

Averages don't characterise the heavy tails of ransoms

1st Éireann Leverett
*Founder of
Concinnity Risks*
Cambridge, United Kingdom
eleverett[AT]concinnity-risks.com

2nd Eric Jardine
Assistant Professor
Department of Political Science, Virginia Tech
Blacksburg, United States
ejardine[AT]vt[DOT]edu

3rd Erin Burns
*Founder of
Concinnity Risks*
Cambridge, United Kingdom
eburns[AT]concinnity-risks.com

4th Ankit Gangwal
University of Padua, Italy
ankit.gangwal[AT]phd.unipd.it

5th Dan Geer
Senior Fellow In-Q-Tel
dan[AT]geer[DOT]org

Abstract—The Bitcoin block-chain is the scoreboard of Ransomware. By mining the data in it and within the malware itself, we can understand the distribution of ransoms and characterise ransomware risk. Ransoms follow the power-law distribution in their amounts. The alpha parameter (α) of those power laws suggest they do not have a well defined average for most years in our study. Indeed, there has not been an α above 2 since 2015 and so there has not been a stable ransomware average since that time. The changing α has strong implications for cyber risk management and policy responses to ransomware attacks.

Index Terms—Ransomware, ransoms, block-chain, Bitcoin, malware, ecrime, forensic-accounting, power-laws, security economics, cyber risk, black swans

I. INTRODUCTION

Cyber insurance is growing rapidly, and needs risk models to continue to efficiently match capital to risk. Ransomware is a particular problem within the space, and one which so far is often characterised by discussion of average ransoms and average losses^{1,2,3}. Average ransoms are a common metric used in both the security industry and the academic literature to talk about ransomware [1–6]. This paper amply demonstrates that averages are not stable nor a true measure of central tendency when the underlying distribution of events follows a power law with an α parameter below 2. To progress risk management of ransomware, we need better models that fit the underlying distribution of ransoms. Using Kolmogorov-Smirnov statistics to test for goodness of fit, we show that ransoms from 2013 to 2019 follow a power law pattern in most years and so we advocate that averages ought to be abolished as a characterisation of outcomes in security reports and actuarial tables of cyber risk. Maximum likelihood distribution fits may find ransoms follow another distribution in the future when more ransoms are available, our goal here is a softer

one than a precise distribution fitting; it is to simply show they have heavy tails.

Power laws and heavy tails in cybersecurity are nothing new and show tremendous potential to better model cybersecurity risk [9]. Others have clearly demonstrated that data breaches follow power laws [10], unpacked how heavy-tailed cybersecurity events affect the collection of data through cybercrime surveys [11]), used power law distributions within a multi-period Monte Carlo simulation to test for Simpson's Paradox in aggregate cybersecurity statistics [12], and tend to be relevant for estimating cyber losses from insurance data [13]. Other heavy-tailed distributions such as lognormal have also been found to apply to things such as phishing site lifetimes [14] and the time to recovery of IT assets [15]. They have even been used to counter the argument that powerlaws fit data breaches [16]. Yet previous studies have not attempted to fit power law distributions to modern ransomware, and we believe our work to be novel in that respect⁴.

Hubbard and Seiersen wrote that risk quantification is essentially just the iterative reduction of uncertainty [17]. The focus here is to reduce uncertainty through measurement by quantifying the year-on-year severity of ransoms and demonstrating the potential fit of a power law distribution as a tool towards uncertainty reduction. Power laws do not fit the data in all years in our sample, but do for most and overall. This on-again-off-again statistical significance is consistent with other work fitting power laws to wealth distributions [18]. Better fitted metrics for ransoms can significantly improve both cyber insurance and global crime management policy. The presence of power laws also imply that ransom prices can exhibit black swan events and present an accumulation and aggregation risk. More importantly the use of averages can lead to inaccurate estimations of costs if data is collected via cybercrime surveys rather than directly from the Bitcoin blockchain [11]. To put it succinctly, a survey of income will have two very different

¹<https://www.coveware.com/blog/q1-2020-ransomware-marketplace-report>

²<https://kivuconsulting.com/ransomware-you-can-prevent-it-but-you-cant-solve-it/>

³<https://www.businesswire.com/news/home/20170515006833/en/Symantec-Blocks-22-Million-Attempted-WannaCry-Ransomware>

⁴The blockchain didn't exist at the time that some of these previous papers were authored.

averages if it includes or doesn't include the top 1 percent of wealth and ransomware payments are largely the same.

II. DATA AND METHODS

The data studied in this paper include over 5.4 million payments to over 26,000 BTC addresses used in malware or ransom notes between 2011 to 2019. It was assembled by finding Bitcoin addresses within ransomware binaries or ransom notes, and enriched by finding copayments from those addresses to other addresses (Simplistically, examining bundles of addresses while ransomware gangs "cash out"). This allows us to find payments associated with ransoms, that might normally be hidden from our collected malware samples, and was done in a manner consistent with one author's previous work on copayments [7]. It is possible that some of those payments are for other things than ransoms, but we do not believe it invalidates our points about heavy tails and averages in this paper. It is one of the largest available datasets on ransomware ransom payments that we know of, which also preserves the forensic connections between binaries and cryptocurrency addresses.

The majority of the data used in this report is already in the public domain within the block-chain, and available from many fine APIs. By using the code we have written, the interested reader can gather bitcoin addresses from any ransomware data set at their disposal. Ours came from the wonderful data at⁵ and BTC Abuse⁶. For obvious reasons we cannot share our malware corpus, but some simple vetting can be done with the administrators of VirusShare to gain access to it yourself. We also believe the results in this paper could be reproduced from other malware databases such as the Malware Bazaar⁷ or PolySwarm⁸. Once you have gathered enough ransomware addresses, you begin by gathering all transactions to those addresses, normalise the value in USD per diem as a BTC daily spot price and then apply the Powerlaw Python library to analysis of such data USD normalised data. [8].

Despite the scope of our malware data, it is important to acknowledge the limitations of this method. Not all ransoms paid as a result of ransomware attacks can be included. First, we only focus on ransoms paid in Bitcoin. The graphical inlay clarifies the data collection process see Fig. 1. Secondly, if an address is not included in a binary or ransomnote and only shows up in negotiation, we are less likely to have found it. We use other methodologies to collect those addresses, but we acknowledge those methodologies as ad hoc. Despite some obvious limitations, we do not see any reason why this collection methodology, particularly when combined with the scale of the data, would be systematically biased in a way that would effect out estimation of the fit of power law distributions. When we extrapolate about other things, we will be careful to caveat what we believe are blinds spots or biases of the data.

⁵<https://virusshare.com/>

⁶<https://www.bitcoinabuse.com/>

⁷<https://bazaar.abuse.ch/>

⁸<https://polyswarm.io/>

No doubt malware enthusiasts will be primarily focused on some metadata of our malware corpus rather than the economics of payments. A table can be found at the end of the Appendix that should satisfy the thirst of those readers. The table that shows the most common filetypes (above 10000 samples of that filetype) that we processed. We choose to publish this slightly abbreviated list as the full table of file types would be more than 1000 lines, with many, many files that are unique in many ways. Also please note that these are not all ransomware files, but rather malware or associated with malware, that we sifted through to find ransomware. Once ransomware was identified we turned our efforts to find ransomware Bitcoin addresses.

It should also be noted that the number of ransoms per year in the data drops by 4 orders of magnitude across time. This decline is probably *not* strong evidence that the occurrence rate of ransomware has necessarily dropped, but rather a side effect of our collection methodology. Putting it simply, if there is not cryptocurrency address in a malicious binary, then we will not discover it. E.g., if the ransom note only has an email address (an increasingly common practice), then we will not be able to discover a BTC address and record the payment of a ransom. Criminal adaptation, in other words, might be affecting the coverage of the data over time. But it is not clear a priori whether changing tactics would simply reduce the volume of recorded attacks at random or do so in a systematic way.

A. Ransomcoin

A variety of tools help us produce the data we analysed in the paper, and can be reused by others to reproduce the work, or indeed challenge its findings in the future.

1) *Code*: The code we used to gather the BTC accounts from ransomware samples is called RansomCoin and is stored and maintained here⁹. We provide the specific commit used for this paper, for reproducibility reasons. Though readers in the future may prefer to use a more up to date version of both the data and the code depending on their goals.

Most of the code used in this project is very simplistic, and specific to the formatting our data set. In the two images below though, we display the general code that would be useful to anyone to explore ransoms or ransomware losses with the Python3 Powerlaw library, for a more generalised dataset than our own.

III. RESULTS

One of the defining features of a power law distribution is that the mean and median diverge. In a normal distribution, descriptive measures such as the mean, median and mode point to the centre. When power laws are at play, mean averages tend to be wildly inflated by the outliers in the heavy tails. Data breaches exhibit such a pattern, with a "A factor 4.5× difference between mean and median [which] is indicative of greater concentration than even the US wealth distribution" [10].

⁹<https://github.com/Concinnity-Risks/RansomCoinPublic/commit/bd554ba41c55e074f79069fe4e14f4762bb71228>



Fig. 1. A quick description of how we collect and process our data on ransoms.

At a basic descriptive level, the ratio of mean to median in our ransomware ransom data set is an excruciating 31.07 (cf. TABLE 1). These descriptive results strongly suggest that heavy tail distributions might best fit the underlying data and, more parenthetically, that using averages [1? –6] to discuss and price ransom risk will both lead to routine over payment and fall prey to the risk of astronomically large events found in the heavy tail of the ransom distribution. To formalise the distributional fit of the data, we employ a series of Kolmogorov-Smirnoff tests and log/log plots.

A. Power-laws apply to Ransoms

Power laws apply in many places and are particularly common in questions of income and wealth [19]. Indeed, the

TABLE I
MATHEMATICAL CHARACTERISTICS OF THE RANSOM DATA IN THIS PAPER.

Parameter	USD
Variance	49,633,044.87
Mean	800.58
Median	25.77
Ratio	31.06

Pareto Principle (i.e. 80 percent of the outcome is produced by just 20 percent of the inputs) was coined in the realm of wealth inequality. Power laws also apply to the concentration of users on websites and other online services, the topography of incoming links, the distribution of academic citations, the magnitude of earthquakes and solar flares, the intensity of wars, and the sizes of cities [20–23].

Power laws likely apply to ransoms not necessarily because of anything to do with malware design or the propagation of infections across networks. Instead, ransoms likely follow a power law distribution because corporate and civic wealth does [19]. Framed differently, if we assume that malicious actors are targeting their ransomware with knowledge of their target’s general size (e.g. a big company versus small or a large city of millions versus a tiny town of a few thousand), then they ought to calibrate their ransom demands based upon their target’s expected ability to pay (willingness to pay in the economics parlance). Since that ability to pay is a function of revenue that is distributed in a power law, then ransoms, it stands to reason, are a sampling themselves of that power-law distributed wealth. Additional measurement going forward can parse the fit of other long tail distributions such as log-normal or parabolic fractals, which may indeed turn out to be the most likely fits. Regardless, we will have succeeded in our goal to have abolished the use of averages by that time. Why shouldn’t averages be used in our ransomware studies?

Formally, power-law distributions are:

$$P(x) = CX^{-\alpha} \quad (1)$$

The α parameter of the power law determines the stability of the average. For example, if:

$$1 < \alpha < 2 \quad (2)$$

Then the mean and average are infinite¹⁰. Also notably, even when:

$$2 < \alpha < 3 \quad (3)$$

There is a finite mean, yet there is still not a well defined standard deviation. This complicates the ability to predict ransoms (and as a follow on, losses if they follow the same distribution). α parameters drifting into or out of this territory could move ransomware from insurable via standard methods

¹⁰http://tuvalu.santafe.edu/~aaronc/courses/7000/csci7000-001_2011_L2.pdf

to insurable via catastrophe bonds. These are literally the "fat tails" of Kousky and Cooke in "Explaining the Failure to Insure Catastrophic Risks" [25].

Quoting them:

"Risks that people colloquially term "catastrophic" are usually characterised by fat tails and dependence. With fat-tailed loss distributions, the probability of an event declines slowly, relative to its severity. Simply, very large losses are possible. The precise mathematical definition of fat tails is rather subtle, but a working notion is that damage variable X has a fat tail if, for sufficiently large values x , the probability that X exceeds x is kx^{-a} , for some constants $a, k > 0$. The variable a is referred to as the tail index and it roughly governs how fat the tail of the distribution will be. Many natural catastrophes, from earthquakes to wildfires, have been shown to be fat tailed."

The distribution fitting was done using the Python Powerlaw package [8]. We chose to plot ransoms above \$10,000 in the tables and graphs because cyber insurance is not interested in lower limit ransoms, as they are unlikely to ever turn into claims. Pooling the data from all years and plotting on a log/log plot with a fitted Kolmogorov-Smirnov power law trend line is indicative Fig. 2. Interpretation of a KS one sample statistic works via rejection of the null, where the null hypothesis is that the sample is drawn from the specified distribution, in our case a power law distribution. We suspect that there might be two populations of ransoms at play, one that applies to an individual who is asked to pay smaller ransom values and where parametric statistics might apply and another where organisations asked to pay higher dollar values that are subject to fat tails. One friendly economist-turned-malware-analyst suggested this might show the price-takers at the bottom and the price-setters at the top. Dropping the economics jargon: only the big companies can afford ransom negotiators who reduce the ransom substantially. So even within 'organisational ransoms' there may be one distribution of un-negotiated ransoms and another of negotiated ransoms.

This muddies the waters of distribution fitting, because we may be dealing with more than one underlying distribution and one may have averages while another does not. Or both may, or neither. Which is why we make it clear that seeing the divergence of the median from the mean suggests even ransoms have fat tails, regardless of our ability to disentangle the underlying distributions. Clearly further work in distribution fitting can be done as more data and losses emerge. Future papers by better statisticians may be able to find the signal in the noise, and we encourage them, while also suggesting we stop using averages to discuss ransoms because they clearly over inflate median severity.

Of course, the pooled sample showing 2013 ransomware alongside 2019 ransomware events could mask underlying temporal variation in the data. It is widely reported that ransomware is "getting worse", though this is not usually articulated with mathematical or actuarial clarity. Yet, as can

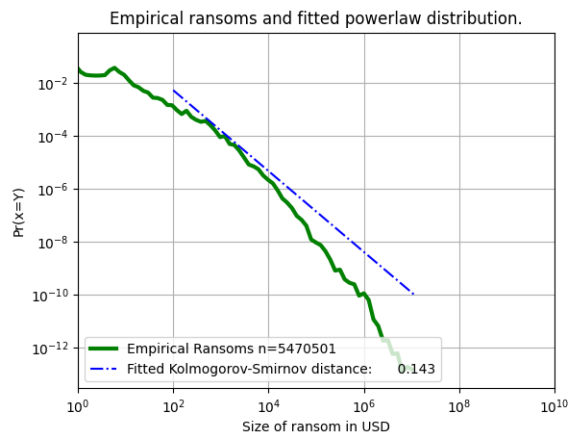


Fig. 2. This figure captures all ransom and copay data analysed, and shows the power-law fit

be seen in the data below that dis-aggregates the results into yearly estimation, the frequency of ransoms was much higher before 2016, though the severity is much higher afterwards. We believe this to be broadly true. This finding somewhat counters the narrative that security companies and insurance companies are pushing: That ransomware is occurring with ever growing frequency. It also highlights a fascinating phase transition from mass ransomware with lower ransoms to targeted corporate ransomware with much higher impacts and associated ransom payments.

B. Yearly analysis of ransoms

In 1989 Joseph L. Popp Jr. invented the first ransomware virus spread on floppy disks¹¹. Unfortunately, we don't have access to the data on how many people paid the PC Cyborg corporation in Panama, though we do know the exposure rate (20,000 diskettes distributed at a WHO AIDS event) and ransom demanded (\$189 USD). Notice the fixed price ransom here, which became variable pricing over time.

However, the game really changed for attackers and defenders with the invention of Bitcoin in 2008, and the realisation of a public ledger on 3 January 2009 [26]. This allowed us an unparalleled view into the transactions of ransomware from 2013 onwards. It also offered exceptional scalability to ransomware operations. The ability to perform all the necessary steps in an automated way online, increased their opportunities in both scale and maintaining anonymity. This led to a wide proliferation of victims paying small amounts. A notable strain of 2013 Ransomware was Cryptolocker, which spread via emails, file sharing, and downloads.

For all the years in our data, we show our other parameters (both set and derived) the TABLE II. Below, we provide 12 month rolling windows of analysis for each of years in the data, which we also scale to be the same size on the advice of a trusted friend. Fitting by year allows us to track changes to the α parameter that might be germane to both the stability

¹¹<https://www.knowbe4.com/aids-trojan>

of ransom averages and the management of cyber risk. The tables speak to that point in the Data and Code Appendix.

TABLE II
YEARLY α AND σ VALUES OF THE FITTED POWER LAWS.

Year	xmin	xmax	α	σ	Observed Ransoms
2013	10000.0	None	2.35	0.010	684,813
2014	10000.0	None	2.42	0.005	4,137,194
2015	10000.0	None	2.16	0.065	155,584
2016	10000.0	None	1.78	0.013	444,030
2017	10000.0	None	1.50	0.012	20,987
2018	10000.0	None	1.77	0.042	20,983
2019	10000.0	None	1.98	0.135	1,980

With that bit of yearly context, let's look at the severity of ransoms paid per year. In each figure, we plot the a Kolmogorov-Smirnov statistic and a baseline threshold value with an α of 2 under neither which traditional instruments of insurability fail to apply. For the data from 2013, as can seen in Fig. 3, the Kolmogorov-Smirnov distance of a power-law distribution fit is both statistically significant, suggesting that the data do not follow a power law, and shows an α of 2.352. The plotted power-law distribution with an α of 2.00 for reference is displayed in the rest of the graphs that progress year by year. This visual aid continues throughout the graphs in the paper to help reflect on questions of averages and variance, and identify the threshold with which their properties shift in power laws.

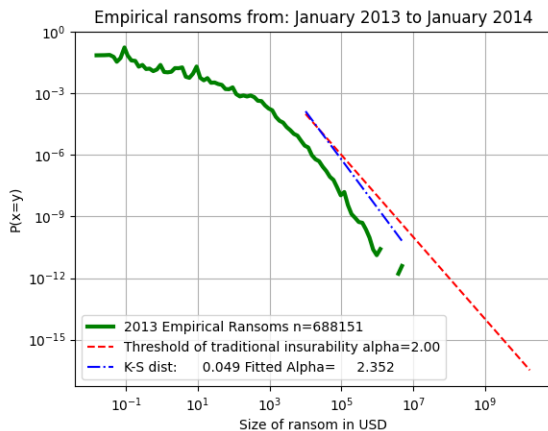


Fig. 3. This figure captures 2013 ransom data, and shows the power-law fit and resulting alpha parameter.

2014 saw the development of Cryptowall I and II, as well as CTB-Locker, and CryptoTorLocker2015. In Fig. 4, power-laws do not fit the data in this year, and the α parameter is similar to that of 2013, notably also similar to studies of wealth distribution over time where the fit varies by time frame and sample. Note particularly Brzezinski's "Do wealth distributions follow power laws? Evidence from 'rich lists'": Figure one of that paper shows the global variance of the α is 3.3 - 2.2 [18].

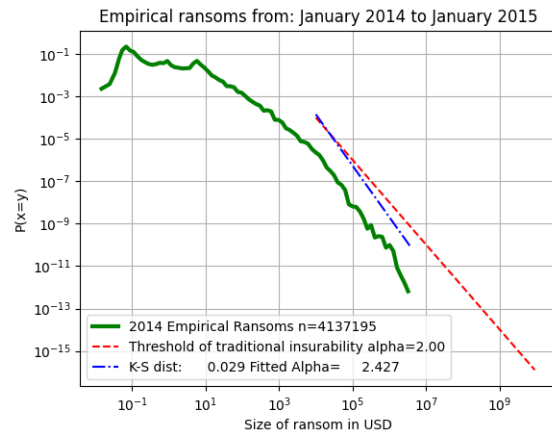


Fig. 4. This figure captures 2014 ransom data, and shows the power law fit and resulting alpha parameter.

In 2015 the α reverses direction and begins approaching 2, suggesting greater mean/median divergence. The KS statistic still suggests that the sample is not drawn from a power law distribution, but averages in this period are trending toward a worse measure. This complicates things, both in modelling terms, but also insurability. The alpha in this period is testing some of the limits of traditional risk mitigation and management. We believe the alpha of losses might be even lower than that of ransoms, though. This follows logically, since when a ransom is paid one presumes it was cheaper to do so than to handle the losses (the principle that victims who pay are rational). The losses contain many more categories of expense which mostly are proportional to organisation size I.E. reputational harm, business interruption and downtime, regulatory fines, future audits, and forensic investigations.

This drift towards an α of 2 can be seen in Fig. 5. We hope that insurance companies aren't building strategies where simply paying the ransom is the approach, because that could lead to severe accumulation risks in a world of powerlaw distributed ransoms¹². It's worth paying special attention to who your incident responders are¹³.

In Fig. 6, we can see that 2016 ransoms have an α of lower than 2, and also supports the statistical fit of our power-laws.

However, the trend in the data through 2015, is in many senses not as troubling as the α parameter we find in 2016. Here, the KS statistic fails to reject the null, suggesting that the data do follow a power law distribution and the α drops below 2, making averages inherently unstable.

In fact we prefer to interpret it as shifting the burden of ransomware insurability onto those who wish to insure it. Let those who can disprove this paper with one of their own build the model of ransomware insurance accumulation. They will no doubt earn their salary many times over by doing so.

¹²<https://www.propublica.org/article/the-extortion-economy-how-insurance-companies-are-fueling-a-rise-in-ransomware-attacks>

¹³<https://features.propublica.org/ransomware/ransomware-attack-data-recovery-firms-paying-hackers/>

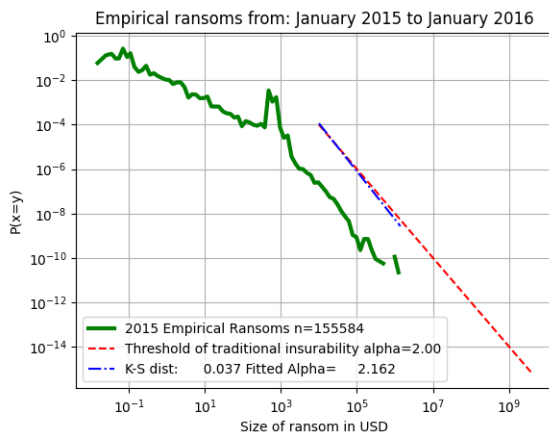


Fig. 5. This figure captures 2015 ransom data, and shows the power-law fit and resulting alpha parameter.

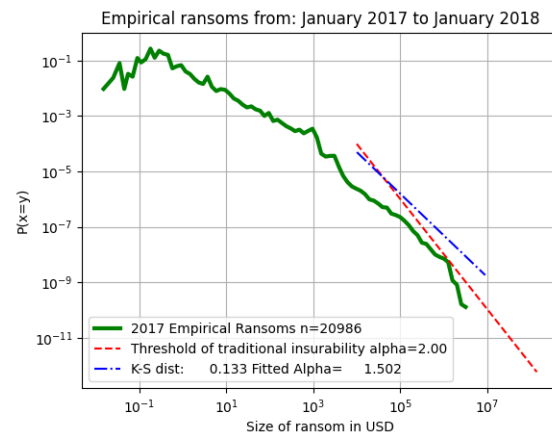


Fig. 7. This figure captures 2017 ransom data, and shows the power-law fit and resulting alpha parameter.

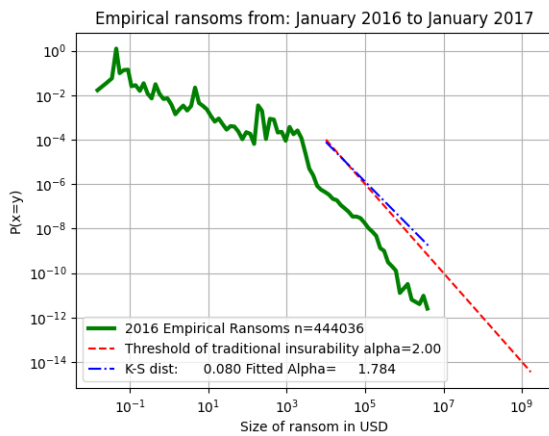


Fig. 6. This figure captures 2016 ransom data, and shows the power-law fit and resulting alpha parameter.

Returning to the data though, perhaps a log-normal or truncated power-law would still fit here. The α lower than 2 though shows us that we are dealing with very heavy tails.

The total costs of historical wars is considered to be uninsurable, and many powerlaws are present in the analysis with an α lower than 2 [31].

In Fig. 7, we see this trend of α reduction continuing. Indeed, this is the minimum α we will witness within the data set we have studied. A historical low in the frequency of ransomware events that we can detect, but a notably high in terms of severity. This also coincides with the FBI stating that there were over 4,000 ransomware attacks per day that year¹⁴. That provides us a rare opportunity to check our number of ransoms for 2016 (444030/365 gives us 1216.5), suggesting our data set is only seeing a little over a quarter of 2016 occurrences of ransomware (though perhaps no ransom was paid in those cases).

¹⁴<https://www.fbi.gov/file-repository/ransomware-prevention-and-response-for-cisoc.pdf/view>

In Fig. 8, the power law distribution still fits the data, but we see a return towards the threshold of alpha to where mathematical conditions change. We wonder if such a trend will continue, or it will turn back again after some new innovation in ransomware policy. Indeed a tantalising and intriguing intellectual question is: Could cyber policy interventions change the alpha of ransoms in some way?

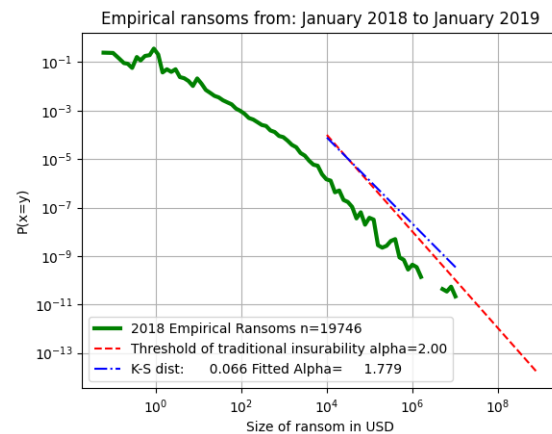


Fig. 8. This figure captures 2018 ransom data, and shows the power-law fit and resulting alpha parameter.

Consider then, if you were simply paying ransoms through insurance since 2016, this would be a strategy doomed to make ransoms too unpredictable for actuarial tables and traditional insurance methods. Ideally no insurance company would consider this strategy simply on the grounds that it creates many perverse incentives and feeds the problem. Secondly it is a legal risk and moral hazard that should be avoided. However, if the moral and ethical arguments fail to motivate, then perhaps we should return to the idea that as they increase in severity, we push the α into places where traditional insurance risk management fails.

Fig. 9 is included for completeness, but contains less ransoms than we believe occurred. This is discussed extensively in the next section.

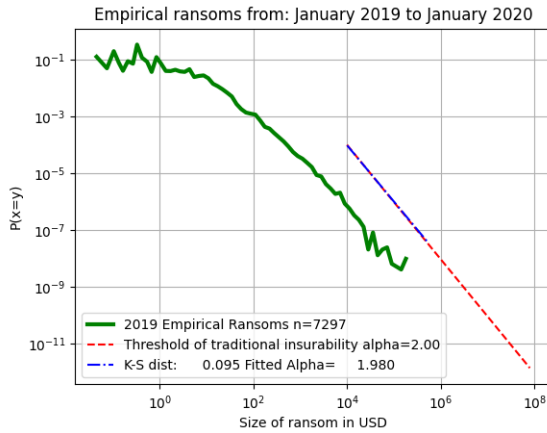


Fig. 9. This figure captures 2019 ransom data, and shows the power-law fit and resulting alpha parameter.

For this final year, we have not created a table of rolling windows as 2020 is not finished at the time of writing, and also the small number of ransoms from which to work make the data quality issues too great to be included. It is left to future authors to examine if these trends continue, and postulate the causation of power-laws in ransomware ransom severity.

IV. DISCUSSION

Let us be succinct here: our main premise throughout this paper is that the α of power-laws fitted to ransoms is *not* related to technical elements of the malware such as the spreading factor of infections, or number of systems vulnerable to a particular CVE. Instead, the α of ransoms reflects the α of financial power-laws located throughout society, such as the α of 2.32 for American personal wealth in documented in [11] or [27]. Thus the α in ransoms reflects a malicious actor’s assessment of their target’s ability to pay. Like in the market for malware more generally, Sutton’s Law prevails¹⁵ and the financially wealthy, be they big firms, large cities, or giant organisations, seem to be reaping the negative outcomes of their concentrated wealth [28].

The empirical findings suggest that ransoms fit a power law distribution in many of the sample years and the underlying causal explorations suggest several relevant points. First, it is likely that all-in losses are even more scale free than direct ransom costs. Ransoms will compose some fraction of costs an enterprise needs to pay, but do not directly include remediation and recovery costs, loss of services, nor any potential intangible costs associated with a tarnished reputation. Generally, though, ransom costs will be less than all-in costs, except in the event that sufficient backups are at hand to eliminate

¹⁵https://en.wikipedia.org/wiki/Sutton%27s_law

the risk without paying. If ransoms alone are pushing the boundaries of risk management through averages, then all-in losses with a lower alpha would be too. Essentially, for the firms, organisations and entities hit with larger ransom demands where the lower α suggests that averages are highly unstable and standard deviations will be misleading.

Second, there are potentially two or more populations of ransoms with different underlying distributions at play. Power laws appear to apply in most years, so averages will break down and potential losses become scale free. But it is possible that for some subset of the ransom payments, typical parametric statistics likely still apply. To the extent that the causal reasoning that ransoms follow a power law because wealth/resources do the same is correct, at least two populations might be at play. One population in the lower ransoms of the distribution likely includes small organisations and individuals, those whose wealth is such that they will only be targeted with a ransom valued at a few hundred dollars. Here, speaking of average ransoms likely captures a somewhat stable measure of central tendency, with little price targeting by ransomware gangs. The second population includes larger firms, governments and organisations. Ransom risk for this population might be non-parametric in most years and is more likely to be above \$10,000.

This second population where averages fail to apply is of especial interest to insurers. Bigger organisations are the most likely to actually purchase insurance for cyber risk, but they may be the least likely to have ransom costs and all-in financial losses due to a security event that can be modelled well with tradition parametric statistics. Reinsurers and cyber policy advocates should be worried about the low α of ransoms for large firms and examine how they might shift it back into a stabler parametric statistical range.

To that end, the cyber-insurance industry may be able to “push” the α of ransoms higher through a variety of mechanisms. A government back stop¹⁶ such as the one provided to poolRE¹⁷, is one option. Another is the construction of Catastrophe Bonds [25]¹⁸. Indeed it might be possible to use these power laws as a parametric risk trigger for ransomware in the future. Risk Pooling¹⁹ is another option and has been used to tackle other reinsurance aggregation and accumulation problems. Some novel approaches within the technical domain might be the creation of bounties for ransomware decryptors and other technical measures, or the construction of a dedicated cyber insurance CERT, perhaps even a non-profit or community interest scheme dedicated to reducing ransomware losses through a variety of methods.

¹⁶<https://www.brookings.edu/blog/techtank/2019/09/27/a-federal-backstop-for-insuring-against-cyberattacks/>

¹⁷<https://www.globalreinsurance.com/a-tale-of-two-systems-terrorism-reinsurance-backstops-in-the-us-and-uk/1428134.article>

¹⁸https://en.wikipedia.org/wiki/Catastrophe_bond

¹⁹<https://www.worldbank.org/en/news/feature/2017/11/14/what-makes-catastrophe-risk-pools-work>

V. CONCLUSIONS

The blockchain is the scoreboard for ransomware and CERT teams everywhere. As such it has been largely ignored for ransom risk analysis, focusing almost exclusively on identifying new strains of ransomware or new addresses being used for crime. We believe this paper is one of the first to identify that the blockchain is also enormously useful for constructing risk and actuarial tables for ransomware. This might benefit everyone, for example if an open source ransomware risk model was built like those in Oasis²⁰, then any business, organization or government could quickly estimate their ransomware risk and then move forward with mitigation, restitutions, or transfer mechanisms. Such an open source collaboration could also benefit insurers, reinsurers, Computer Emergency Response Teams (CERTs), and anti-ransomware projects such as the wonderful NOMORERANSOMS project alike²¹. Continual measurement is also key, as attackers change their way of doing things in directions that could affect the underlying distributional type.

Ransoms can probably be split into a population of small payments and comparatively estimable risk and a second population that might be pushing the boundaries of risk management or insurability by averages. This might mean we need to energeise communities such as global CERTs²², insurers, reinsurers, digital forensics companies, and even governments to collaborate on methods that reduce the α parameter of losses and ransoms. What these methods may achieve remains open to discussion, but at least we have a methodology year on year to see how we are progressing. Only when reinsurers insist that insurers begin the process of **active risk management** as well as passive risk transfer, will we see a manageable risk of ransomware in the global economy.

VI. CONTRIBUTIONS

Éireann Leverett conceived of the RansomCoin software with Jurriann Bremer in a hotel bar in Poland back in 2017, but went on to do the coding required to gather the data over 3 years. He applied the distribution fitting of power-laws to ransoms over the last year. Eric Jardine applied quantification rigour around data biases, offered deep insights into cyber risk metrics and cyber policy interventions. Erin Burns helps maintain the diverse code bases that produce this dataset, visualise and analyse the results, and advises us on financial forensics and fraud, generally. Ankit Gangwal contributed his share of ransomware addresses, and taught us how to classify co-payments. Dan Geer lent us his prodigious knowledge of past publications, and provided thoughtful reviews of this paper while bringing it to publication.

“The data curation, data science and writing for this project was supported by ‘Incrementally Tailoring a Better Cyber Risk Score’ funded by a Comcast Innovation Grant. Grant Number 2019-145”

²⁰<https://oasislmf.org/>

²¹<https://www.nomoreransom.org/>

²²<https://www.first.org/global/signs/cyberinsurance/>

The authors wish to thank their loved ones for the time they spent away working on data and maths. Éireann wants to thank the curiosity of the penguins for inspiring him. You know who you are.

“The authors declare no conflict of interest.”

REFERENCES

- [1] Paquet-Clouston, Masarah, Bernhard Haslhofer, and Benoit Dupont. “Ransomware payments in the bitcoin ecosystem.” *Journal of Cybersecurity* 5, no. 1 (2019).
- [2] Alhawi, Omar MK, James Baldwin, and Ali Dehghan-tanha. “Leveraging machine learning techniques for windows ransomware network traffic detection.” In *Cyber Threat Intelligence*, pp. 93-106. Springer, Cham, (2018).
- [3] Sgandurra, Daniele, Luis Muñoz-González, Rabih Mohsen, and Emil C. Lupu. “Automated dynamic analysis of ransomware: Benefits, limitations and use for detection.” arXiv preprint arXiv:1609.03020 (2016).
- [4] Cabaj, Krzysztof, and Wojciech Mazurczyk. “Using software-defined networking for ransomware mitigation: the case of cryptowall.” *IEEE Network* 30, no. 6 (2016): 14-20.
- [5] Everett, Cath. “Ransomware: to pay or not to pay?.” *Computer Fraud Security* 2016, no. 4 (2016): 8-12.
- [6] Nadir, Ibrahim, and Taimur Bakhshi. “Contemporary cybercrime: A taxonomy of ransomware threats mitigation techniques.” In *2018 International Conference on Computing, Mathematics and Engineering Technologies (iCoMET)*, pp. 1-7. IEEE, (2018).
- [7] Conti, Mauro, Ankit Gangwal, and Sushmita Ruj. “On the economic significance of ransomware campaigns: A Bitcoin transactions perspective.” *Computers Security* 79 (2018): 162-189.
- [8] Jeff Alstott, Ed Bullmore, Dietmar Plenz. *powerlaw: a Python package for analysis of heavy-tailed distributions*. ‘PLoS ONE 9’ (2014).
- [9] Stefan Prandl and Mihai Lazarescu and Subhash Kak. “PEIMA: Harnessing Power Laws to Detect Malicious Activities from Denial of Service to Intrusion Detection Traffic Analysis and Beyond”, (2017).
- [10] Maillart, Thomas, and Didier Sornette. “Heavy-tailed distribution of cyber-risks.” *The European Physical Journal B* 75, (2010).
- [11] Florêncio, Dinei and Herley, Cormac. “Sex, lies and cyber-crime surveys”, *Economics of information security and privacy III*, Springer, (2013).
- [12] Jardine, Eric. “Sometimes three rights really do make a wrong: Measuring cybersecurity and Simpson’s paradox.” In *Workshop on the Economics of Information Security*, La Jolla, CA. (2017).
- [13] Woods, Daniel, Tyler Moore, and A. Simpson. “The County Fair Cyber Loss Distribution: Drawing Inferences from Insurance Prices.” In *Workshop on the Economics of Information Security*, (2019)

APPENDIX A
CODE AND DATA APPENDIX

- [14] Moore, Tyler, and Richard Clayton. "Examining the impact of website take-down on phishing." In Proceedings of the anti-phishing working groups 2nd annual eCrime researchers summit, pp. 1-13. ACM, (2007).
- [15] Franke, Ulrik, Hannes Holm, and Johan König. "The distribution of time to recovery of enterprise IT services." *IEEE Transactions on Reliability* 63, no. 4 (2014): 858-867.
- [16] Edwards, Benjamin, Steven Hofmeyr, and Stephanie Forrest. "Hype and heavy tails: A closer look at data breaches." *Journal of Cybersecurity* 2, no. 1 (2016).
- [17] , Hubbard, Douglas W., and Richard Seiersen. *How to measure anything in cybersecurity risk*. Hoboken: Wiley, (2016).
- [18] Brzezinski, Michal. "Do wealth distributions follow power laws? Evidence from 'rich lists'." *Physica A: Statistical Mechanics and its Applications* 406 (2014).
- [19] Gabaix, Xavier. "Power laws in economics: An introduction." *Journal of Economic Perspectives* 30, no. 1 (2016).
- [20] Adamic, L. A. Huberman, B.A. "Power-Law Distribution of the World Wide Web." *Science*, 287, 2115a. (2000).
- [21] Barabasi, A.L. Albert, R. "Emergence of Scaling in Random Networks." *Science*, 286, 509-512. (1999).
- [22] Faloutsos, M. Faloutsos, P. Faloutsos, C. "On Power-Law Relationships of the Internet Topology." *Proceedings of the Conference on Applications, Technologies, Architectures, and Protocols For Computer Communication* (1999). Retrieved from: <http://www.cs.cmu.edu/~christos/PUBLICATIONS/sigcomm99.pdf>
- [23] Newman, Mark EJ. "Power laws, Pareto distributions and Zipf's law." *Contemporary physics*, (2005).
- [24] Clauset, Aaron, Cosma Rohilla Shalizi, and Mark EJ Newman. "Power-law distributions in empirical data." *SIAM review* 51, (2009).
- [25] Kousky, Carolyn, and Roger Cooke. "Explaining the failure to insure catastrophic risks." *The Geneva Papers on Risk and Insurance-Issues and Practice* 37, no. 2 (2012): 206-227.
- [26] Nakamoto, Satoshi. "Bitcoin: A peer-to-peer electronic cash system." (2008).
- [27] Piketty, Thomas. "Capital in the 21st century." In *Inequality in the 21st Century*, Routledge, (2018).
- [28] Arce, Daniel G, "Malware and market share", *Journal of Cybersecurity*, Vol 4, Issue 1. (2018).
- [29] Blank, Aharon and Solomon, Sorin. "Power laws in cities population, financial markets and internet sites (scaling in systems with a variable number of components)." In *Physics a A: Statistical Mechanics and its Applications*, Vol 287, 1-2, Elsevier, (2000).
- [30] Toda, Alexis Akira. "The double power law in income distribution: Explanations and evidence." *Journal of Economic Behavior & Organization* 84, (2012).
- [31] Richardson, Lewis Fry. "Statistics of deadly quarrels"

Fig. A.1 contains a crucial snippet of the code used to fit power laws for scientific clarity. This should make it possible to reproduce the work with other datasets, or claim the code was misused in some way. A necessity of scientific publication is verifiable and reproducible results, after all.

The other tables display changes to the rolling windows of study to verify if the KS distance was a fluke of data shaping. These tables also serve to document the powerlaws continue to fit in a variety of different time intervals, which is an important property of power laws generally. TABLE A.1 explores one year rolling windows between 2013-2014. TABLE A.2 explores one year rolling windows between 2014-2015. Note particularly some of the fluctuations in KS distance, showing some windows of time that powerlaws are not found to be statistically significant within. We don't believe this invalidates our results, but does capture why tables such as these must be displayed. TABLE A.3 explores one year rolling windows between 2015-2016. This table is particularly interesting as it identifies the year the fat tails intensified. Could this be the year that ransomware gangs shifted from targeting personal computers to corporate networks? TABLE A.4 explores one year rolling windows between 2016-2017. The KS distance of many of these windows raises some questions. It is our belief that as more ransoms from these periods are gathered, the heavy tails will bear out again. However, it is also possible changes in ransom negotiations or law enforcement changed something substantial within the ecosystem. Determining what it was and why it broke the power law relationship would be a ground breaking bit of policy work. TABLE A.5 explores one year rolling windows between 2017-2018. TABLE A.6 explores one year rolling windows between 2018-2019. Not that these last two tables probably contain many fewer ransoms than actually occurred in the wild. Collecting this data carefully for the future will ensure that time will tell.

```

1 import powerlaw as powerlaw
2 import matplotlib.pyplot as plt
3 import pandas as pd
4
5 data = pd.read_csv('All-TemporalRansoms.csv',header=0)
6 data.time = pd.to_datetime(data.time, yearfirst=True)
7 data.set_index('time',inplace=True)
8 years = ['2013','2014','2015','2016','2017','2018','2019']
9 for year in years:
10     yeardata = data.loc[year+'-01-01':year+'-12-31']
11     flat_list = yeardata['USD']
12     #Remove zeros from the list for powerlaw lib and plotting pdf
13     flat_list = list(filter(lambda a: a != 0.0, flat_list))
14     if len(flat_list) > 150:
15         tailresults = powerlaw.Fit(flat_list, xmin=10000.00, discrete=False,fit_method="KS")
16         uninsuredable = powerlaw.Power_Law(xmin=10000.00,parameters=[2.0])

```

Fig. A.1. This figure shows the code used for fitting power-laws to ransoms.

TABLE A.1
ALTERNATIVE ROLLING WINDOWS 2013-2014

Time span	xmin	xmax	α	σ	KS distance
Jan 2013 to Jan 2014	10000.0	None	2.35224	0.01076	0.04851
Feb 2013 to Feb 2014	10000.0	None	2.32300	0.00957	0.05166
Mar 2013 to Mar 2014	10000.0	None	2.32445	0.00910	0.05008
Apr 2013 to Apr 2014	10000.0	None	2.43384	0.00650	0.04246
May 2013 to May 2014	10000.0	None	2.43761	0.00581	0.03746
Jun 2013 to Jun 2014	10000.0	None	2.44709	0.00542	0.03505
Jul 2013 to Jul 2014	10000.0	None	2.42149	0.00493	0.02580
Aug 2013 to Aug 2014	10000.0	None	2.41287	0.00467	0.02703
Sep 2013 to Sep 2014	10000.0	None	2.41343	0.00467	0.02711
Oct 2013 to Oct 2014	10000.0	None	2.41342	0.00467	0.02713
Nov 2013 to Nov 2014	10000.0	None	2.41571	0.00468	0.02724
Dec 2013 to Dec 2014	10000.0	None	2.42333	0.00499	0.02675

TABLE A.2
ALTERNATIVE ROLLING WINDOWS 2014-2015

Time span	xmin	xmax	α	σ	KS distance
Jan 2014 to Jan 2015	10000.0	None	2.42658	0.00517	0.02904
Feb 2014 to Feb 2015	10000.0	None	2.43867	0.00534	0.03151
Mar 2014 to Mar 2015	10000.0	None	2.44249	0.00543	0.03175
Apr 2014 to Apr 2015	10000.0	None	2.39119	0.00667	0.02820
May 2014 to May 2015	10000.0	None	2.36747	0.00776	0.04242
Jun 2014 to Jun 2015	10000.0	None	2.30978	0.00902	0.06795
Jul 2014 to Jul 2015	10000.0	None	2.33608	0.01354	0.08392
Aug 2014 to Aug 2015	10000.0	None	2.41757	0.07566	0.10755
Sep 2014 to Sep 2015	10000.0	None	2.25862	0.07645	0.05271
Oct 2014 to Oct 2015	10000.0	None	2.11551	0.06839	0.04081
Nov 2014 to Nov 2015	10000.0	None	2.10693	0.06591	0.04103
Dec 2014 to Dec 2015	10000.0	None	2.15541	0.06551	0.03761

TABLE A.3
ALTERNATIVE ROLLING WINDOWS 2015-2016

Time span	xmin	xmax	α	σ	KS distance
Jan 2015 to Jan 2016	10000.0	None	2.16231	0.06538	0.03732
Feb 2015 to Feb 2016	10000.0	None	2.16474	0.06562	0.03697
Mar 2015 to Mar 2016	10000.0	None	2.17525	0.05966	0.02780
Apr 2015 to Apr 2016	10000.0	None	2.18867	0.04328	0.02262
May 2015 to May 2016	10000.0	None	2.02853	0.03193	0.02375
Jun 2015 to Jun 2016	10000.0	None	1.94477	0.02590	0.04099
Jul 2015 to Jul 2016	10000.0	None	1.90214	0.02346	0.04616
Aug 2015 to Aug 2016	10000.0	None	1.85518	0.02082	0.05980
Sep 2015 to Sep 2016	10000.0	None	1.83482	0.01904	0.06646
Oct 2015 to Oct 2016	10000.0	None	1.82726	0.01720	0.06877
Nov 2015 to Nov 2016	10000.0	None	1.80561	0.01529	0.07469
Dec 2015 to Dec 2016	10000.0	None	1.79450	0.01454	0.07932

TABLE A.4
ALTERNATIVE ROLLING WINDOWS 2016-2017

Time span	xmin	xmax	α	σ	KS distance
Jan 2016 to Jan 2017	10000.0	None	1.78386	0.01385	0.07996
Feb 2016 to Feb 2017	10000.0	None	1.76831	0.01309	0.08193
Mar 2016 to Mar 2017	10000.0	None	1.73211	0.01209	0.08824
Apr 2016 to Apr 2017	10000.0	None	1.68462	0.01151	0.10738
May 2016 to May 2017	10000.0	None	1.67362	0.01149	0.11176
Jun 2016 to Jun 2017	10000.0	None	1.66259	0.01146	0.11147
Jul 2016 to Jul 2017	10000.0	None	1.63081	0.01099	0.11134
Aug 2016 to Aug 2017	10000.0	None	1.61096	0.01094	0.10748
Sep 2016 to Sep 2017	10000.0	None	1.58945	0.01079	0.10884
Oct 2016 to Oct 2017	10000.0	None	1.56693	0.01120	0.10994
Nov 2016 to Nov 2017	10000.0	None	1.54150	0.01185	0.11897
Dec 2016 to Dec 2017	10000.0	None	1.52084	0.01208	0.12635

TABLE A.5
ALTERNATIVE ROLLING WINDOWS 2017-2018

Time span	xmin	xmax	α	σ	KS distance
Jan 2017 to Jan 2018	10000.0	None	1.50240	0.01223	0.13342
Feb 2017 to Feb 2018	10000.0	None	1.48852	0.01263	0.14053
Mar 2017 to Mar 2018	10000.0	None	1.48473	0.01393	0.12442
Apr 2017 to Apr 2018	10000.0	None	1.47737	0.01517	0.09764
May 2017 to May 2018	10000.0	None	1.45170	0.01577	0.10293
Jun 2017 to Jun 2018	10000.0	None	1.43037	0.01693	0.09390
Jul 2017 to Jul 2018	10000.0	None	1.45530	0.02093	0.08675
Aug 2017 to Aug 2018	10000.0	None	1.46989	0.02548	0.11490
Sep 2017 to Sep 2018	10000.0	None	1.52679	0.03429	0.08446
Oct 2017 to Oct 2018	10000.0	None	1.58965	0.03480	0.04851
Nov 2017 to Nov 2018	10000.0	None	1.61834	0.03445	0.03769
Dec 2017 to Dec 2018	10000.0	None	1.70260	0.03504	0.05331

TABLE A.6
ALTERNATIVE ROLLING WINDOWS 2018-2019

Time span	xmin	xmax	α	σ	KS distance
Jan 2018 to Jan 2019	10000.0	None	1.77870	0.04204	0.06623
Feb 2018 to Feb 2019	10000.0	None	1.82802	0.04636	0.07462
Mar 2018 to Mar 2019	10000.0	None	1.84693	0.04889	0.07772
Apr 2018 to Apr 2019	10000.0	None	1.92612	0.05495	0.06630
May 2018 to May 2019	10000.0	None	2.00447	0.06182	0.06792
Jun 2018 to Jun 2019	10000.0	None	2.00387	0.06261	0.07619
Jul 2018 to Jul 2019	10000.0	None	2.03304	0.06456	0.08716
Aug 2018 to Aug 2019	10000.0	None	2.01844	0.06377	0.09136
Sep 2018 to Sep 2019	10000.0	None	2.03407	0.06475	0.09215
Oct 2018 to Oct 2019	10000.0	None	1.98531	0.06932	0.06957
Nov 2018 to Nov 2019	10000.0	None	1.99236	0.07796	0.05507
Dec 2018 to Dec 2019	10000.0	None	1.96680	0.12586	0.08643

Number of files in corpus	File Type
10,590	HTML document ASCII text with very long lines with no line terminators
10,928	HTML document ISO-8859 text with very long lines with CRLF LF NEL line terminators
11,122	XML 1.0 document Non-ISO extended-ASCII text with very long lines
11,799	HTML document Non-ISO extended-ASCII text with very long lines with CRLF line terminators with overstriking
12,900	XML 1.0 document ISO-8859 text with very long lines with CRLF LF line terminators
13,011	HTML document ISO-8859 text
13,253	ASCII text
13,499	XML 1.0 document ISO-8859 text with CRLF LF line terminators
14,886	PE32 executable (GUI) Intel 80386 Mono/.Net assembly for MS Windows
15,097	ASCII text with very long lines with no line terminators
16,035	HTML document Non-ISO extended-ASCII text with very long lines with CRLF CR NEL line terminators
17,587	HTML document ISO-8859 text with CRLF LF line terminators
17,795	XML 1.0 document ASCII text with very long lines
19,213	XML 1.0 document ISO-8859 text with CRLF CR LF line terminators
19,618	HTML document ISO-8859 text with very long lines with CRLF NEL line terminators
24,570	Bourne-Again shell script ASCII text executable
24,719	HTML document UTF-8 Unicode text with CRLF line terminators
31,863	Java archive data (JAR)
31,991	HTML document ASCII text with CRLF LF line terminators
32,029	HTML document ISO-8859 text with CRLF line terminators
32,076	HTML document UTF-8 Unicode (with BOM) text with very long lines with CR LF line terminators
33,475	XML 1.0 document Non-ISO extended-ASCII text with very long lines with CRLF LF line terminators
36,333	PE32 executable (GUI) Intel 80386 for MS Windows UPX compressed
40,915	HTML document Non-ISO extended-ASCII text with very long lines with LF NEL line terminators
42,063	XML 1.0 document UTF-8 Unicode text with very long lines with CRLF CR LF line terminators
42,280	PE32 executable (GUI) Intel 80386 for MS Windows Nullsoft Installer self-extracting archive
46,890	HTML document Non-ISO extended-ASCII text with very long lines with CRLF CR line terminators
47,505	UTF-8 Unicode text with very long lines
49,164	HTML document Non-ISO extended-ASCII text with very long lines with CRLF CR LF NEL line terminators
49,909	PE32 executable (DLL) (GUI) Intel 80386 for MS Windows
55,625	XML 1.0 document UTF-8 Unicode text with very long lines with CRLF LF line terminators
55,728	ASCII text with very long lines
60,420	HTML document UTF-8 Unicode (with BOM) text with very long lines with CRLF CR LF line terminators
61,236	HTML document UTF-8 Unicode text with very long lines with CRLF CR line terminators
64,568	FGDC-STD-001-1998
72,059	Dalvik dex file version 035
77,197	XML 1.0 document UTF-8 Unicode text with very long lines
83,215	HTML document UTF-8 Unicode text with CRLF LF line terminators
84,585	HTML document ISO-8859 text with very long lines with CRLF CR LF line terminators
85,299	HTML document UTF-8 Unicode (with BOM) text with very long lines
91,387	HTML document ASCII text with CRLF line terminators
95,065	HTML document UTF-8 Unicode text
112,427	HTML document ASCII text with very long lines with CR LF line terminators
123,561	Non-ISO extended-ASCII text with CRLF NEL line terminators
129,410	HTML document ASCII text
137,262	HTML document Non-ISO extended-ASCII text with very long lines with CRLF CR LF line terminators
168,488	HTML document UTF-8 Unicode (with BOM) text with very long lines with CRLF line terminators
169,299	HTML document UTF-8 Unicode text with very long lines with no line terminators
239,686	XML 1.0 document Non-ISO extended-ASCII text with very long lines with LF NEL line terminators
240,103	Zip archive data at least v2.0 to extract
333,889	HTML document ISO-8859 text with very long lines with CRLF CR line terminators
352,830	HTML document Non-ISO extended-ASCII text with very long lines with CRLF LF NEL line terminators
411,263	HTML document ASCII text with very long lines with CRLF line terminators
481,029	HTML document ASCII text with very long lines with CRLF CR LF line terminators
511,081	HTML document UTF-8 Unicode text with very long lines with CR LF line terminators
703,792	data
861,712	PE32 executable (GUI) Intel 80386 for MS Windows
895,169	HTML document Non-ISO extended-ASCII text with very long lines with CRLF NEL line terminators
1,384,847	ASCII text with very long lines with CRLF line terminators
1,478,909	HTML document UTF-8 Unicode (with BOM) text with very long lines with CRLF LF line terminators
1,928,664	HTML document ISO-8859 text with very long lines with CRLF LF line terminators
2,226,704	HTML document Non-ISO extended-ASCII text with very long lines with CRLF LF line terminators
2,751,110	HTML document ISO-8859 text with very long lines
3,524,466	HTML document ISO-8859 text with very long lines with CRLF line terminators
6,551,316	HTML document ASCII text with very long lines with CRLF LF line terminators
9,092,887	HTML document Non-ISO extended-ASCII text with very long lines
12,062,927	HTML document UTF-8 Unicode text with very long lines with CRLF CR LF line terminators
13,097,315	HTML document ASCII text with very long lines
18,345,823	HTML document UTF-8 Unicode text with very long lines with CRLF line terminators
22,905,879	HTML document UTF-8 Unicode text with very long lines with CRLF LF line terminators
26,003,413	HTML document Non-ISO extended-ASCII text with very long lines with CRLF line terminators
48,203,226	HTML document UTF-8 Unicode text with very long lines