

The Difference-of-Log-Normals Distribution is Fundamental in Nature

Robert Parham (✉ robertp@virginia.edu)

University of Virginia <https://orcid.org/0000-0003-1462-1839>

Social Sciences - Article

Keywords:

DOI: <https://doi.org/>

License:   This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

Additional Declarations: There is **NO** Competing Interest.

The Difference-of-Log-Normals Distribution is Fundamental in Nature

Robert Parham*

Friday 10th February, 2023

Summary

The growth of many natural and social phenomena including pandemics¹, firms,² cities,³ and various economic indices,^{4,5} is known to be heavy-tailed. Most growth is modest, but we often observe explosive growth rates, such as firms doubling or halving in size within a short period. Neither a simple explanation nor a well-fitting distributional form for these growth phenomena is known. Here we show that a hitherto obscure statistical distribution — the Difference-of-Log-Normals (DLN) — describes a plethora of growth phenomena remarkably well, and discuss why it arises as a natural consequence of the Central Limit Theorem (CLT). Our results demonstrate how growth phenomena subject to opposing random exponential forces are likely to distribute DLN. This provides both a framework for scientifically modeling these phenomena and a simple distributional form to be used when empirically modeling observed heavy-tailed growth. We hence posit that the DLN is a *fundamental distribution in nature*, in the sense that it emerges in many disparate natural phenomena, especially growth phenomena, similar to the repeated disparate emergence of the Normal and log-Normal distributions.

*University of Virginia (robertp@virginia.edu).

21 What do the growth rates of such varied phenomena as new COVID-19 infection cases,
 22 tests conducted, and vaccinations administered; firm sales, capital, income, and stock values;
 23 and city populations and regional GDPs have in common? They all appear to distribute as
 24 the Difference-of-Log-Normals (Figures 1- 3).

25 The Difference-of-Log-Normals distribution, henceforth DLN, is the distribution arising
 26 when one subtracts a log-Normal random variable (RV) from another. To define the DLN,
 27 consider an RV W such that

$$W = Y_p - Y_n = \exp(X_p) - \exp(X_n) \quad \text{with} \quad \mathbf{X} = (X_p, X_n)^T \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \quad (1)$$

28 in which \mathbf{X} is a bi-variate Normal with

$$\boldsymbol{\mu} = \begin{bmatrix} \mu_p \\ \mu_n \end{bmatrix} \quad \boldsymbol{\Sigma} = \begin{bmatrix} \sigma_p^2 & \sigma_p \cdot \sigma_n \cdot \rho_{pn} \\ \sigma_p \cdot \sigma_n \cdot \rho_{pn} & \sigma_n^2 \end{bmatrix} \quad (2)$$

29 We say W follows the five-parameter DLN distribution, i.e. $W \sim \text{DLN}(\mu_p, \sigma_p, \mu_n, \sigma_n, \rho_{pn})$,
 30 and fully derive its properties elsewhere.⁶

31 The *sum* of log-Normal RVs has been used in several disciplines including telecommu-
 32 nication, actuary, insurance, and derivative valuation. The DLN, in contrast, is almost
 33 completely unexplored. At the time of writing, we were unable to find instances of using it
 34 anywhere in the sciences, and only two statistical works considering it.^{7,8} Both papers concen-
 35 trate on the sum of log-Normals but show their results hold for the difference of log-Normals
 36 as well, under some conditions. Nevertheless, we posit that the DLN is a *fundamental dis-*
 37 *tribution in nature*, likely describing a plethora of natural and economic phenomena.⁹

38 To see this, consider first the central limit theorems (CLTs), which state that

$$Y^+ = \lim_{K \rightarrow \infty} \frac{1}{K} \sum_{i=1}^K X_i^+ \sim \mathcal{N} \quad (3)$$

39 for $X_i^+ \sim \Omega_i^+$ under mild regularity conditions on the Ω_i^+ depending on the version of the

40 CLT used. Put differently, the CLTs state that a phenomenon in nature which is an additive
 41 combination of many latent random forces will tend to distribute Normally.

42 Consider next the multiplicative CLT, sometimes known as ‘‘Gibrat’s law’’,¹⁰ which states
 43 that

$$Y^* = \lim_{K \rightarrow \infty} \left(\prod_{i=1}^K X_i^* \right)^{\frac{1}{K}} \sim \log\mathcal{N} \quad (4)$$

44 for $X_i^* > 0 \sim \Omega_i^*$ under similarly mild regularity conditions. Put differently, the multiplica-
 45 tive CLT states that a phenomenon in nature which is a *product* of many latent random
 46 forces will tend to distribute log-Normally. Many physical and economic non-negative quan-
 47 tities, such as mass, population count, epidemic spread, interest rates, firm sales, and firm
 48 value are products of latent random factors and are approximately log-Normally distributed.

49 Finally, consider a natural phenomenon impacted by two main forces operating in oppo-
 50 site directions, i.e., $W = Y_p - Y_n$. If the two main forces are additive combinations of latent
 51 random forces,

$$W^+ = Y_p^+ - Y_n^+ = \lim_{K_p \rightarrow \infty} \frac{1}{K_p} \sum_{i=1}^{K_p} X_i^+ - \lim_{K_n \rightarrow \infty} \frac{1}{K_n} \sum_{j=1}^{K_n} X_j^+ \sim \mathcal{N} \quad (5)$$

52 then the natural phenomenon will tend to distribute Normally as well, because the difference
 53 of two Normal RV is itself Normal, under mild conditions. But the same is not true if the
 54 two main forces are multiplicative combinations of latent random factors. In this case,

$$W^* = Y_p^* - Y_n^* = \lim_{K_p \rightarrow \infty} \left(\prod_{i=1}^{K_p} X_i^* \right)^{\frac{1}{K_p}} - \lim_{K_n \rightarrow \infty} \left(\prod_{j=1}^{K_n} X_j^* \right)^{\frac{1}{K_n}} \sim \text{DLN} \quad (6)$$

55 because the difference between two log-Normal RVs does not collapse to a log-Normal RV.

56 To fix ideas, Figure 4 presents several instances of the DLN distribution. Panel (a)
 57 presents and contrasts the standard Normal, standard DLN, and standard log-Normal. The
 58 standard DLN is defined as $\text{DLN}(0,1,0,1,0)$, i.e. the difference between two exponentiated
 59 uncorrelated standard Normal RVs. Panel (b) shows the role of the correlation coefficient ρ_{pn} ,

60 controlling tail-weight vs. peakedness. Panel (c) repeats the analysis of Panel (b) for a dif-
 61 ferent parametrization common in practical applications,¹¹ exhibiting the problem of dealing
 62 with the DLN’s characteristic heavy tails in both the positive and negative directions. Panel
 63 (d) presents the data of panel (c) after taking an Inverse Hyperbolic Sine (asinh) transform
 64 of the data. The asinh acts as a log transform in both the positive and negative directions,
 65 allowing us to observe the characteristic “double Normal” shape of the transformed DLN.

66 Possibly the most intuitive example of the DLN’s emergence is in the context of a simple
 67 population dynamics (“birth-death”) model.¹² Denote $N(t)$ the size of the population in
 68 some closed natural habitat (with no immigration or emigration) at time t . The population
 69 dynamics of the system are described by the ordinary differential equation:

$$\frac{dN(t)}{dt} = b(t) \cdot N(t) - d(t) \cdot N(t) = N(t) \cdot [b(t) - d(t)] \quad (7)$$

70 in which $b(t) \geq 0$ and $d(t) \geq 0$ are the instantaneous birth and death rates. Generally, $b(t)$
 71 and $d(t)$ are stochastic, depending on some underlying latent forces such as food availability,
 72 climate, predation, etc. Because negative birth or death rates are inadmissible, we cannot
 73 assume they are jointly Normal. The next-simplest hypothesis (in the maximum entropy
 74 sense) for their distribution is hence the bi-variate log-Normal. This means the distribution
 75 of their difference, or the distribution of population growth in the model, is DLN — providing
 76 an intuitive explanation to the emergence of DLN in the COVID-19 data, described in
 77 Figure 1.

78 Moving on to the realm of finance, consider the most fundamental “sources and uses”
 79 equation of the firm: $income = sales - expenses$. Both sales and expenses are approximately
 80 log-Normally distributed, and it is standard practice in neo-classical economics to model
 81 income as a controlled AR(1) stochastic Markov process in logs, with Normal innovations.
 82 In such models, growth is counter-factually Normally distributed. If we instead model sales
 83 and expenses *separately* as two co-controlled AR(1) stochastic Markov processes in logs,

84 with (possibly correlated) Normal innovations, the model then predicts firm income (and
 85 consequently, growth) will distribute DLN. In essence, we replace the neo-classical log-linear,
 86 or ‘‘Cobb-Douglas’’, production function

$$\mathbf{Y}_z(K_t, Z_t) = \underbrace{Z_t \cdot K_t^{\theta_Z}}_{Income} = \exp(z_t + \theta_Z \cdot k_t) \quad (8)$$

87 with a difference-of-log-linears production function explicitly modeling sales and expenses

$$\mathbf{Y}_{sx}(K_t, S_t, X_t) = \underbrace{S_t \cdot K_t^{\theta_S}}_{Sales \equiv \mathbb{S}_t} - \underbrace{X_t \cdot K_t^{\theta_X}}_{Expenses \equiv \mathbb{X}_t} = \exp(s_t + \theta_S \cdot k_t) - \exp(x_t + \theta_X \cdot k_t) \quad (9)$$

88 With the stochastic log-productivity variables z_t, s_t, x_t following AR(1) with Normal innova-
 89 tions, $\theta_Z, \theta_S, \theta_X$ returns-to-scale coefficients, and k_t logged firm capital. Importantly, Equa-
 90 tion 9 can be decomposed such that

$$\mathbf{Y}_{sx}(K_t, S_t, X_t) = 2 \cdot \exp(\lambda_t) \cdot \sinh(\tau_t)$$

$$\lambda_t = \lambda(k_t, s_t, x_t) = \frac{s_t + x_t}{2} + \frac{\theta_S + \theta_X}{2} \cdot k_t = \log(\sqrt{\mathbb{S}_t \cdot \mathbb{X}_t}) \quad (10)$$

$$\tau_t = \tau(k_t, s_t, x_t) = \frac{s_t - x_t}{2} + \frac{\theta_S - \theta_X}{2} \cdot k_t = \log(\sqrt{\mathbb{S}_t / \mathbb{X}_t})$$

91 with λ_t firm *scale*, and τ_t firm *efficiency*, thus justifying our use of the asinh transform.

92 To test these predictions of the model, we empirically analyze the data on all public
 93 US firms in the 50-year period 1970-2019. Figure 2 graphically presents several of our
 94 findings. Panels (a)-(c) show the distribution of firm income, with Panel (a) showing the
 95 un-transformed but truncated data, Panel (b) showing the asinh-transformed data, and
 96 Panel (c) showing its q-q plot vs. the DLN, with excellent fit. The next six panels exhibit
 97 similarly excellent visual fits for: the growth of firm sales, the Fama-French factor-adjusted
 98 equity returns at monthly frequency, and net total yearly investment. The fit between the

99 DLN and equity returns is especially noteworthy, given the voluminous literature on the
100 determinants, fat-tails, and statistical properties of equity returns.

101 In a battery of tests, the DLN is shown to be the core distribution driving firm dynamics.¹¹

102 The DLN is not rejected for

- 103 • Firm income: including both free cash flows and disbursements to/from stakeholders.
- 104 • Firm growth in: sales, expenses, capital, total value, and market value of equity.
- 105 • Firm returns: yearly, monthly, daily, both raw and adjusted for Fama-French factors.
- 106 • Firm income growth: growth in cash flows and disbursements.
- 107 • Firm (net) investment: both total and physical investment net of asset sales.

108 while the typical candidates in the literature — the Normal, Laplace, and Lévy-Stable (aka
109 Pareto-Stable or Power-Law) are generally strongly rejected. In likelihood-based information
110 criteria “horse-race” tests, such as Akaike’s Information Criterion (AIC) or the Bayesian
111 Information Criterion (BIC), the DLN is overwhelmingly favored over the other candidates.

112 Returning to population dynamics, Figure 3 presents data on population growth,¹³ and
113 on economic activity growth by county and by metropolitan area and industry from the US
114 Bureau of Economic Analysis. The DLN arises again, as can be seen visually in the q-q plots.
115 A “horse-race” with the other typical candidate distributions again strongly favors the DLN.
116 This finding is a natural outcome of a simple model of city dynamics¹⁴ in which cities grow
117 subject to the interplay between agglomeration benefits and congestion costs, both of which
118 exert exponential influence on the flow of aggregate economic value created by cities. This
119 economic flow is captured by city inhabitants, firms operating in the city, the government,
120 or in general the social planner. Both benefits and costs increase with city size, just as both
121 sales and expenses increase with firm size. But the interplay between them may give rise to
122 positive net economic flow (i.e., $Y_{sx} > 0$) leading to immigration into the city, or to negative
123 net economic flow (i.e., $Y_{sx} < 0$) leading to emigration out of the city.

124 A plethora of phenomena in nature arise as the balance of two opposing forces. When
125 these two forces are themselves multiplicative combinations of underlying latent random
126 forces, the phenomena will tend to distribute DLN. Growth phenomena are especially likely
127 to be DLN, as growth is at its essence a multiplicative (i.e. exponential) process. Hence, the
128 forces supporting growth and the forces opposing it are likely to be log-Normal, and growth
129 itself is likely to distribute DLN. This insight is useful both when constructing models to
130 describe these phenomena, as briefly outlined above for firms and cities, and when empirically
131 modeling such phenomena by providing a simple distributional form to estimate and use.

References

- [1] Kris V. Parag, Christl A. Donnelly, and Alexander E. Zarebski. “Quantifying the Information in Noisy Epidemic Curves”. In: *Nature Computational Science* 2.9 (Sept. 2022), pp. 584–594. ISSN: 2662-8457. DOI: [10.1038/s43588-022-00313-1](https://doi.org/10.1038/s43588-022-00313-1).
- [2] T. S. Ashton. “The Growth of Textile Businesses in the Oldham District, 1884-1924”. In: *Journal of the Royal Statistical Society* 89.3 (1926), pp. 567–583. ISSN: 0952-8385. DOI: [10.2307/2341720](https://doi.org/10.2307/2341720).
- [3] Jan Eeckhout. “Gibrat’s Law for (All) Cities”. In: *The American Economic Review* 94.5 (2004), pp. 1429–1451. ISSN: 0002-8282.
- [4] Rosario N. Mantegna and H. Eugene Stanley. “Scaling Behaviour in the Dynamics of an Economic Index”. In: *Nature* 376.6535 (July 1995), p. 46. ISSN: 1476-4687. DOI: [10.1038/376046a0](https://doi.org/10.1038/376046a0).
- [5] Michael H. R. Stanley et al. “Scaling Behaviour in the Growth of Companies”. In: *Nature* 379.6568 (Feb. 1996), p. 804. ISSN: 1476-4687. DOI: [10.1038/379804a0](https://doi.org/10.1038/379804a0).
- [6] Robert Parham. “The Difference-of-Log-Normals Distribution: Properties, Estimation, and Growth”. In: *arXiv:2302.02486 [stat.ME]* (2023). arXiv: [2302.02486 \[stat.ME\]](https://arxiv.org/abs/2302.02486).
- [7] C. F. Lo. “The Sum and Difference of Two Lognormal Random Variables”. In: *Journal of Applied Mathematics* 2012 (2012), pp. 1–13. ISSN: 1110-757X, 1687-0042. DOI: [10.1155/2012/838397](https://doi.org/10.1155/2012/838397).
- [8] Archil Gulisashvili and Peter Tankov. “Tail Behavior of Sums and Differences of Log-Normal Random Variables”. In: *Bernoulli* 22.1 (Feb. 2016), pp. 444–493. ISSN: 1350-7265. DOI: [10.3150/14-BEJ665](https://doi.org/10.3150/14-BEJ665). arXiv: [1309.3057 \[math, q-fin\]](https://arxiv.org/abs/1309.3057).
- [9] Geoffrey West. *Scale: The Universal Laws of Growth, Innovation, Sustainability, and the Pace of Life in Organisms, Cities, Economies, and Companies*. First Edition. New York: Penguin Press, May 2017. ISBN: 978-1-59420-558-3.

- 157 [10] R Gibrat. *Les Inégalités Économiques*. Paris: Librairie du Recueil Sirey, 1931.
- 158 [11] Robert Parham. “Facts of US Firm Scale and Growth 1970-2019: An Illustrated Guide”.
159 In: *arXiv:2302.02485 [q-fin.GN]* (2023). arXiv: [2302.02485 \[q-fin.GN\]](https://arxiv.org/abs/2302.02485).
- 160 [12] Thomas Malthus. *An Essay on the Principle of Population*. London: J. Johnson, 1798.
- 161 [13] Hernán D Rozenfeld et al. “The Area and Population of Cities: New Insights from a
162 Different Perspective on Cities”. In: *American Economic Review* 101.5 (Aug. 2011),
163 pp. 2205–2225. ISSN: 0002-8282. DOI: [10.1257/aer.101.5.2205](https://doi.org/10.1257/aer.101.5.2205).
- 164 [14] J. V. Henderson. “The Sizes and Types of Cities”. In: *The American Economic Review*
165 64.4 (1974), pp. 640–656. ISSN: 0002-8282.

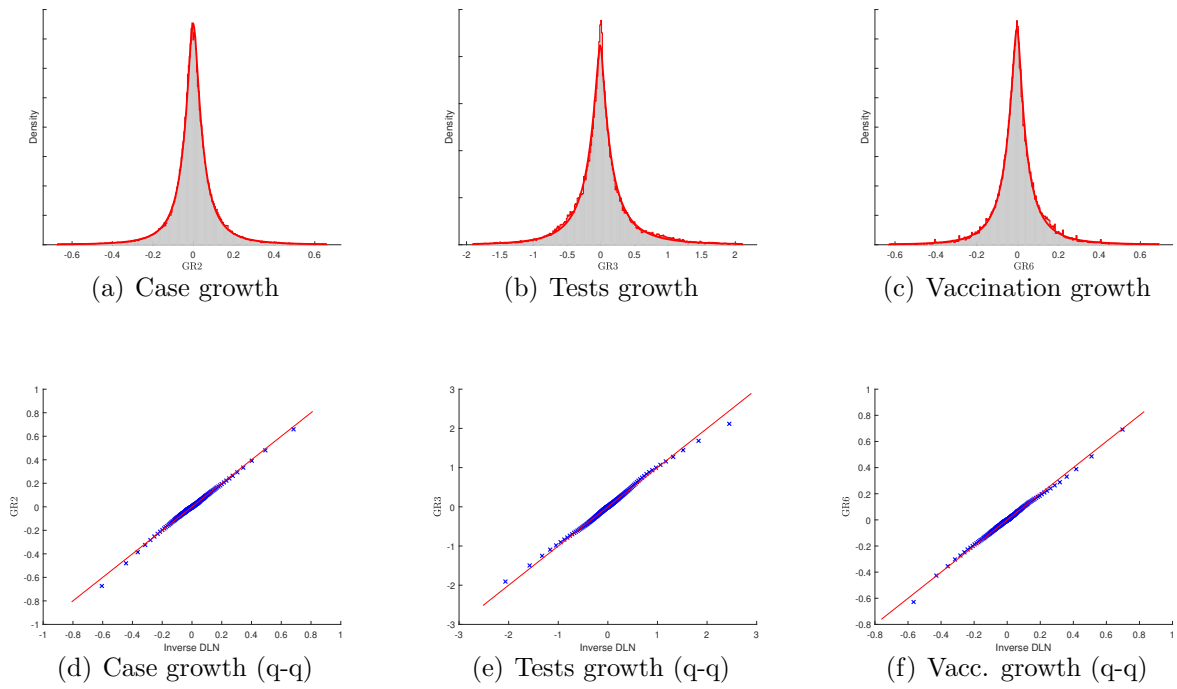


Fig. 1. Covid distributions stylized facts. Panels (a)-(c) present the distributions of daily new cases growth, new tests conducted growth, and new vaccinations growth, respectively. Panels (d)-(f) present the respective q-q plots.

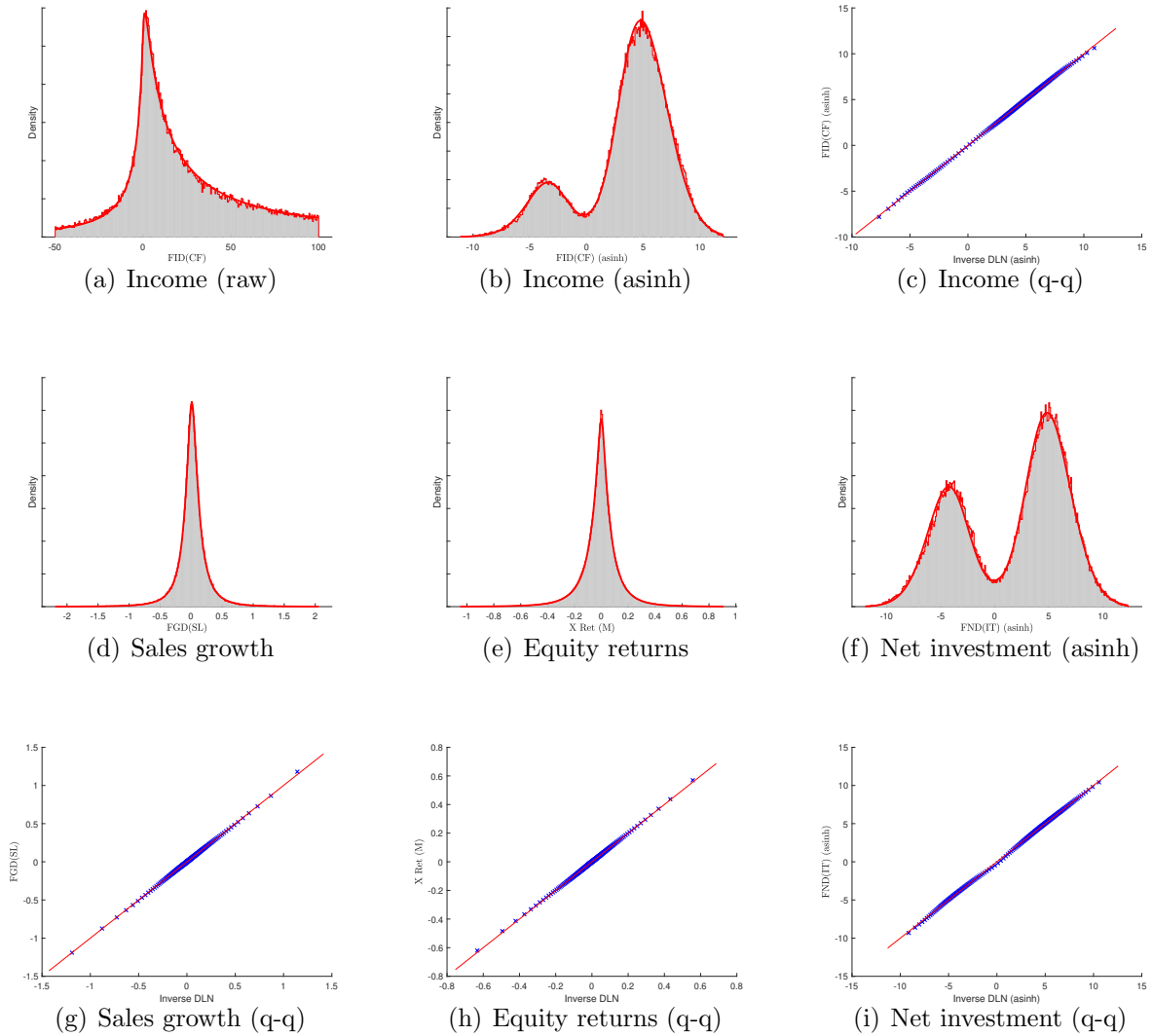


Fig. 2. Firm distributions stylized facts. Panels (a)-(c) present the distribution of income, with raw but truncated values in (a), asinh-transformed values in (b), and q-q plot vs. the DLN in (c). Panels (d) and (g) present the distribution and q-q vs. DLN for sales growth. Panels (e) and (h) repeat for Fama-French factor-adjusted equity returns at monthly frequency, Panels (f) and (g) repeat for asinh-transformed net investment.

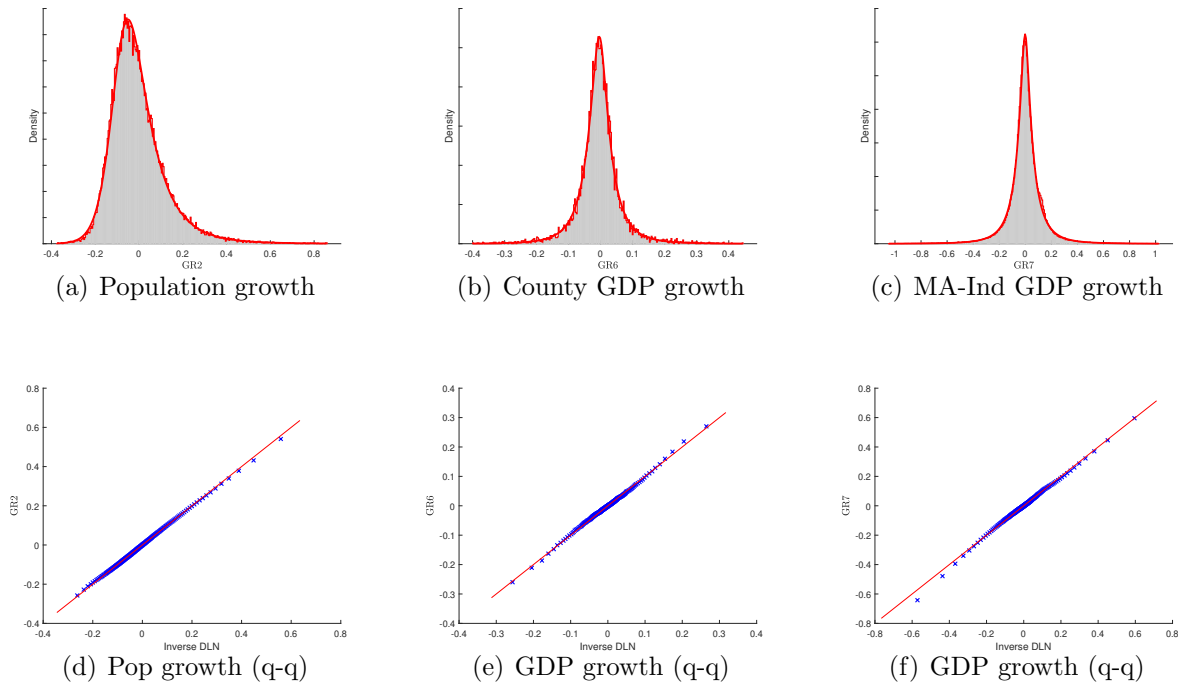


Fig. 3. City distributions stylized facts. Panels (a)-(c) present the distributions of population growth, GDP growth by county, and GDP growth by metropolitan area and industry, respectively. Panels (d)-(f) present the respective q-q plots.

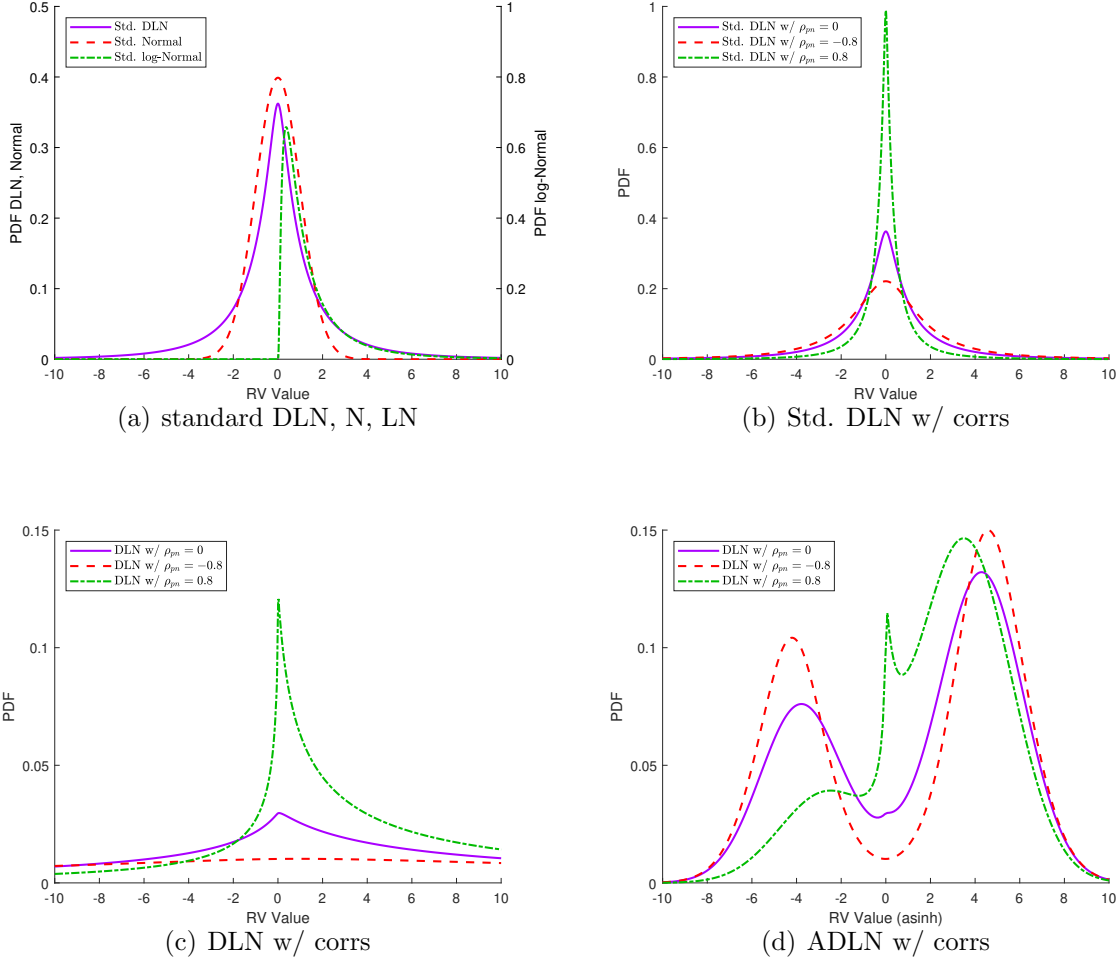


Fig. 4. DLN Examples. Panel (a) graphs the PDFs of the standard Normal, log-Normal, and DLN. Panel (b) graphs the PDFs of standard DLN with different correlation coefficients ρ_{pn} . Panel (c) presents the PDFs of a DLN with parameters $(3, 2, 2, 2)$, common in practice, and varying correlation coefficients ρ_{pn} . Panel (c) presents the PDF for the range ± 10 , which is a significant truncation due to the long tails of this DLN. Panel (d) presents the same PDFs as Panel (c), but the x-axis is asinh-transformed, such that it spans the range $\sinh(-10) \approx -11,000$ to $\sinh(10) \approx 11,000$.

166 **A Methods**

167 **A.1 Data**

168 COVID-19 data analyzed in Figure 1 are from Our World In Data (ourworldindata.org/coronavirus).
169 We use daily worldwide data, and plot all available growth observations when the base value
170 is higher than 10 (e.g., more than 10 infections per day or more than 10 vaccines given), and
171 the growth rate is different from 0 (as growth being exactly 0 usually indicates stale data).
172 The data were downloaded on 2/5/2022 and cover 143K observations for case growth, 71K
173 observations for tests growth, and 112K observations for vaccination growth.

174 Firm data analyzed in Figure 2 are from the Compustat/CRSP-combined dataset, ac-
175 cessed via Wharton’s WRDS. The data cover 164K firm-year observations on 15,797 firms
176 between 1970-2019. Sample selection criteria, exact variable definitions, and descriptive
177 statistics are reported in [11].

178 City data presented in 3 are from three sources. The data in Panels (a) and (d) are from
179 [13], and pertain to population growth from 1991 to 2000 in 46K locales identified by the
180 clustering algorithm of [13]. The data in Panels (b) and (e) are from the US Bureau for
181 Economic Analysis (Gross Domestic Product by County, 2017-2020), and pertain to 9,333
182 county-year observations on the growth of per-county GDP for 3,111 counties from 2017 to
183 2020. The data in Panels (c) and (f) are again from the US BEA (Gross Domestic Product
184 by Metropolitan Area and Industry, 2001-2017), and pertain to 270K growth observations
185 on 87 industries within 384 MAs over 17 years.

186 **A.2 Analysis**

187 For each growth measure, the data are first fit to the DLN distribution using the MLE
188 estimator described in [6]. The empirical distribution of the data, along with the fitted
189 DLN (in red) are presented as figures, along with a q-q plot of the empirical CDF vs. the
190 theoretical DLN CDF, also developed in [6].

191 Next, three statistical goodness-of-fit tests are used to verify whether the empirical data
192 indeed stem from the DLN distribution. The three tests used are the Kolmogorov-Smirnov
193 test, the Anderson-Darling test, and the Chi-square test. Details on conducting these tests
194 and on constructing the p-values for the tests are available in [6] and [11]. The tests generally
195 do not reject the DLN for the growth data at the 5% confidence level. This is in contrast with
196 the two other candidate distributions discussed in the literature: the Laplace distribution
197 (itself a difference of exponentially distributed variates), and the Lévy-Stable (aka Pareto-
198 Stable or Power-Law) distribution. Both of these distributions are generally rejected at the
199 5% level by the goodness-of-fit tests.

200 Finally, for each growth measure, we conduct likelihood-based information criteria “horse-
201 race” tests between the DLN and the two other candidate distributions. The tests are based
202 on Akaike’s Information Criterion (AIC) and the Bayesian Information Criterion (BIC). In
203 all such tests, the DLN is overwhelmingly favored over the other candidates.

204 **B Data Availability**

205 COVID data are freely and publicly available from Our World In Data (ourworldin-
206 data.org/coronavirus).

207 Firm data are publicly available from the CRSP/Compustat Merged Database, which
208 is subscription-fee-based (crsp.org). The firm data were normalized by the year’s nominal
209 GDP from the St. Louis Fed (fred.stlouisfed.org).

210 City data are freely and publicly available from two sources: the US Bureau for Economic
211 Analysis (bea.gov) and a replication package for [13] from the AER’s website (aeaweb.org).

212 **C Code Availability**

213 Code to reproduce all figures reported in this analysis, including code implementing the
214 CDF, PDF, and DLN parameter estimation is publicly available from the author and will

215 be provided to Nature.