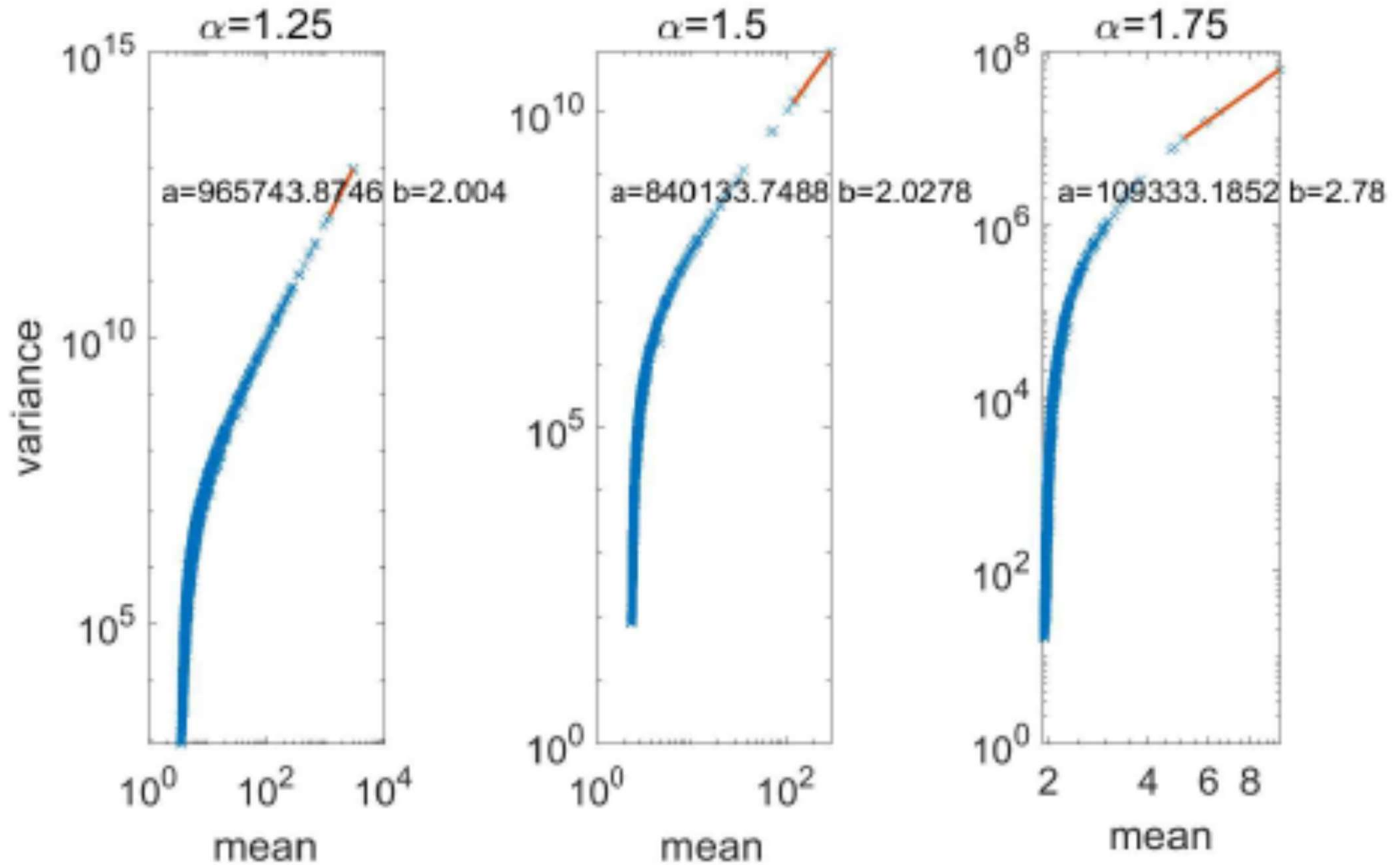
10^4 samples, each sample size 10^4 , for $\alpha \in (1,2)$



10^5 samples, each sample size 10^5 ,
for $\alpha \in (1,2)$



10^6 samples, each sample size 10^6 ,
for $\alpha \in (1,2)$



Estimated slope b of log variance as a function of log mean

for 5 largest values of log mean in c samples, each of size $r = c$

$r = c$	Tail index		
	$\alpha = 1.25$	$\alpha = 1.5$	$\alpha = 1.75$
10^3	2.2353	2.5797	4.2446
10^4	2.0186	2.3340	3.4629
10^5	2.0100	2.2190	2.6676
10^6	2.0040	2.0278	2.7800

Outline

1. COVID-19 U.S. data analysis:
Taylor's law & heavy upper tails
2. Simulations of an idealization
3. Mathematics

Upper tail of multiple samples obeys Taylor's law with slope 2.



Richard A. Davis



Gennady Samorodnitsky

TL holds when $0 < \alpha < 1$: w.h.p.,
 $\log \text{var}(X_{:j}) \approx \log r + 2 \log \bar{X}_{:j}$

Theorem 4.1. *Let $\varepsilon > 0$. Let every element $X_{ij}, i = 1, \dots, r; j = 1, \dots, c$ of the array \mathbf{X} have the same distribution in $RV(\alpha)$, $0 < \alpha < 1$, and assume that all columns are equally distributed. Assume also that the entries within each column are, conditionally on a σ -field \mathcal{G} , independent. Let $x(r)$ satisfy*

$$\lim_{r \rightarrow \infty} x(r) = \infty \quad \text{and} \quad \lim_{r \rightarrow \infty} r \Pr(X_{11} > x(r)) = 0. \quad (4.1)$$

Let the number c of columns depend on the number r of rows in such a way that the function $c = c(r)$ satisfies

$$\lim_{r \rightarrow \infty} c(r)r^2 \max_{i=1, \dots, r} E \left(\left[\frac{1}{x(r)} \int_0^{x(r)} \Pr(X_{i,1} > x \mid \mathcal{G}) dx \right]^2 \right) = 0. \quad (4.2)$$

Then

$$\lim_{r \rightarrow \infty} \Pr \left(\left| \log \frac{\text{var}(X_{:j})}{r(\bar{X}_{:j})^2} \right| > \varepsilon \text{ for some } j = 1, \dots, c(r) \text{ such that } r|\bar{X}_{:j}| > x(r) \right) = 0. \quad (4.3)$$

Assumptions that all X_{ij} are iid,
& that different columns are
mutually independent, can be
relaxed, & same TL holds.

Theorem 4.4. Let $\{Z_{ij}\}$ be an iid space-time process with a regularly varying distribution with index $\alpha \in (0, 1)$. Let $\{G_{ij}\}$ be a stationary space-time Gaussian process with mean 0 and covariance function $\gamma(\cdot, \cdot)$, i.e., $\text{Cov}(G_{ij}, G_{i'j'}) = \gamma(i - i', j - j')$. Given $\mathcal{G} = \sigma(G_{ij}, i = 1, \dots, r; j = 1, \dots, c)$, the process $X_{ij} := Z_{ij}|G_{ij}$ is conditionally independent, but not identically distributed. Under the assumptions (4.1) and (4.2), (4.3) holds.

TL holds when $1 < \alpha < 2$: w.h.p.,
 $\log \text{var}(X_{:j}) \approx \log r + 2 \log \bar{X}_{:j}$

Theorem 4.2. *Let $\varepsilon > 0$. Suppose the elements of the matrix are iid, with the right tail in $RV(\alpha)$, $1 < \alpha < 2$, and the left tail with a finite second moment. Let $x(r)$ satisfy*

$$\lim_{r \rightarrow \infty} x(r) = \infty \quad \text{and} \quad \lim_{r \rightarrow \infty} r \Pr(X_{11} > x(r)) = 0. \quad (4.9)$$

Let the number c of columns depend on the number r of rows in such a way that the function $c = c(r)$ satisfies

$$\lim_{r \rightarrow \infty} c(r) (r \Pr(X_{11} > x(r)))^2 = 0. \quad (4.10)$$

Then

$$\lim_{r \rightarrow \infty} \Pr \left(\left| \log \frac{\text{var}(X_{:j})}{r(\bar{X}_{:j})^2} \right| > \varepsilon \text{ for some } j = 1, \dots, c(r) \text{ such that } r(\bar{X}_{:j} - EX_{11}) > x(r) \right) = 0. \quad (4.11)$$

Relevant prior papers

Brown, M., Cohen, J. E. & de la Peña, V. 2017 Taylor's law, via ratios, for some distributions with infinite mean. *Journal of Applied Probability* 54(3):1-13.
doi:10.1017/jpr.2017.25

Cohen, J. E., Davis, R.A. & Samorodnitsky, G. 2020 Heavy-tailed distributions, correlations, kurtosis, and Taylor's law of fluctuation scaling. *Proceedings of the Royal Society A* 476:20200610.
<https://doi.org/10.1098/rspa.2020.0610>.

Thank you!
Questions?
cohen@rockefeller.edu

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

2021-09-15

