

Dynamic Models for File Sizes and Double Pareto Distributions

Michael Mitzenmacher

Abstract. In this paper, we introduce and analyze a new, dynamic generative user model to explain the behavior of file size distributions. Our Recursive Forest File model combines multiplicative models that generate lognormal distributions with recent work on random graph models for the web. Unlike similar previous work, our Recursive Forest File model allows new files to be created and old files to be deleted over time, and our analysis covers problematic issues such as correlation among file sizes. Moreover, our model allows natural variations where files that are copied or modified are more likely to be copied or modified subsequently.

Previous empirical work suggests that file sizes tend to have a lognormal body but a Pareto tail. The Recursive Forest File model explains this behavior, yielding a double Pareto distribution, which has a Pareto tail but close to a lognormal body. We believe the Recursive Forest model may be useful for describing other power law phenomena in computer systems as well as other fields.

1. Introduction

In this paper, we attempt to provide a simple generative user model that provides a good approximation for file size distributions. Accurate models for file size distributions are important for both our current understanding of and simulation of file systems and the Internet. In the case of file systems, the problem of capacity planning requires estimating when additional storage space will become necessary. An accurate model for how file systems develop over time might allow

more accurate predictions, easing the burdens of system managers. Similarly, simple generative models may enhance simulation tools for file system behavior. For the Internet, many studies have shown that traffic patterns in the Internet appear to have self-similarity (see, e.g., [Barford and Crovella 98, Barford et al. 99, Crovella and Bestavros 97, Crovella et al. 98, Leland et al. 94]). This self-similarity can possibly be explained if the underlying distribution of file sizes obeys an appropriate power law [Crovella and Bestavros 97]. Understanding why a power law distribution for files might or might not arise naturally is therefore important. Tools used to generate web workloads such as SURGE [Barford and Crovella 98], which can be useful in testing or simulating web servers, may also require a suitable model for simulating file size distributions and how they change over time.

We emphasize that providing a generative model is a fundamentally different task than fitting a model to data, which has been the primary focus of most previous work. In particular, determining possible dynamic generative models is important if one wants to determine what the distribution might look like in the future, as the system changes over time. Without a justified underlying generative model, extrapolating future behavior based on fitting models to data is a risky proposition.

We provide a model that combines long-known multiplicative models and recent work on models for the web graph [Barabási et al. 99, Broder et al. 00, Drinea et al. 01, Kleinberg et al. 99, Krapivsky and Redner 01, Kumar et al. 00]. Our work was inspired by recent work by Downey [Downey 01]. Downey suggests the following idea: one way that users create new files is by taking old files and performing modifications on them, including possibly editing, copying, translating, or filtering. The size of such a new file can be modeled by taking the size of an old file and multiplying it by a random variable. Downey suggests that this model yields a lognormal distribution for file sizes, which arguably counters other previous work that has suggested file size distributions have a lognormal body, but a heavy tail [Barford and Crovella 98, Barford et al. 99].¹

Downey's model suffers from the weakness that all files derive from a single initial file. Files not derived from extant files cannot enter the file system, and old files are not deleted. We expand to a *dynamic* model; that is, we allow additions and deletions in a natural way. As a result, we obtain a family of models, which we refer to generally as the Recursive Forest File model. What is most interesting is that our changes have a dramatic effect in the analysis. The resulting distribution of file sizes is a double Pareto distribution, which we

¹We believe that there are minor problems with Downey's analysis, as we describe in Section 3.

define and describe in Section 2.3. Double Pareto distributions have recently been suggested to describe income distributions and other power law phenomena [Reed 03, Reed and Jorgensen 01]. As we show, such distributions have a lognormal body and a Pareto tail, which matches some previous studies of empirical data for file sizes. We believe that such distributions may be useful for modeling other power law phenomena in computer systems, and we believe our generative model may prove useful for other applications.

We provide a detailed analysis of the Recursive Forest File model that is interesting in its own right. In particular, we find several connections to the theory of random graphs that we expect will provide a useful framework for future work. We also show how to cope with the effects of correlation that are implicit in a file system model where new files are derived from existing files, using a martingale analysis.

In related prior work, the Highly Optimized Tolerance (HOT) model provides another generative model for file size distributions which uses an optimization framework [Carlson and Doyle 99, Zhu et al. 01]. Fabrikant, Koutsoupias, and Papadimitriou specifically utilize this framework to develop a model for file sizes [Fabrikant et al. 02]. Downey suggests (and we concur) that applying this framework to web file systems requires strong assumptions about how web sites are designed, and does not explain why local file systems have similar file size distributions [Downey 01]. Downey's simpler framework appears more intuitively appealing, and therefore we have focused on improving it. We caution, however, that any simple user model is necessarily only approximate, and certainly various models may apply in different situations. Indeed, it may be that our model is useful for describing some types of file systems while HOT-based models are better for other types of systems.

It is also worth noting that this potential confusion between whether file size distributions obey a power law or follow a lognormal distribution is not surprising. Similar discussions have arisen in many fields over several decades. Indeed, there is a rich history of models that generate power law and lognormal distributions, and many models that have been recently proposed to explain such distributions in computer systems have historical antecedents in other fields. Moreover, there are extremely close connections between generative models for power law distributions and lognormal distributions. Rather than dwell on these issues here, we refer the reader to a related historically oriented survey [Mitzenmacher 04].

A natural question not tackled in this paper is the question of verifying our model. We have not focused on this issue because we believe the primary contribution in this paper is the description and analysis of a simple, general model that yields double Pareto distributions. We believe that our model is interesting

in its own right and expect that it will find uses explaining other phenomena besides file size distribution.

It is worth pointing out, however, that subsequent to this work, Mitzenmacher and Tworetzky have performed a careful empirical study of file size distributions, examining how various models fit various data sets [Mitzenmacher and Tworetzky 03]. To summarize this work, double Pareto distributions do appear to fit data sets roughly as well as lognormal distributions, although they appear slightly worse than lognormal Pareto hybrid distribution and another prospective distribution, the log- t distribution. We emphasize that none of these other distributions currently have natural generative models of which we are aware. Further experiments yield that the double Pareto distribution fits better for HTML files than for GIF or JPEG files, a conclusion which is understandable in light of our generative model described below. Another possible approach, besides testing the fit of the distribution, would be to empirically validate the underlying assumptions of our model. Such validation would be an interesting area for future work.

The paper proceeds as follows. In Section 2.1, we provide an extensive review of the relevant terminology. This review includes definitions of Pareto, lognormal, and the more recent double Pareto distributions. In Section 3, we consider Downey's model, examining its motivation and potential problems. We develop the Recursive Forest File Model in Section 4, demonstrating interesting connections to random graph theory. We present simulation results in Section 5.

2. Review of Definitions

We briefly review the relevant definitions. For greater details, we recommend references [Aitchison and Brown 57, Crow and Shimura 88, Li 99, Mitzenmacher 04].

2.1. Power Law Distributions

For our purposes, a nonnegative random variable X is said to have a power law distribution if the *complementary cumulative distribution function* (ccdf), or $\Pr[X > x]$, satisfies

$$\Pr[X > x] \sim cx^{-\alpha}$$

for constants $c > 0$ and $\alpha > 0$. Here, $f(x) \sim g(x)$ denotes that the limit of the ratios goes to 1 as x grows large. One specific commonly used power law distribution is the Pareto distribution, which satisfies

$$\Pr[X > x] = \left(\frac{x}{k}\right)^{-\alpha}$$

for some $\alpha > 0$ and $k > 0$. Note the Pareto distribution requires $X \geq k$. If α falls in the range $0 < \alpha \leq 2$, then X has infinite variance. If $\alpha \leq 1$, then X also has an infinite mean. The density function for the Pareto distribution is $f(x) = \alpha k^\alpha x^{-\alpha-1}$.

If X has a power law distribution, then in a log-log plot of the cdf, asymptotically the behavior is a straight line. This is the basis for many tests for power law behavior. The same is true for the density function, which we find easier to work with mathematically. For example, for the Pareto distribution, the log of the density function is exactly linear:

$$\ln f(x) = (-\alpha - 1) \ln x + \alpha \ln k + \ln \alpha.$$

2.2. Lognormal Distributions

A random variable X has a lognormal distribution if the random variable $Y = \ln X$ has a normal (i.e., Gaussian) distribution. The density function for a lognormal distribution satisfies

$$f(x) = \frac{1}{\sqrt{2\pi}\sigma x} e^{-(\ln x - \mu)^2 / 2\sigma^2},$$

where μ is the mean and σ is the standard deviation of the associated normal distribution. We will say that X has parameters (μ, σ^2) when the associated normal Y has mean μ and variance σ^2 , where the meaning is clear. The lognormal distribution is skewed, with mean $e^{\mu + \frac{1}{2}\sigma^2}$, median e^μ , and mode $e^{\mu - \sigma^2}$. Although the lognormal distribution, in contrast to the Pareto distribution, has finite moments, it is extremely similar in shape to power law distributions, in that a large portion of the body of the density function and the cdf can appear linear [Mitzenmacher 04, Montroll and Shlesinger 83]. Specifically, for a lognormal distribution, we have

$$\ln f(x) = -\ln x - \ln \sqrt{2\pi}\sigma - \frac{(\ln x - \mu)^2}{2\sigma^2}. \quad (2.1)$$

If σ is sufficiently large, then the quadratic term above is small for a large range of x values, and hence the logarithm of the density function will appear linear for a large range of values. (The same is also therefore true for the cdf.)

Recall that normal distributions have the property that the sum of two independent normal random variables Y_1 and Y_2 with means μ_1 and μ_2 and variances σ_1^2 and σ_2^2 , respectively, is a normal random variable with mean $\mu_1 + \mu_2$ and variance $\sigma_1^2 + \sigma_2^2$. It follows that the *product* of independent random variables with lognormal distributions is a random variable with a lognormal distribution.

Lognormal distributions can be naturally generated by *multiplicative processes*. We start with a biological example. Suppose we start with an organism of size

X_0 . At each step j , the organism may grow or shrink by a certain percentage, according to a random variable F_j , so that

$$X_j = F_j X_{j-1}.$$

If the $F_k, 1 \leq k \leq j$, are all governed by independent lognormal distributions, then so is each X_j , inductively, since the product of lognormal random variables is again a lognormal random variable. More generally, approximately lognormal distributions may be obtained even if the F_j are not themselves lognormal. Specifically, consider

$$\ln X_j = \ln X_0 + \sum_{k=1}^j \ln F_k.$$

Assuming the random variables $\ln F_k$ satisfy appropriate conditions, the Central Limit Theorem says that $\sum_{k=1}^j \ln F_k$ converges to a normal distribution, and hence for sufficiently large j , X_j is well approximated by a lognormal distribution. In particular, if the $\ln F_k$ are independent and identically distributed variables with finite mean and variance, then asymptotically X_j will approach a lognormal distribution. Lognormal distributions are natural for describing growth of organisms, growth in options prices, and any process where over a time-step the underlying growth is a random factor independent of the current size [Crow and Shimura 88, Mitzenmacher 04].

2.3. From Lognormal to Power Law: Double Pareto Distributions

Before presenting our model, we explain how a natural mixture of lognormal distributions yields a power law distribution. This result provides the foundation for much of our later analysis, and is interesting in its own right. We therefore present it first and show how it arises in the context of the model subsequently.

Suppose we have a system $X_t = F_t X_{t-1}$, where $X_0 = 1$ and F_t is a lognormal distribution with parameters (μ, σ^2) . We think of the index t as referring to time. If we let the system run and stop it at some fixed time k , we obtain a random variable from the lognormal distribution with parameters $(k\mu, k\sigma^2)$. Suppose instead we run the process until some random time k . Then we obtain a random variable that comes from a mixture of lognormal distributions. Specifically consider the case where we have a geometric mixture of lognormal distributions: we stop the process at time $k \geq 1$ with probability $\gamma(1 - \gamma)^{k-1}$, where γ is the parameter for the geometric distribution. Hence, with probability $\gamma(1 - \gamma)^{k-1}$, we obtain random numbers from the lognormal distribution with parameters $(k\mu, k\sigma^2)$. We claim that the resulting distribution from this mixture will have a power law.

To see this, we present a result of Reed [Reed 03, Reed and Jorgensen 01] for the continuous analogue where the “mixture” of an exponentially distributed number of lognormal distributions is considered, in a sense clarified below.² Suppose that we choose a random number X from a lognormal distribution with parameters $(k\mu, k\sigma^2)$, where k itself is a random variable with an exponential distribution with mean $1/\lambda$. The resulting density function is

$$f(x) = \int_{k=0}^{\infty} \lambda e^{-\lambda k} \frac{1}{\sqrt{2\pi k\sigma x}} e^{-(\ln x - k\mu)^2 / 2k\sigma^2} dk. \quad (2.2)$$

Using the substitution $k = u^2$ gives

$$f(x) = \frac{2\lambda e^{\mu \ln x / \sigma^2}}{\sqrt{2\pi x\sigma}} \int_{u=0}^{\infty} e^{-(\lambda + \mu^2 / 2\sigma^2)u^2} e^{-(\ln x)^2 / 2\sigma^2 u^2} du.$$

An integral table gives us the identity

$$\int_{z=0}^{\infty} e^{-az^2 - b/z^2} = \frac{1}{2} \sqrt{\frac{\pi}{a}} e^{-2\sqrt{ab}},$$

which allows us to solve for the resulting form. Note that in the exponent $2\sqrt{ab}$ of the identity we have $b = (\ln x)^2 / 2\sigma^2$. Because of this, there are two different behaviors, depending on whether $x \geq 1$ or $x \leq 1$. Let $C_1 = \lambda / (\sigma (\sqrt{(\mu/\sigma)^2 + 2\lambda}))$ and let $C_2 = (\sqrt{(\mu/\sigma)^2 + 2\lambda}) / \sigma$. For $x \geq 1$, $f(x) = C_1 x^{-1 + \mu/\sigma^2 - C_2}$, so the result is a power law distribution. For $x \leq 1$, $f(x) = C_1 x^{-1 + \mu/\sigma^2 + C_2}$. In particular, a case we use later is when $\mu = 0$ and $\sigma = 1$. In this case, for $x \geq 1$, $f(x) = (\sqrt{\lambda/2}) x^{-1 - \sqrt{2\lambda}}$, and for $x \leq 1$, $f(x) = (\sqrt{\lambda/2}) x^{-1 + \sqrt{2\lambda}}$. Reed therefore suggests the following stringent definition:

Definition 2.1. A double Pareto distribution defined over $x > 0$ with parameters $\alpha, \beta > 0$ has $f(x) = \frac{\alpha\beta}{\alpha+\beta} x^{\beta-1}$ for $0 < x \leq 1$ and $f(x) = \frac{\alpha\beta}{\alpha+\beta} x^{-\alpha-1}$ for $x > 1$.

A key characteristic of the double Pareto distribution is that it has a power law at both tails. That is, if we look at the cumulative distribution function (cdf) on a log-log plot, it will also have a linear tail (for the small files). This provides a test for seeing whether a distribution has a double Pareto distribution; look at both the ccdf and the cdf on log-log plots for linear tails.

²Huberman and Adamic [Huberman and Adamic 99, Huberman and Adamic 00] also examine this distribution and conclude that it has a power law distribution. Their earlier work, however, fails to note that the behavior of the distribution goes through a phase shift, which Reed clarifies.

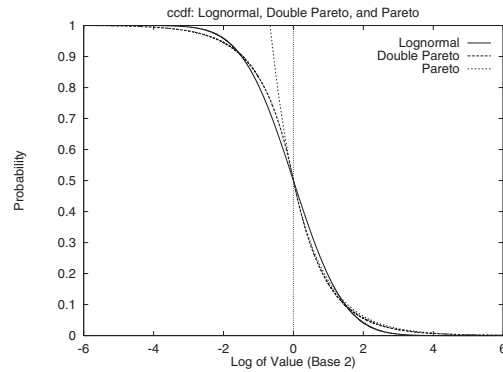


Figure 1. Shapes of lognormal, double Pareto, and Pareto distributions.

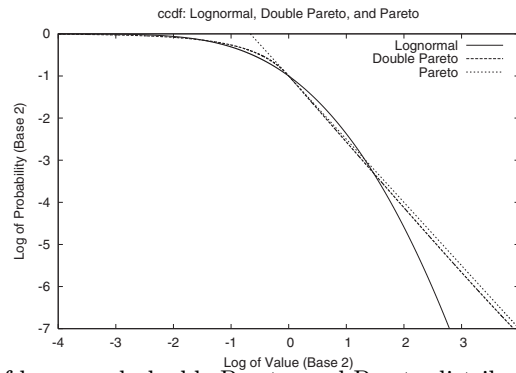


Figure 2. Shapes of lognormal, double Pareto, and Pareto distributions—log-log plot.

The double Pareto distribution falls nicely between the lognormal distribution and the Pareto distribution. Like the Pareto distribution, it is a power law distribution. But while the log-log plot of the density of the Pareto distribution is a single straight line, for the double Pareto distribution the log-log plot of the density consists of two straight line segments that meet at a transition point. This is similar to the lognormal distribution, which has a transition point around its median e^μ due to the quadratic term, as shown in (2.1). Hence, an appropriate double Pareto distribution can closely match the body of a lognormal distribution and the tail of a Pareto distribution. For example, Figure 1 shows the complementary cumulative distribution function for a lognormal, double Pareto, and Pareto distribution. (These graphs have only been minimally tuned to give a reasonable pictorial match; they could be made to match more closely.) The lognormal and double Pareto distributions match quite well with a standard scale for probabilities, but on the log-log scale in Figure 2 one can see the difference in the tail behavior, where the double Pareto more closely matches the Pareto.

When we have the discrete geometric mixture instead of the continuous exponential mixture, the proper equation for the density function is

$$f(x) = \sum_{k=1}^{\infty} (\gamma(1-\gamma)^{k-1}) \left(\frac{1}{\sqrt{2\pi kx\sigma}} e^{-(\ln x - k\mu)^2/2k\sigma^2} \right). \tag{2.3}$$

The summation is well approximated when $\ln x$ is very large or very small by the corresponding integral

$$f(x) \approx \int_{k=1}^{\infty} \frac{\gamma}{\sqrt{2\pi kx\sigma}(1-\gamma)} e^{k \ln(1-\gamma) - (\ln x - k\mu)^2/2k\sigma^2} dk. \tag{2.4}$$

Comparing (2.2) and (2.4), we see that essentially the same tail behaviors from the geometric mixture as the exponential mixture (although we do not obtain such a nice closed form). That is, we have the following theorem:

Theorem 2.2. *There exist positive constants $\alpha, \beta, c_1, c_2, c_3, c_4, m$, and ϵ such that the density function in (2.3) satisfies $c_1 x^{\beta-1} \leq f(x) \leq c_2 x^{\beta-1}$ for $x < \epsilon$, and $c_3 x^{-\alpha-1} \leq f(x) \leq c_4 x^{-\alpha-1}$ for $x > m$. (Here the c_i may depend on γ, μ , and σ but not on x .)*

Proof. For the proof, let

$$f_1(x) = \sum_{k=1}^{\infty} (\gamma(1-\gamma)^{k-1}) \left(\frac{1}{\sqrt{2\pi kx\sigma}} e^{-(\ln x - k\mu)^2/2k\sigma^2} \right)$$

and

$$f_2(x) = \int_{k=0}^{\infty} \lambda e^{-\lambda k} \frac{1}{\sqrt{2\pi k\sigma x}} e^{-(\ln x - k\mu)^2/2k\sigma^2} dk.$$

As we have shown that $f_2(x) = x^{-\alpha-1}$ for $x > 1$ and $f_2(x) = x^{\beta-1}$ for $x < 1$ for appropriate α and β , it suffices to show that for sufficiently large and small values of x that when $\lambda = \ln(1-\gamma)$, f_1 and f_2 differ by at most constant factors. After separating out constant factors, we find

$$f_1(x) = \frac{C_1}{x} \sum_{k=1}^{\infty} e^{-ak - b/k - (\ln k)/2}$$

and

$$f_2(x) = \frac{C_2}{x} \int_{k=0}^{\infty} e^{-ak - b/k - (\ln k)/2} dk,$$

where C_1 and C_2 are positive constants, $b = (\ln x)^2/(2\sigma^2)$, and a is a positive constant independent of x . Hence, it suffices to show that

$$C_3 \int_{k=0}^{\infty} e^{-ak - b/k - (\ln k)/2} dk \leq \sum_{k=1}^{\infty} e^{-ak - b/k - (\ln k)/2} \leq C_4 \int_{k=0}^{\infty} e^{-ak - b/k - (\ln k)/2} dk$$

for some positive constants C_3 and C_4 . The important point in showing this is to keep track of the important term $e^{-b/k}$, which is increasing in k . Hence, it is easy to show that for $k > 0$,

$$e^{-ak-b/k-(\ln k)/2} \leq e^{-a[k]-b/[k]-(\ln[k])/2} e^{a+1};$$

this yields

$$\sum_{k=1}^{\infty} e^{-ak-b/k-(\ln k)/2} \leq C_4 \int_{k=0}^{\infty} e^{-ak-b/k-(\ln k)/2} dk.$$

Similarly, for $k \geq 1$,

$$e^{-a[k]-b/[k]-(\ln[k])/2} e^{-a-1} \leq e^{-ak-b/k-(\ln k)/2}.$$

From this, we have that

$$\sum_{k=1}^{\infty} e^{-ak-b/k-(\ln k)/2} \geq C'_3 \int_{k=1}^{\infty} e^{-ak-b/k-(\ln k)/2} dk,$$

which is almost the desired result. It now suffices to note that

$$\int_{k=0}^1 e^{-ak-b/k-(\ln k)/2} dk \leq \int_{k=1}^2 e^{-ak-b/k-(\ln k)/2} dk$$

for values of b sufficiently large, which occurs whenever $\ln x$ is sufficiently large. Hence, using $C_3 = 2C'_3$, we obtain that there exists ϵ and m such that for $x < \epsilon$ and $x > m$,

$$\sum_{k=1}^{\infty} e^{-ak-b/k-(\ln k)/2} \geq C_3 \int_{k=0}^{\infty} e^{-ak-b/k-(\ln k)/2} dk,$$

yielding the result. \square

Technically, the geometric mixture of lognormal distributions yields an approximate double Pareto distribution, and not a true double Pareto distribution according to Reed's stringent definition. For convenience, we ignore this distinction henceforth in the paper, and refer to both the result of the exponential mixture of lognormal distributions and the geometric distribution of lognormal distributions as double Pareto. In particular, from Theorem 2.2 it follows at the tails (for $x < \epsilon$ and $x > m$) that the cumulative distribution function and complementary cumulative distribution function of the geometric mixture are each bounded by two power law distributions that differ only by constant factors.

This fact, that a geometric mixture of lognormal distributions yields a double Pareto distribution, plays an important role in the development of our Recursive Forest File model throughout the rest of the paper.

The requirement that the exponential mixture be of lognormal distributions can be weakened substantially without changing the tail behaviors [Reed and Hughes 02]. We use exponential mixtures of lognormal distributions throughout the paper for convenience. Also, Reed has suggested a generalization of the double Pareto distributions called double Pareto-lognormal distributions with similar properties [Reed and Jorgensen 01]. The double Pareto-lognormal distribution has more parameters, but might allow closer matches with empirical distributions.

3. Downey's Multiplicative File Size Model

3.1. The Basic Model

We now present Downey's model to provide appropriate background. In particular, we point out weaknesses in Downey's model that we ameliorate and introduce features of analysis that prove useful subsequently for our dynamic model.

Downey's model for file sizes is based on the following idea: users tend to create new files from old files, by copying, editing, or filtering in some way. Downey therefore suggests the following model. The system begins with a single file S_0 , and the user repeatedly performs the following actions.

- Select a file S to modify uniformly at random. Let the size of S be s .
- Choose a multiplicative factor f from a given distribution \mathcal{D} .
- Create a new file S' with size fs .

The assumption behind this model is that creating a new file from a template file from processes such as copying, editing, translating, or filtering yields a file whose size differs from the template file by a factor that is independent of the size of the template. With filtering, for example, a fraction of the input may be recorded. For editing, if the amount of changes made is proportional to the size of the file (three edits per page), then this assumption appears reasonable. (Arguably, in many cases edits are additive rather than multiplicative; a constant number of changes are made. This can be modeled in a way reasonably consistent with the assumption by giving the distribution \mathcal{D} a strong mode around 1.)

Looking at any individual file, there is a history of j steps that created all the previous versions, or predecessors, of that file. That is, a file S_j was created

from a file S_{j-1} and so on back to the root S_0 . Let X_0 represent the size of S_0 and let F_k represent the random multiplicative factor chosen from \mathcal{D} in the creation of S_k . Then $\ln X_j = \ln X_0 + \sum_{k=1}^j \ln F_k$, and hence if \mathcal{D} is lognormal the distribution of the size of any specific individual file is lognormal. Alternatively, even if \mathcal{D} is not lognormal, X_j will be approximately lognormal if j is sufficiently large. Downey therefore suggests that the entire file size distribution resulting from this process is lognormal. This is not entirely accurate, as we explain below.

We note that preliminary empirical studies by Downey suggest that the right distribution for \mathcal{D} is roughly lognormal, although it is more leptokurtotic; that is, there are more values near the mode, which is close to 1 since the most common operation on a file is a copy or a small change [Downey 01]. Downey finds that this has little effect on the overall results; again, this is justified by the analysis in [Reed and Hughes 02].

3.2. Random Tree Models

We provide an alternative view of the generative file process above by embedding it into a tree structure. Initially, we start with a root node, corresponding to the initial file. For convenience let us here take the size of the original file to be 1.

At each step, a random node of the current tree is chosen, and a new child of that node is created. Each node therefore corresponds to a new file that was created from the file corresponding to its parent, and the path from the root to the node corresponds to the file history. From here on, we use the terms node and file interchangeably. If we think of each edge as being labeled by a multiplicative factor, then by multiplying the numbers on the path from the root to a node we obtain its size (relative to the root node). Alternatively, we consider each edge as being labeled with the log of the multiplicative factor; then summing the weight along each edge gives the logarithm of the file size. As in Downey's model, let us suppose the multiplicative factor is always chosen from a given distribution \mathcal{D} .

When we say the distribution of file sizes of a file system with t files, we mean the following. From some initial starting state, we generate new files according to the process above, until there are t files. The file size distribution is the distribution obtained by choosing a file uniformly at random from the t resulting files.

This tree model emphasizes that files have varying depths. While nodes at the same depth have the same size distribution, the size distribution varies for nodes at different depths. Assuming that the distribution of the growth factor is lognormally distributed, a node at depth $k \geq 1$ has a lognormally distributed size with parameters $(k\mu, k\sigma^2)$. Hence, if the file sizes were independent, the dis-

tribution would be a mixture of lognormal distributions, derived from weighing the distribution for each depth with the proportion of nodes at each depth.

The tree developed under this model is well-studied in the combinatorial literature. It is known as a *uniform random recursive tree*, since the process looks the same to each node in the tree. Results regarding the height of tree, the distribution of depths of nodes, and so on are known. We provide a brief summary, based on [Smythe and Mahmoud 95]. An exact formula for the average number of nodes of depth k in a tree with n nodes is

$$\frac{1}{(n-1)!} \left[\begin{matrix} n \\ k+1 \end{matrix} \right],$$

where $\left[\begin{matrix} n \\ k \end{matrix} \right]$ is the Stirling number of the first kind, or the number of ways to arrange n objects into k nonempty cycles. Asymptotically the distribution of the depths of the nodes is sharply concentrated around $\ln n$. This explains why empirically Downey's model yields close to a lognormal distribution for file sizes; most nodes are at approximately the same depth and therefore have sizes governed by almost the same lognormal distribution, with additional symmetry to smooth out the effects of deep and shallow nodes. It is not clear that in practice we would expect the average depth of a file should be dependent on n , the number of files in the system, suggesting another problem with this model. (Arguing that the maximum depth depends on n is more clear; perhaps some file, such as a script file, is used and modified occasionally as new files arise.)

One obvious way to generalize the file model is to use a different recursive tree model, such as plane-oriented recursive trees [Devroye 98, Smythe and Mahmoud 95]. In this model, the probability that a new node is the child of a node x is proportional to $c(x) + 1$, where $c(x)$ is the number of existing children of x . (Adding one avoids problems at the leaves and root.) This model is entirely similar to current models for the web graph, which use this sort of preferential attachment in order to obtain power law distributions [Barabási et al. 99, Drinea et al. 01, Kleinberg et al. 99, Kumar et al. 00]. Such a model could apply if a user is more likely to modify versions of files that have already been modified several times. This may be quite possible—a useful shell script, for instance, may be more likely to be modified multiple times for various situations.

Specifically, in this tree model, the fraction of nodes with k children is roughly proportional to $1/k^3$, a power law distribution. In this case, in a tree with n nodes the depth of the nodes are sharply concentrated around $\frac{1}{2} \ln n$.

One can generalize this model by having the probability that a new node is the child of a node x be proportional to $b \cdot c(x) + 1$ for some constant $b > 0$. A larger constant b strengthens the effect that nodes with children get more

children; as b approaches 0, the model becomes more like the uniform random recursive tree. We revisit this possibility in the context of our Recursive Forest File model. We also note that further variations can be created by using different probabilities for the generation of children at each depth; however, such models seem excessively complex to be useful, and we avoid them here.

3.3. Correlations

The tree model also clarifies that file sizes are necessarily correlated: a child is clearly correlated to the size of its parent. Because of this, it is not clear what the resulting overall distribution of file sizes will be in this model. For example, one large multiplicative factor close to the root will affect several nodes, changing the overall distribution for an entire subtree. We emphasize that while the distribution of individual nodes is not affected by correlation, because of correlation it is difficult to make statements about the resulting joint distribution of the entire file system determined by the model.

We attempt to highlight the problem of correlation with a simple experiment. We simulated Downey's model, placing weights chosen from a normal distribution with mean 0 and variance 1 on each edge. Recall the logarithm of the ratio of the file size at a node to the initial file size is the sum of the weights on the edges along the path from the root; using this distribution, the average of these values, or the *average log ratio*, should be 0. Over 1,000 different runs generating 10,000 files, we found the average log ratio varied significantly, between -4.2 and 5.2 . The absolute value of the average log ratio was greater than 2 more than 150 times. These high average log ratios occur even though the sample variance is small; it is generally between 5 and 10. Moreover, similar experiments generating 100,000 and 1,000,000 files yield the same high average log ratios; over 1,000 trials, the range was roughly the same, and about 15% of the trials have average log ratio with absolute value at least 2. This effect is entirely due to the fact that a single large edge near the root can affect many nodes, moving the entire average log ratio. For a comparison, we performed 1,000 trials of taking the average of 10,000 independent normal random variables with mean 0 and an extremely large variance of 100. The distribution of the average is a random variable with mean 0 and a standard deviation of 0.1; over 1,000 trials, the averages ranged between -0.33 and 0.32 .

Such correlations are problematic both philosophically and practically. At an abstract level, we do not expect that the behavior of the distribution resulting from the model should potentially vary so significantly from trial to trial. More practically, such correlations are problematic because they render the model much less useful for predicting behavior based on the results.

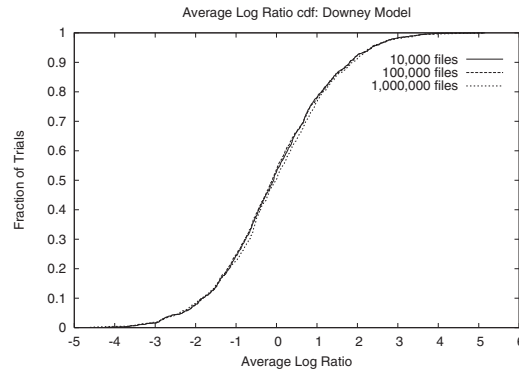


Figure 3. The cdf of the average log ratio using Downey’s model.

3.4. Minimum File Sizes

A further potential argument against the multiplicative model is that it allows files to grow arbitrarily small as well as arbitrarily big. In practice, there is generally a natural lower bound to a file size (for instance, one byte). It is therefore worth asking how the multiplicative process behaves when there is a lower bound on the minimum size. That is, suppose that we have a (near) multiplicative process

$$X_j = \max\{F_j X_{j-1}, \epsilon\}$$

for some constant ϵ . In this case, the limiting distribution of X_j is not lognormal, but instead a power law [Gabaiz 99]. This close connection between the lognormal and power law distributions is discussed more fully in [Mitzenmacher 04], but it suggests that attempting to distinguish strictly between file size models that yield lognormal distributions from models that yield power law distributions may be a futile exercise. We avoid further focus on this issue in the analysis, however; generally we believe the effect on the model is relatively minor.

4. The Recursive Forest File Model

4.1. Insertions

We now suggest a new class of dynamic models, based on similar dynamic models for modeling web graphs. We call our models dynamic because they allow the introduction of new files into the system as well as the deletion of old files. We begin by handling the insertion of new files only. We also temporarily ignore the problems of correlation until Section 4.3.

The ability to handle the insertion of new files is clearly important for modeling current systems, where new content (such as audio, video, and text) are often created or downloaded from external sources, such as the Internet. Our model begins with a collection of one or more files, whose sizes are drawn from a distribution \mathcal{D}_1 . Repeatedly, new files are generated as follows:

- with probability γ , add a new file with size chosen from a given distribution \mathcal{D}_1 .
- with probability $1 - \gamma$: select a file S (with size denoted by s) uniformly at random, choose a multiplicative factor f from a given distribution \mathcal{D}_2 , and create a new file S' with size fs .

This generalizes the uniform random recursive tree model, so that the model produces a random recursive forest [Balinska et al. 94]. This explains why we refer to our class of models as Recursive Forest File models. Also, we have given each file an initial size. Implicitly, we may think of an edge to each root giving its initial size.

We first ask in this model how many nodes of each depth k there are when n files are in the system. Note that we could write exact recurrences for the expected value of these variables and use martingale arguments to obtain high probability results. In the interest of space and highlighting the idea of the model, we present here a more intuitive limiting argument. Let $X_{t,j}$ be the number of nodes at depth j at time t . Since new nodes of depth 0, or roots, enter the system with probability γ , it is clear that $X_{t,0}/t \rightarrow \gamma$, where \rightarrow signifies convergence with probability 1 in the limit as t goes to infinity. Now for $X_{t,1}$ to increase, a new node that is the child of an existing node must enter; this happens with probability $1 - \gamma$. Its parent must be a root; if $X_{t,0}$ is γt , this occurs with probability γ . Hence, nodes of depth 1 arise at a rate of $\gamma(1 - \gamma)$, so $X_{t,1}/t \rightarrow \gamma(1 - \gamma)$. Continuing inductively, we find (asymptotically) that $X_{t,j}$ approaches $\gamma(1 - \gamma)^j t$; that is, node depths have a geometric distribution.

Lemma 4.1. *In the Random Recursive Forest File model,*

$$\lim_{t \rightarrow \infty} \frac{X_{t,j}}{t} = \gamma(1 - \gamma)^j.$$

Martingale arguments, quite similar to those in [Kumar et al. 00], can be used to yield high probability results. Alternatively, the framework relating differential equations and martingales used by Kurtz [Kurtz 81] and enhanced by the work of Wormald [Wormald 95] allow one to state concentration results

for nodes of any constant depth j . (See, for example, Theorem 1 of [Wormald 95].)

This model has several appealing implications. Starting from a collection of roots, the average depth of a node is bounded above by a constant independent of the number of files in the system [Balinska et al. 94], which seems more reasonable than the logarithmic average depth in Downey's model. The maximum depth still depends on the number of nodes. The most likely depth of a file is 0, which means it is not derived from other files. The forests themselves demonstrate preferential attachment: a forest with several nodes is more likely to produce new children. Hence, the forest sizes obey a power law, and in particular a constant fraction of the nodes are roots that have no children. These features appear realistic.

The geometric distribution of the depths is also appealing considering our results of Section 2.3. If \mathcal{D}_1 is a lognormal distribution with parameters (μ, σ^2) and \mathcal{D}_2 is a lognormal distribution with parameters (μ, σ^2) , then the results of Section 2.3 imply that the resulting distribution of file sizes is (approximately) double Pareto, since the size of a node of depth k has a lognormal distribution with parameters $((k+1)\mu, (k+1)\sigma^2)$. Indeed, we take advantage of this fact repeatedly in this section; with the above assumptions on \mathcal{D}_1 and \mathcal{D}_2 , *as long as the resulting depth distribution is geometric, the resulting file size distribution is double Pareto.*

It is clear that in this model the choice of distributions for \mathcal{D}_1 and \mathcal{D}_2 can have an important effect. If \mathcal{D}_1 and \mathcal{D}_2 are both lognormal (but do not necessarily have the same distribution), the resulting distribution is what Reed calls a double lognormal-Pareto distribution, which has properties similar to the double Pareto distribution [Reed and Jorgensen 01]. Similarly, if \mathcal{D}_1 is double Pareto or double lognormal-Pareto and \mathcal{D}_2 is lognormal, we still expect a distribution similar to the double Pareto (with Pareto tails and an approximately lognormal body).

If \mathcal{D}_2 is not lognormal, then nodes with sufficiently large depth will appear approximately lognormal (by the Central Limit Theorem argument of Section 2.2), but shallow nodes will not. The resulting distribution may therefore depend on how deep the nodes are and how quickly the product of random variables chosen from \mathcal{D}_2 converges to a lognormal distribution; however, we again emphasize that \mathcal{D}_2 does not strictly need to be lognormal for our results to hold [Reed and Hughes 02, Reed 01]. Specifically, the deepest nodes in the forest have the largest variation, and hence the small number of nodes with large depth are sufficient to yield a power law tail; the shape of the body of the distribution may be more complex. As mentioned previously, Downey's preliminary results suggest that \mathcal{D}_2 appears to be close enough to a lognormal distribution that it

quickly converges to an almost lognormal distribution after a small number of multiplicative steps, which is favorable for our analysis. Further experimental analysis and understanding of both the initial file size distribution and the multiplicative growth distribution would be an excellent starting point for future work. Also, a stronger result demonstrating the robustness of our model to deviations in the distribution of \mathcal{D}_2 would be useful, but outside the scope of this work.

4.2. Deletions

We now consider the addition of deletions to the Recursive Forest File model. Suppose at each step that a new root enters with probability γ , a file chosen uniformly at random is deleted with probability η , and a new child node is introduced as before with probability $1 - \gamma - \eta$. The introduction of deletions into the model has a surprisingly small overall effect on our previous analysis. We again give an intuitive argument for the limiting distribution, using a mean-field limit approach; these results can easily be made more rigorous using standard martingale arguments (see [Motwani and Raghavan 95]). Let $X_{t,j}$ be the number of nodes at depth j at time t , and $n(t)$ be the number of nodes at time t . It is important to clarify that the depth of a node is still computed by taking the deleted files into account, which is appropriate in our model, since the depth is meant to account for the number of modifications the file corresponding to that node has undergone. Then

$$\frac{dX_{t,0}}{dt} = \gamma - \eta \frac{X_{t,0}}{n(t)},$$

and for $j \geq 1$

$$\frac{dX_{t,j}}{dt} = (1 - \gamma - \eta) \frac{X_{t,j-1}}{n(t)} - \eta \frac{X_{t,j}}{n(t)}.$$

Now $n(t) = (1 - 2\eta)t$ in the limit as t goes large, since a node is added with probability $(1 - \eta)$ and deleted with probability η at each time-step, and inductively we can solve for the limiting values of $X_{t,j}$. The fraction of nodes at time t with depth j is then $X_{t,j}/n(t)$, and a simple induction yields $X_{t,j}/n(t) \rightarrow \gamma(1 - \gamma - \eta)^j / (1 - \eta)^{j+1}$.

Lemma 4.2. *In the Random Recursive Forest File model,*

$$\lim_{t \rightarrow \infty} \frac{X_{t,j}}{t} = \frac{\gamma(1 - 2\eta)}{(1 - \eta)} \left(\frac{1 - \gamma - \eta}{1 - \eta} \right)^j.$$

Hence the final distribution is again a geometric mixture of lognormal distributions, with the parameters slightly changed to account for deletions. As a

result, the incorporation of deletions into the model does not disrupt the resulting double Pareto distribution of file sizes.

More complex models can naturally be introduced in this framework. For example, in some situations it might be reasonable to suppose that the probability of a file being deleted is related to its depth; shallower (older) nodes may be more likely to disappear. This approach can be generalized to handle such situations, although it will affect the distribution of node depths, and again such models may be too complex to be useful.

4.3. Correlations

In our model, the file system is represented by a forest, instead of single tree. There are still correlations between file sizes; a file is still related to the size of its parent. However, the effect of these correlations is smaller, since the number of files descended from a single node is generally small compared to the size of the file system.

We can make this statement rigorous with a martingale argument. For convenience throughout, we consider the case where there are no deletions; the argument generalizes naturally.

Theorem 4.3. *Consider the Random Recursive Forest File model starting with only a single root node. For a specific value z (which may depend on n), let Z_n be the number of files with size greater than z when there are n nodes in the system. Then*

$$\Pr[|Z_n - E[Z_n]| \geq \epsilon n] \leq 2e^{-\epsilon^2 f(n)},$$

where $f(n)$ is a polynomial in n dependent on γ .

Proof. Let Y_j be the expected number of nodes with size greater than z once the first j nodes and values of the corresponding edges from their parents are revealed. (Recall that we may think of the root node of a tree as having an edge providing the size of the node.) Then $Y_0, Y_1, Y_2, \dots, Y_n$ is a martingale, with $Y_0 = E[Z_n]$ being the expected number of nodes with value at least z before any information is revealed, and $Y_n = Z_n$ being the actual number of nodes with value at least z . Let ν_j be the expected number of nodes in the subtree rooted at the j th node, where the nodes are numbered in the order of arrival (initial root nodes may be ordered arbitrarily). Notice that ν_j is independent of where the j th node is placed in the forest. Since the value of the edge corresponding to the j th node only affects the nodes in this subtree, ν_j gives an upper bound on the expected number of nodes whose final value depends on the revelation of the edge corresponding to the j th node, which implies that ν_j is an upper bound on $|Y_j - Y_{j-1}|$.

Using Azuma's inequality (see [Motwani and Raghavan 95]), we have

$$\Pr[|Y_n - Y_0| \geq \epsilon n] \leq 2e^{-\epsilon^2 n^2 / (2 \sum_{j=1}^n \nu_j^2)}.$$

Suppose that we can show that $\sum_{j=1}^n \nu_j^2$ is $O(n^{2-\zeta})$ for some $\zeta > 0$. Then we have that for any value z , the fraction of nodes with value greater than z is within ϵ of its expectation with very high probability; specifically, the probability is exponential in n^ζ . This would demonstrate that the effect of correlation is very small when looking at the ccdf.

Hence, we need an upper bound on $\sum_{j=1}^n \nu_j^2$. One approach is to simply use ν_1 as an upper bound on ν_j , so $\sum_{j=1}^n \nu_j^2 \leq n\nu_1^2$. To bound ν_1 , let $\nu_{1,k}$ be the expected number of nodes in the tree of the initial root when there are k total nodes. If we begin with a single root node, then $\nu_{1,1} = 1$ and $\nu_{1,k} = \nu_{1,k-1} \left(1 + \frac{1-\gamma}{k-1}\right)$. Using $1 + x \leq e^x$, we obtain

$$\begin{aligned} \nu_{1,n} &= \prod_{j=1}^{n-1} \left(1 + \frac{1-\gamma}{j}\right) \\ &\leq e^{(1-\gamma) \sum_{j=1}^{n-1} 1/j} \\ &= e^{(1-\gamma)(\ln n + O(1))}. \end{aligned}$$

This gives us that ν_1 is $O(n^{1-\gamma})$. This is only sufficient for Azuma's inequality if $\gamma > 1/2$, which is fairly limiting.

One way to cope with this problem is to use more initial nodes at the beginning of the process. For example, suppose that we begin with \sqrt{n} root nodes in the file system originally. The expected size of the tree rooted at any of these nodes follows the same recurrence, but now the initial condition is $\nu_{1,\sqrt{n}} = 1$. Hence,

$$\begin{aligned} \nu_{1,n} &= \prod_{j=\sqrt{n}}^{n-1} \left(1 + \frac{1-\gamma}{j}\right) \\ &\leq e^{(1-\gamma) \sum_{j=\sqrt{n}}^{n-1} 1/j} \\ &= e^{(1-\gamma)(\ln n - \ln \sqrt{n} + O(1))} \\ &= e^{(1-\gamma)(\ln n)/2 + O(1)}. \end{aligned}$$

Now for any $\gamma > 0$, ν_1 is $O(n^{(1-\gamma)/2})$, and Azuma's inequality applies.

Using the above analysis, however, we can obtain a tighter bound on ν_j , even if we begin with a single root node. Let $\nu_{j,k}$ be the expected number of nodes in the subtree of the j th node when there are k total nodes. Then $\nu_{j,j} = 1$, and

$\nu_{j,n} = \nu_j$ satisfies

$$\begin{aligned}\nu_{j,n} &= \prod_{k=j}^{n-1} \left(1 + \frac{1-\gamma}{k}\right) \\ &\leq e^{(1-\gamma) \sum_{k=j}^{n-1} 1/k} \\ &= e^{(1-\gamma) \ln(n/j) + O(1)}.\end{aligned}$$

In the above the $O(1)$ term can be taken to be independent of j for sufficiently large n . Hence, ν_j^2 is $O((n/j)^{2(1-\gamma)})$. Algebra now yields that for $\gamma < 1/2$, $\sum_{j=1}^n \nu_j^2$ is $O(n^{2(1-\gamma)})$; for $\gamma = 1/2$, $\sum_{j=1}^n \nu_j^2$ is $O(n \ln n)$; and for $\gamma > 1/2$, $\sum_{j=1}^n \nu_j^2$ is $O(n)$. In all cases, Azuma's inequality gives strong probabilistic bounds. \square

We may conclude that the fraction of node values greater than any particular value is very close to its expectation with high probability. In broader terms, the effects of correlation are small for large enough systems with small enough trees. Note that this argument demonstrates that correlation can be substantially reduced if we have more initial nodes to start the process.

Experiments using the average log ratio demonstrate that the unusual effects of correlation evident in Downey's original model do not occur in the Recursive Forest File model, as we show in Section 5.

4.4. Variations on the Derivation of New Nodes

In our dynamic Recursive Forest File model, it is again possible to consider variations on how new nodes derive from old nodes, just as it was in the recursive tree model. The variety of possibilities is rather broad, so we content ourselves here to variations of the plane-oriented recursive forest. We call this variation the Recursive Forest File model with preferential attachment. In this setting a new root is introduced at each step with probability γ ; otherwise, a new child node is introduced, and the probability that the new node is the child of a node x is proportional to $b \cdot c(x) + 1$, where $b > 0$ is a constant and $c(x)$ is the number of children of x .

Again we may begin by looking at the number of children of each depth. As before, let $X_{t,j}$ be the number of nodes of depth j at time t . Let $w(t)$, or the *weight* at time t , be the sum of $b \cdot c(x) + 1$ over all nodes. In the mean field limit, with one node added per unit time,

$$\frac{dX_{t,0}}{dt} = \gamma,$$

so that $X_{t,0}/t \rightarrow \gamma$. The case for $j \geq 1$ simplifies once we use the fact that the total number of children of nodes of depth j equals the number of nodes of depth $j + 1$. Hence, the probability of creating a child at depth j is proportional to $bX_{t,j} + X_{t,j-1}$, since this is the sum of $b \cdot c(x) + 1$ over all nodes of depth $j - 1$. Hence,

$$\frac{dX_{t,j}}{dt} = (1 - \gamma) \frac{bX_{t,j} + X_{t,j-1}}{w(t)}.$$

In the limit for large t , $w(t)$ grows to $((1 + b)(1 - \gamma) + \gamma)t$, since every new root node contributes 1 to the weight and every other node contributes $1 + b$. Now if $X_{t,j}/t$ approaches x_j asymptotically, we find from the above that

$$x_j = (1 - \gamma) \frac{bx_j + x_{j-1}}{((1 + b)(1 - \gamma) + \gamma)}.$$

Simplifying the above yields

$$x_j = (1 - \gamma)x_{j-1},$$

so a simple induction again yields $X_{t,j}/t \rightarrow \gamma(1 - \gamma)^j$.

Lemma 4.4. *In the Random Recursive Forest File model with preferential attachment,*

$$\lim_{t \rightarrow \infty} \frac{X_{t,j}}{t} = \gamma(1 - \gamma)^j.$$

Surprisingly, this is the same result as in the Random Recursive Forest File model, regardless of the value of b !

The value of b therefore does not affect the resulting geometric distribution of the depths of the nodes, and hence the double Pareto analysis still applies. We believe this demonstrates substantial robustness for this model in the face of changes.

The value of b does affect the model, however, in how the nodes are distributed among the trees in the forest. As a concrete example, comparing the uniform case ($b = 0$) with the plane-oriented recursive forest model ($b = 1$), we find for the larger b value that there are a substantially greater number of trees consisting of just a single vertex and there is greater variance in the number of offspring from a root node. Hence, the choice of b might be used to fine-tune the underlying model to various file systems.

To see how b affects the distribution of the size of trees in the forest, we again describe an asymptotic mean field argument. Let $Y_{t,j}$ be the number of trees with j nodes at time t . Note that the total weight corresponding to a tree with

j nodes is $(j - 1)b + j$, since every node contributes $1 + b$ to the weight except for the root. Hence, we obtain the following equations:

$$\frac{dY_{t,1}}{dt} = \gamma - (1 - \gamma) \frac{Y_{t,1}}{w(t)},$$

and for $j \geq 2$,

$$\frac{dY_{t,j}}{dt} = (1 - \gamma) \frac{Y_{t,j-1}((j - 2)b + j - 1) - Y_{t,j}((j - 1)b + j)}{w(t)}.$$

The asymptotic behavior of this system is easy to solve for, and the distribution of tree sizes in the forest follows a power law with the exponent in the power law depending on b [Drinea et al. 01, Krapivsky and Redner 01].

We also note that a similar derivation shows that the distribution of the depths of the nodes remains geometric under these variations when random deletions occur as in Section 4.2.

5. Simulations

In this section, we examine simulations using the Recursive Forest File model to compare it to the theory. In particular, we examine the issues of correlation and convergence to the limiting depth distribution. Overall, we find that simulations match the theory well. Rather than compare with actual data sets, we refer the reader to [Mitzenmacher and Tworetzky 03] for a detailed evaluation.

Consider first the problem of correlation. Recall that we simulated Downey's model by placing weights chosen from a normal distribution with mean 0 and variance 1 on each edge. In 1,000 runs of generating 10,000 files, the average log ratio varied between -4.2 and 5.2 . We repeated the experiment using our dynamic model with $\gamma = 0.1$; also, the original size of each root node is lognormally distributed, so that there is an implicit edge with mean 0 and variance 1 into each root. Starting initially with 1 root node, the average log ratio varied between -2.26 and 2.52 ; starting with 10 root nodes, it varied between -0.95 and 1.22 ; and starting with 100 root nodes, it varied between -0.38 and 0.49 . While it is clear that there are still correlations in the file sizes, they are dramatically reduced over Downey's model. Similarly, increasing the number of files leads to sharper concentration of the average log ratio, as our analysis would predict.

A second issue is convergence in the depth distribution. While asymptotically the depths will converge to a geometric distribution, it is not clear how many files are necessary for this to occur, especially if one starts with multiple roots at the beginning. Indeed, we find that the convergence in the depth distribution is slow, but it does not dramatically change the characteristics of the distribution shapes produced.

A representative example is instructive. We generated sets with 10,000 and 100,000 nodes, using $\gamma = 0.1$ and beginning with 1 and 100 initial roots. The results are presented in Figure 4. The resulting distribution does not match the theoretical geometric distribution; there is a bump in the distribution depending on the number of nodes generated and the number of initial roots. The more nodes generated, the closer to equilibrium.

Despite this deviation from the theory, examining plots on a log-log scale reveals that the cdf and the ccdf of the file sizes generated by the Recursive Forest File model still have essentially linear bodies and tails, as shown in Figures 5 and 6. The deviation of the model from the theoretical double Pareto distribution appears to add a small curvature to the distribution. Also, the linear tails break down somewhat at the extremes, because of the small number of samples and because the distribution has not reached the theoretical equilibrium. Indeed, part of the argument for using a lognormal distribution over a Pareto distribution in previous work has been the curvature at the tail of the distribution [Downey 01]. The Recursive Forest File model demonstrates that this curvature could be arising simply because the snapshot of the dynamically changing distribution is taken at some specific finite point in time, before the long-term equilibrium has been reached.

To test that our results hold even when the multiplicative distribution is not lognormal, we have performed similar simulations using other multiplicative distributions, for example, using edge weights that are 1 and -1 with probability $1/2$ for all edges (except those into the root, so that there is some asymmetry). This distribution yielded entirely similar curves for 100,000 nodes.

6. Conclusions

We have provided and analyzed a new generative user model, the Recursive Forest File model, for file size distributions. Understanding the behavior of file size distributions is an important building block for understanding both file systems and Internet behavior. Our model is extremely simple and well suited for simulation tools.

The underlying idea behind the model is to combine a multiplicative generating process with a dynamic insertion and deletion process reminiscent of recent web graph models. A fundamental point in the analysis is to connect the file size model with corresponding random tree and forest models. We have shown that for many natural model variations the depth distributions are asymptotically geometrically distributed, and this in turn yields a double Pareto distribution for the file sizes.

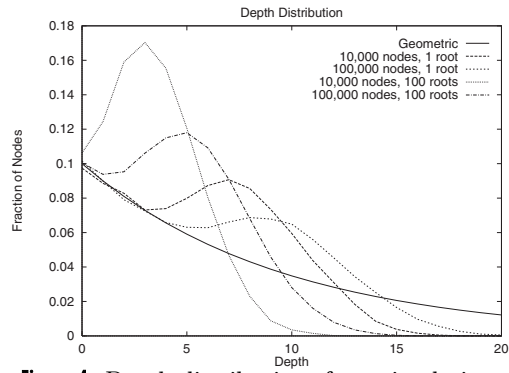


Figure 4. Depth distributions from simulations.

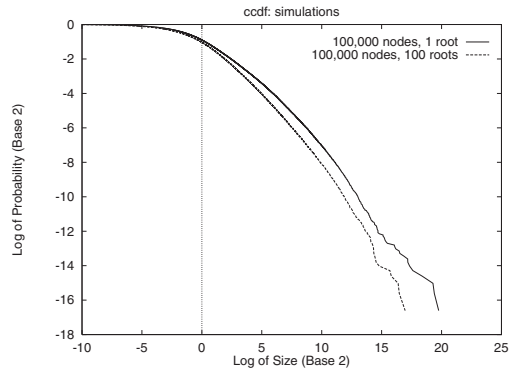


Figure 5. ccdfs for the simulations.

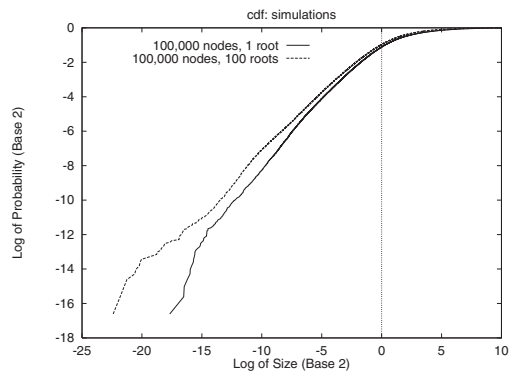


Figure 6. cdfs for the simulations.

From a practical standpoint, this model explains why file size distributions may appear to have a lognormal body and a Pareto tail. (In fairness, we point out that the shape of these distributions is still a subject of debate.) While previous work has suggested using specific hybrid distributions to model file sizes, our generative model appears sufficiently accurate, and has the advantage that it can be used to simulate dynamic systems where files may change over times. An open question for future work is how to design tools to fit properly parametrized double Pareto (or double Pareto-lognormal) distributions to empirically observed distributions.

From a theoretical standpoint, a Recursive Forest model provides a general mechanism for producing power law distributions that may apply to other natural systems. The robustness of the model to deletions and to changes in how elements produce offspring appears to be an extremely appealing feature. The flexibility and simplicity of the random graph framework should allow for further variations worthy of study.

There remain many open problems to pursue. On the practical side, there does not appear to be experimental work that considers how files change or are generated over time. Such data might validate this model or lead to other dynamic models for file sizes. Specifically, understanding how files are created and deleted over time, knowing the distribution of file sizes when they are created, and determining whether modifications truly lead to multiplicative changes in the file size would be useful information for studying the dynamic behavior of file systems. Dynamic traces covering long time spans are important for further research in this area.

On the theoretical side, perhaps the most interesting question is the rate of convergence to the double Pareto distribution. Our simulations have shown that it takes significant time for the node depths to converge to a geometric distribution. The general shape of the corresponding file size distribution does not seem to change significantly, however; the major difference appears to be that the distribution appears more like a lognormal distribution, in that the tail dies off somewhat more quickly than expected. It is an open problem to formalize these findings theoretically. Another issue is to provide a better understanding of the sensitivity of the Recursive Forest File model to the underlying distributions \mathcal{D}_1 and \mathcal{D}_2 . Finally, determining alternative generative models that could justify a lognormal distribution or another distribution for file sizes could lead to new debates on the appropriateness of various models for file sizes.

Acknowledgments. The author would like to thank Laura Serban for help coding simulations; Steve Lumetta for help processing HTTP logs; Mark Crovella and Paul Barford for making their data available; and Mark Crovella, John Byers, Steve Lumetta, and

Allen Downey for helpful comments. The author was supported in part by an Alfred P. Sloan Research Fellowship and NSF grants CCR-9983832, CCR-0118701, and CCR-0121154.

References

- [Aitchison and Brown 57] J. Aitchison and J. A. C. Brown. *The Lognormal Distribution*. Cambridge, UK: Cambridge University Press, 1957.
- [Arlitt and Williamson 96] M. F. Arlitt and C. L. Williamson. “Web Server Workload Characterization: The Search for Invariants.” *Performance Evaluation Review* 24:1 (1996), 125–137.
- [Balinska et al. 94] K. T. Balinska, L. V. Quintasem, and J. Szymański. “Random Recursive Forests.” *Random Structures and Algorithms* 5 (1994), 3–12.
- [Barabási et al. 99] A. -L. Barabási, R. Albert, and H. Jeong. “Mean-Field Theory for Scale-Free Random Networks.” *Physica A* 272 (1999), 173–189.
- [Barford and Crovella 98] P. Barford and M. Crovella. “Generating Representative Web Workloads for Network and Server Performance Evaluation.” *Performance Evaluation Review* 26:1 (1998), 151–160.
- [Barford et al. 99] P. Barford, A. Bestavros, A. Bradley, and M. Crovella. “Changes in Web Client Access Patterns: Characteristics and Caching Implications.” *World Wide Web* 2 (1999), 15–28.
- [Broder et al. 00] A. Broder, R. Kumar, F. Maghoul, P. Raghavan, S. Rajagopalan, R. Stata, A. Tomkins, and J. Wiener. “Graph Structure in the Web: Experiments and Models.” In *Proceedings of the 9th World Wide Web Conference*. Available in *Computer Networks* 33:1–6 (2000), 309–320.
- [Carlson and Doyle 99] J. M. Carlson and J. Doyle. “Highly Optimized Tolerance: A Mechanism for Power Laws in Designed Systems.” *Physics Review E* 60:2 (1999), 1412–1427.
- [Crovella and Bestavros 97] M. Crovella and A. Bestavros. “Self-Similarity in World Wide Web Traffic: Evidence and Possible Causes.” *IEEE/ACM Transactions on Networking* 5:6 (1997), 835–846.
- [Crovella et al. 98] M. Crovella, M. S. Taqqu, and A. Bestavros. “Heavy-Tailed Probability Distributions in the World Wide Web.” In *A Practical Guide to Heavy Tails*, edited by R. J. Adler, R. E. Feldman, and M. S. Taqqu, pp. 3–26. London: Chapman and Hall, 1998.
- [Crow and Shimura 88] E. L. Crow and K. Shimizu, editors. *Lognormal Distributions: Theory and Applications*. New York: Marcel Dekker, Inc., New York, 1988.
- [Devroye 98] L. Devroye. “Branching Processes and Their Applications in the Analysis of Tree Structures and Tree Algorithms.” In *Probabilistic Methods for Algorithmic Discrete Mathematics*, edited by M. Habib, C. McDiarmid, J. Ramirez-Alfonsin, and B. Reed, pp. 249–314. Berlin: Springer-Verlag, 1998.

- [Downey 01] A. B. Downey. “The Structural Causes of File Size Distributions.” In *Proceedings of the Ninth International Symposium on Modeling, Analysis, and Simulation of Computer and Telecommunication Systems*, pp. 361–370. Los Alamitos, CA: IEEE Computer Society, 2001.
- [Drinea et al. 01] E. Drinea, M. Enachescu, and M. Mitzenmacher. “Variations on Random Graph Models of the Web.” Harvard Computer Science Technical Report TR-06-01, 2001.
- [Fabrikant et al. 02] A. Fabrikant, E. Koutsoupias, and C. Papadimitriou. “Heuristically Optimized Trade-Offs: A New Paradigm for Power Laws on the Internet.” In *Proceedings of the Twenty-Ninth International Colloquium on Automata, Languages, and Programming*, pp. 110–122, Lecture Notes in Computer Science 2380. Berlin: Springer-Verlag, 2002.
- [Gabaiz 99] X. Gabaix. “Zipf’s Law for Cities: An Explanation.” *Quarterly Journal of Economics* 114 (1999), 739–767.
- [Huberman and Adamic 99] B. A. Huberman and L. A. Adamic. “Evolutionary Dynamics of the World Wide Web.” Technical Report, Xerox Palo Alto Research Center, 1999. Appears as a brief communication in *Nature* 399 (1999), 130.
- [Huberman and Adamic 00] B. A. Huberman and L. A. Adamic. “The Nature of Markets in the World Wide Web.” *Quarterly Journal of Economic Commerce* 1 (2000), 5–12.
- [Kleinberg et al. 99] J. Kleinberg, R. Kumar, P. Raghavan, S. Rajagopalan, and A. Tomkins. “The Web as a Graph: Measurements, Models, and Methods.” In *Proceedings of the 5th International Conference on Combinatorics and Computing*, pp. 1–17, Lecture Notes in Computer Science 1627. Berlin: Springer-Verlag, 1999.
- [Krapivsky and Redner 01] P. L. Krapivsky and S. Redner. “Organization of Growing Random Networks.” *Physical Review E* 63 (2001), 066123-1–066123-14.
- [Kumar et al. 00] R. Kumar, P. Raghavan, S. Rajagopalan, D. Sivakumar, A. Tomkins, and E. Upfal. Stochastic Models for the Web Graph.” In *Proceedings of the 41st Annual Symposium on Foundations of Computer Science*, pp. 57–65. Los Alamitos, CA: IEEE Computer Society, 2000.
- [Kurtz 81] T. G. Kurtz. *Approximation of Population Processes*. CBMS-NSF Regional Conference Series in Applied Mathematics, 36. Philadelphia: SIAM, 1981.
- [Leland et al. 94] W. E. Leland, M. S. Taqqu, W. Willinger, and D. V. Wilson. “On the Self-Similar Nature of Ethernet Traffic.” *IEEE/ACM Transactions on Networking* (1994), 1–15.
- [Li 99] W. Li. “References on Zipf’s Law.” Available from World Wide Web (<http://linkage.rockerfeller.edu/wli/zipf/>), 1999.
- [Mitzenmacher 04] M. Mitzenmacher. “A Brief History of Generative Models for Power Law and Lognormal Distributions.” *Internet Mathematics* 1:2 (2004), 226–251.
- [Mitzenmacher and Tworetzky 03] M. Mitzenmacher and B. Tworetzky. “New Models and Methods for File Size Distributions.” In *Proceedings of the 41st Annual Allerton Conference on Communication, Control, and Computing*, pp. 603–612. Urbana, IL: University of Illinois at Urbana-Champaign, 2003.

- [Montroll and Shlesinger 83] E. W. Montroll and M. F. Shlesinger. “Maximum Entropy Formalism, Fractals, Scaling Phenomena, and $1/f$ Noise: A Tale of Tails.” *Journal of Statistical Physics* 32 (1983), 209–230.
- [Motwani and Raghavan 95] R. Motwani and P. Raghavan. *Randomized Algorithms*. Cambridge, UK: Cambridge University Press, 1995.
- [Reed 01] W. J. Reed. Personal communication, 2001.
- [Reed 03] W. J. Reed. “The Pareto Law of Incomes - An Explanation and an Extension.” *Physica A* 319 (2003), 469–485.
- [Reed and Jorgensen 01] W. J. Reed and M. Jorgensen. “The Double Pareto-Lognormal Distribution - A New Parametric Model for Size Distributions.” To appear in *Communications in Statistics: Theory and Methods* 33:8 (2004).
- [Reed and Hughes 02] W. J. Reed and B. D. Hughes. “From Gene Families and Genera to Incomes and Internet File Sizes: Why Power-Laws Are So Common in Nature.” *Physical Review E* 66 (2002), 067103.
- [Smythe and Mahmoud 95] R. Smythe and H. Mahmoud. “A Survey of Recursive Trees.” *Theoretical Probability and Mathematical Statistics* 51 (1995), 1–27.
- [Wormald 95] N. C. Wormald. “Differential Equations for Random Processes and Random Graphs.” *Annals of Applied Probability* 5:4 (1995), 1217–1235.
- [Zhu et al. 01] X. Zhu, J. Yu, and J. Doyle. “Heavy Tails, Generalized Coding, and Optimal Web Layout.” In *Proceedings of IEEE INFOCOM*, pp. 1617–1626. New York: IEEE Computer and Communications Society, 2001.

Michael Mitzenmacher, Harvard University, Division of Engineering and Applied Science, 33 Oxford Street, Cambridge, MA 02138 (michaelm@eecs.harvard.edu)

Received April 22, 2003; accepted August 4, 2003.