

1 Fully Specified Estimation Plan for
2 Optimal Static Parametric Estimation of Arbitrary Distributions
3 (OSPEAD)

4 Rucknium

5 September 21, 2022

6 Delivered to the OSPEAD Scientific Review Panel for Milestone 1 of the OSPEAD CCS proposal

7 **1 Introduction**

8 Research on Monero and similar systems has established the vital importance of constructing a decoy selection
9 algorithm that closely resembles the real spend age distribution. A sample of researchers' conclusions and recom-
10 mendations follows:

11 As we have seen the current [2017] sampling strategy for mix-ins fails drastically in preventing temporal
12 analysis. There are two possible strategies towards mitigating the ensuing risks: (a) mimic users'
13 spending behavior or, (b) force mix-ins to be picked according to some "unknown" distribution.

14 - [Kumar et al., 2017]

15 We have provided evidence that the strategy by which Monero mixins are sampled [circa 2017] results in
16 a different time distribution than real spends, significantly undermining the untraceability mechanism.
17 To correct this problem, the mixins should ideally be sampled in such a way that the resulting time
18 distributions match.

19 - [Möser et al., 2018]

20 The mixin sampling distribution has since been replaced with a gamma distribution (from [Möser et al., 2018])
21 fitted to the empirical spend-time distribution. Our results show that these sampling distribution changes
22 have made a significant impact in reducing the accuracy of the guess-newest heuristic.

23 - [Ye et al., 2020]

24 [A mimicking decoy selection algorithm's] anonymity depends on how well $\hat{\mathcal{S}}$ [the decoy selection al-
25 gorithm] estimates \mathcal{S} [the real spend age distribution]....It is therefore reasonable to expect that if the
26 mimicking sampler has access to the true source distribution \mathcal{S} , its anonymity should be close to optimal.
27 In the following, we give an [sic] evidence that this is the case.

28 - [Ronge et al., 2021]

29 The problem of creating a decoy selection algorithm that minimizes the usefulness of timing information to an
30 adversary was recognized by the Monero Research Lab very early [Mackenzie et al., 2015]:

31 One solution to this problem is to determine a non-uniform method of choosing transaction outputs
32 for ring signatures; choose transaction outputs based on their age such that the probability that these
33 outputs are chosen for a ring signature is inversely related to the probability that they have been
34 spent already. This would suggest that we, the developers of Monero, must estimate the probability
35 distribution governing the age of transaction outputs.

36 A reliable estimator of the real spend age probability distribution has been elusive. [Möser et al., 2018] estimated
 37 the distribution based on partial knowledge of real spends obtained from other de-anonymization techniques. The
 38 introduction of RingCT and other improvements to Monero have rendered the de-anonymization techniques of
 39 [Möser et al., 2018] ineffective ([Ye et al., 2020], [Vijayakumaran, 2021]). Therefore, an updated and improved es-
 40 timate would require an estimator to use only the fully anonymized ring data on the Monero blockchain.

41 In this document I lay out such an estimator and justify it in a fairly rigorous manner.

42 Part I

43 Undisclosed Portion

44 2 Overview of the plan

45 OSPEAD is a two-part statistical technique. The two parts can each be separated into two steps. The first part
 46 produces an estimate of the historical real spend age distribution. The first part is potentially useful to an adversary
 47 who wishes to harm Monero users' privacy, so a consideration of disclosure is necessary before releasing it publicly.
 48 The second part takes the real spend age distribution as an input and produces a usable parametric distribution
 49 for Monero's decoy selection algorithm, taking into account variability in the real spend age distribution over time.
 50 The second part can be released publicly. Therefore, I intend to detach the second part of this document and release
 51 it publicly so the community can examine progress of OSPEAD and give feedback.

52 The two steps of the first part are:

- 53 1. Estimate the sub-distributions formed by the diverse decoy selection algorithms that various wallet software
 54 use.
- 55 2. Use the sub-distribution associated with the standard `wallet2` decoy selection algorithm to estimate the real
 56 spend age distribution.

57 The two steps of the second part are:

- 58 1. Fit a parametric distribution to the real spend age distribution estimated in part 1 by using a loss function
 59 minimization approach.
- 60 2. Use the time variability of the historical real spend age distribution and forecast error risk as feedback into a
 61 final parametric decoy selection algorithm.

62 The estimation procedure I propose in this document follows the outline in my HackerOne submission closely, with
 63 two main exceptions:

64 In the submission I noted “the fact that $f_S(x)$ [the real spend age distribution] is literally a moving target” since
 65 it evolves over time, but I did not propose any solution to this problem. In this document I propose some forecast
 66 validation methods.

67 In my HackerOne submission I suggested that a Maximum Likelihood Estimator could be developed to handle
 68 the multiple decoy selection algorithms (DSAs), modeling them as a mixture distribution. In this document I take
 69 a different path. There is still the notion of DSAs forming a mixture distribution, but with my proposed estimator
 70 we do not need to know the form of all the DSAs being used in the wild, as would have been required for my
 71 original idea. My new proposed approach explicitly leverages the hierarchical grouped ring structure of the data,

72 which provides much greater ability to identify and estimate all the necessary statistical objects. I am not sure if
 73 my original model is even identifiable.¹

74 I only considered methods within the frequentist statistical paradigm rather than the Bayesian paradigm since
 75 I am not very familiar with Bayesian methods.

76 3 The Meaning of OSPEAD

77 OSPEAD is an acronym for Optimal Static Parametric Estimation of Arbitrary Distributions. “Estimation of
 78 Arbitrary Distributions” is the main term. “Static” and “Parametric” are qualifiers for the “Optimal” descriptor.
 79 The qualifiers are used in the same way that Ordinary Least Squares (OLS) is a Best Linear Unbiased Estimator
 80 (BLUE) of the conditional mean of a random variable, as stated by the Gauss-Markov Theorem. BLUE means that
 81 among estimators that are both linear and unbiased, OLS is the “best”, i.e. has the lowest variance, of any of them.
 82 I will work backwards through the terms:

83 3.1 “Estimation of Arbitrary Distributions”

84 OSPEAD is an estimator of arbitrary distributions. In this case, the distributions are the weekly real spend age
 85 distributions of transaction outputs on the Monero blockchain. I say “arbitrary” here since I do not impose strict
 86 assumptions on the form of the distributions.

87 3.2 “Parametric”

88 OSPEAD is parametric in the sense that the final decoy selection algorithm will be defined by a parametric rather
 89 than nonparametric probability density function. The current decoy selection algorithm is formed by a parametric
 90 probability density function. It may be the case that a nonparametric probability density function would be better
 91 than a parametric one, but the advisability and feasibility of a nonparametric one is left for future research. This
 92 is consistent with a Monero development principle: “The Monero Core Team and the Monero Research Lab would
 93 like to follow the development philosophy that it is wise to start with smaller changes at first and then ramp those
 94 changes up over time, rather than start with drastic changes and try to scale them back.” [Mackenzie et al., 2015]
 95 Even though the final form of the OSPEAD decoy selection algorithm will be parametric, there are intermediate
 96 steps that will involve nonparametric estimation.

97 3.3 “Static”

98 OSPEAD is static rather than dynamic. OSPEAD will recommend that wallet software draws ring members in
 99 the same way regardless of any shifts in the blockchain data over time, as it does currently. This is consistent
 100 with the cautious development philosophy of Monero. A dynamic decoy selection algorithm could re-estimate the
 101 “current” real spend age distribution at any point in time on-the-fly. The advisability and feasible of a dynamic
 102 decoy selection algorithm would require further research. However, OSPEAD will take into account the historical
 103 variability of the real spend age distribution over time and will attempt to produce good performance despite the
 104 fact that the real spend age distribution will change in the future.

¹In the $f(x) = (1 - \alpha)f_S(x) + \sum_{i=1}^N \alpha_i f_{D,i}(x)$ model with $\sum_{i=1}^N \alpha_i = \alpha$ and each $f_{D,i}(x)$ known, it could be true that the mixing proportions α_i are identifiable by the argument of [Hall, 1981] and the unknown nonparametric $f_S(x)$ might be identifiable by the argument of [Patra & Sen, 2016]. Then the whole model might be identifiable. After finding the literature on repeated measurements in finite mixture models, I decided not to pursue this model further.

105 **3.4 “Optimal”**

106 I use the term “optimal” somewhat loosely. I do not claim that my proposed estimator of the real spend age
 107 distribution has the lowest variance of any possible estimator. I argue that my proposed estimator is consistent (i.e.
 108 approaches the true value as the number of rings in the sample rises) and has acceptable variance. See **Section**
 109 **14 Future Work to Improve the Two-Step Estimator** for a discussion about how my proposed estimator
 110 could be improved. The parametric fitting part is optimal in the sense that I search over a large space of possible
 111 parametric distributions and their parameters values to minimize (optimize) certain specified loss functions.

112 If we have a near-optimal estimate and parametric fit of the real spend age distribution, do we have near-optimal
 113 protection of Monero user privacy, at least regarding timing-based statistical attacks? [Ronge et al., 2021] argue
 114 “yes”:

115 Furthermore, in Section 5 we show that the definition is robust, in the sense that the anonymity changes
 116 only slightly when the signer distribution changes slightly. In particular, if a ring sampler is shown to
 117 be good with respect to a close estimate $\hat{\mathcal{S}}$ of the real signer distribution \mathcal{S} , then it should also be good
 118 with respect to the real signer distribution \mathcal{S}

119 Later in this work, we will show that for 1-signer distributions \mathcal{S} [i.e. only a single real spend within
 120 a ring, which is the case for Monero,] the optimal anonymity is always almost achievable. More con-
 121 cretely, in Section 6.2 we show that there exists a “mimicking” sampler Π_{Mimic} which achieves anonymity
 122 $\alpha(\mathcal{S}, \Pi_{\text{Mimic}}) \gtrapprox \frac{1}{2} \lg n$, which is only a constant fraction away from the optimum, assuming minimally
 123 that \mathcal{S} has at least $\lg n$ bits of min-entropy. Although the mimicking sampler achieves near-optimal
 124 anonymity, it is mostly theoretical as it requires the knowledge of the distribution \mathcal{S} .

125 Of course, I disagree with the statement of [Ronge et al., 2021] that this fact is only of theoretical interest. A good
 126 estimate of \mathcal{S} is feasible.

127 **4 What to Read and What to Skip**

128 This document is lengthy. It is not strictly necessary to read certain sections to arrive at a judgment of what
 129 estimation strategy OSPEAD should use if the reader has a very limited time budget. An abridged version of this
 130 document would be:

- 131 • Section 5 What to Decide
- 132 • Section 7 Finite Mixture Models
- 133 • Section 8 Modeling the Real Spend and Decoy Distributions as a Mixture
- 134 • Section 10 First Step: Bonhomme-Jochmans-Robin Estimator
- 135 • Section 11 Second Step: Patra-Sen Inversion Estimator
- 136 • Section 12.4 Options for Confidence Interval Estimation
- 137 • Section 21 Criteria for Best Fit
- 138 • Section 23 Options for Loss Functions \mathcal{L} and Parametric Distributions \mathcal{F}
- 139 • Section 24 Dynamic Risk and Forecasting
- 140 • Section 26 Options for Forecasting

141 These section can be skipped if desired:

- Section 6 Foundational Statistical Concepts, but I would still recommend reading Section 6.4 Consistency and Section 6.5 Identifiability and Identification if unfamiliar with those concepts
- Section 9 Two-Step Estimation
- Section 12 Confidence Interval Estimation
- Section 13 Sampling Strategy and Data Features that Support Estimation
- Section 14 Future Work to Improve the Two-Step Estimator
- Section 15 Rucknium Ratio Attack/MAP Decoder/Discriminant Analysis
- Section 16 Postestimation Classification of On-Chain Rings into DSAs
- Section 17 Disclosure Considerations
- Section 18 OSPEAD Requests to Monero C++ Developer(s)
- Section 19 Plain English Explanation of the Problem
- Section 20 The Solution
- Section 22 Dry Run with Old M oser et al. (2018) Data
- Section 25 Inter-Temporal Stability of the Spent Output Age Distribution for BTC, BCH, LTC, and DOGE
- Section 27 Acknowledgments

5 What to Decide

The purpose of this document is to give the OSPEAD scientific review panel an opportunity to give input on the estimation plan. The estimation is designed to give a range of candidates for the new decoy selection algorithm, so it is expected that multiple estimation options can be done at this point in time. The OSPEAD CCS has two further milestones: “Deliver initial probability density function to scientific review panel” and “Deliver final version of probability density function to Monero developers.” So we can do the estimations, see what we get, and then decide on the winning candidate. I have distilled the choices into these sections:

- Section 8 Modeling the Real Spend and Decoy Distributions as a Mixture, regarding whether Equation (7) should be used as the OSPEAD decoy selection algorithm.
- Section 10.5 Options for the First Step Estimator
- Section 11.1 Options for the Second Step Estimator
- Section 12.4 Options for Confidence Interval Estimation
- Section 23 Options for Loss Functions \mathcal{L} and Parametric Distributions \mathcal{F}
- Section 26 Options for Forecasting

It would be good to get early feedback on whether the Bonhomme-Jochmans-Robin Estimator in **Section 10.5 Options for the First Step Estimator** should be used so that work can begin on writing a C++ implementation of it (point 2 of **Section 18 OSPEAD Requests to Monero C++ Developer(s)**).

¹⁷⁴ 6 Foundational Statistical Concepts

¹⁷⁵ Here I will review some statistical concepts that will be needed to compare different estimators and describe the
¹⁷⁶ assumptions required for their validity.

¹⁷⁷ 6.1 Cumulative Distribution Function, Probability Density Function, Probability ¹⁷⁸ Mass Function, and Their Support

¹⁷⁹ The following are some definitions of statistical objects, from [Casella & Berger, 2002]:

¹⁸⁰ **Definition 1.5.1** The *cumulative distribution function* or *cdf* of a random variable X , denoted $F_X(x)$,
¹⁸¹ is defined by

$$F_X(x) = P_X(X \leq x), \quad \text{for all } x.$$

¹⁸² $P_X(X \leq x)$ is the probability that a single value of the random variable X is below x . A CDF monotonically
¹⁸³ rises from 0 to 1. All random variables have a CDF regardless of whether they are continuous or discrete.
¹⁸⁴ Therefore, CDFs are often used when statistical theory attempts to be more general. Another definition from
¹⁸⁵ [Casella & Berger, 2002]:

¹⁸⁶ **Definition 1.5.7** A random variable X is *continuous* if $F_X(x)$ is a continuous function of x . A random
¹⁸⁷ variable X is *discrete* if $F_X(x)$ is a step function of x .

¹⁸⁸ The actual concept of time is continuous. Monero transactions are measured in discrete time because they are
¹⁸⁹ confirmed on the block chain every two minutes on average. Therefore their exact time distribution is governed by
¹⁹⁰ a probability mass function:

¹⁹¹ **Definition 1.6.1** The *probability mass function* (*pmf*) of a discrete random variable X is given by

$$f_X(x) = P(X = x) \quad \text{for all } x.$$

¹⁹² The value of a PMF at any point x is the probability that the random variable X will have value x .

¹⁹³ **Definition 1.6.3** The *probability density function* or *pdf*, $f_X(x)$, of a continuous random variable X is
¹⁹⁴ the function that satisfies

$$F_X(x) = \int_{-\infty}^x f_X(t) dt \quad \text{for all } x.$$

¹⁹⁵ In essence the PDF is the first derivative of the CDF. Even though ring member ages are discrete due to the
¹⁹⁶ discrete nature of blocks, continuous PDFs can be useful in approximating distributions for creating decoy selection
¹⁹⁷ algorithms. [Casella & Berger, 2002] continues:

¹⁹⁸ [Let] $\mathcal{X} = \{x : f_X(x) > 0\}$. The pdf of random variable X is positive only on the set \mathcal{X} and is 0
¹⁹⁹ elsewhere. Such a set is called the *support set* of a distribution or, more informally, the *support* of a
²⁰⁰ distribution. This terminology can also apply to a pmf or, in general, to any nonnegative function.

²⁰¹ Support is an important factor to consider when attempting to distinguish different distributions and different
²⁰² random variables from each other. Intuitively, if the support of two variables overlaps either partially or completely,
²⁰³ then it will be more difficult to distinguish real data being generated by them.

204 6.2 Estimator, Estimate, and Estimand

205 OSPEAD is an estimator. [Casella & Berger, 2002] give a very general definition of a point estimator, which is the
 206 most common type of estimator. First their definition of “statistic”:

207 **Definition 5.2.1** Let X_1, \dots, X_n be a random sample of size n from a population and let $T(x_1, \dots, x_n)$
 208 be a real-valued or vector-valued function whose domain includes the sample space of (X_1, \dots, X_n) .
 209 Then the random variable or random vector $Y = T(X_1, \dots, X_n)$ is called a *statistic*. The probability
 210 distribution of a statistic Y is called the *sampling distribution of Y* .

211 [...]

212 **Definition 7.1.1** A *point estimator* is any function $W(X_1, \dots, X_n)$ of a sample; that is, any statistic
 213 is a point estimator.

214 Notice that the definition makes no mention of any correspondence between the estimator and the
 215 parameter it is to estimate. While it might be argued that such a statement should be included in the
 216 definition, such a statement would restrict the available set of estimators.. Also, there is no mention in
 217 the definition of the range of the statistic $W(X_1, \dots, X_n)$. While, in principle, the range of the statistic
 218 should coincide with that of the parameter, we will see that this is not always the case.

219 There is one distinction that must be made clear, the difference between an estimate and an estimator.
 220 An *estimator* is a function of the sample, while an *estimate* is the realized value of an estimator (that
 221 is, a number) that is obtained when a sample is actually taken. Notationally, when a sample is taken,
 222 an estimator is a function of the random variables X_1, \dots, X_n , while an estimate is a function of the
 223 realized values x_1, \dots, x_n .

224 Besides point estimators, there are interval estimators (e.g. estimates of confidence intervals), estimators of prob-
 225 ability density functions, and so forth. Generally the purpose of an estimator is to generate a close approximation
 226 of the *estimand*, i.e. the object to be estimated, from the available data in the sample.

227 6.3 Parametric and Nonparametric Models

228 The distinction between parametric and nonparametric statistical models is crucial for understanding the problem
 229 that OSPEAD aims to solve. The Wikipedia entry for “Parametric statistics” adequately explains the distinction:²

230 **Parametric statistics** is a branch of statistics which assumes that sample data comes from a popu-
 231 lation that can be adequately modeled by a probability distribution that has a fixed set of parameters.
 232 Conversely a **non-parametric model** does not assume an explicit (finite-parametric) mathematical
 233 form for the distribution when modeling the data. However, it may make some assumptions about that
 234 distribution, such as continuity or symmetry.

235 A normal distribution, for example, is completely specified once its mean μ and variance σ^2 parameter are set. At
 236 every point x , its probability density is then

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

237 If a statistician *assumed* that a given random variable had a normal distribution, then the statistician would
 238 just need to estimate the distribution’s mean and variance to completely describe the random variable. Assuming
 239 that variables are distributed in a particular parametric way is done for convenience, computational and theoretical
 240 tractability, and sample size reasons.

241 Parametric distributions often can describe an idealized physical or social process. A normal distribution is the
 242 limiting distribution of the sum of an infinite number of random variables. An exponential distribution describes the

²https://en.wikipedia.org/wiki/Parametric_statistics

time interval between events for memoryless processes, like proof-of-work block discovery. A binomial distribution defines the probability of n successful Bernoulli trials with probability p . And so forth.

The real spend age distribution is fundamentally an economic process determined by human decision making. The point of OSPEAD is not to develop a strong theoretical economic model for Monero users' behavior. I will let the data speak for itself rather than asking a theoretical model to speak for it. There is, however, an implicit theoretical model: A current or potential Monero user goes about their daily life. At some point the neurons in the user's brain form into a configuration that is interpreted as "I think I will spend some Monero now." The user then broadcasts a Monero transaction. A miner includes it in (usually) the next block. At any given point in time for each Monero user on the planet there is some unknown latent probability that their neurons will attain the "spend Monero now" configuration. That probability may be different for each user — and within each user the probability may change over time.

There is no particular reason to think that the real spend age distribution conforms to a specific tidy parametric distribution. We do not really know why users spend when they spend, nor do we know which "types" of users tend to spend more (or less) often at particular times. Therefore, I take it as given that the real spend age distribution is nonparametric. A basic example of a nonparametric estimator of a distribution is a histogram. A kernel density estimator, which will appear later, is basically a more sophisticated histogram.

6.4 Consistency

Consistency is essentially a requirement of an adequate estimator. [Casella & Berger, 2002] explain:

The property of consistency seems to be quite a fundamental one, requiring that the estimator converges to the "correct" value as the sample size becomes infinite. It is such a fundamental property that the worth of an inconsistent estimator should be questioned (or at least vigorously investigated).

Consistency (as well as all asymptotic properties) concerns a sequence of estimators rather than a single estimator, although it is common to speak of a "consistent estimator." If we observe X_1, X_2, \dots according to a distribution $f(x|\theta)$, we can construct a sequence of estimators $W_n = W_n(X_1, \dots, X_n)$ merely by performing the same estimation procedure for each sample size n . For example, $\bar{X}_1 = X_1$, $\bar{X}_2 = (X_1 + X_2)/2$, $\bar{X}_3 = (X_1 + X_2 + X_3)/3$, etc. We can now define a consistent sequence.

Definition 10.1.1 A sequence of estimators $W_n = W_n(X_1, \dots, X_n)$ is a *consistent sequence of estimators* of the parameter θ if, for every $\epsilon > 0$ and every $\theta \in \Theta$,

$$\lim_{n \rightarrow \infty} P_\theta(|W_n - \theta| < \epsilon) = 1. \quad (1)$$

Informally, (1) says that as the sample size becomes infinite (and the sample information becomes better and better), the estimator will be arbitrarily close to the parameter with high probability, an eminently desirable property. Or, turning things around, we can say that the probability that a consistent sequence of estimators misses the true parameter is small. An equivalent statement to (1) is this: For every $\epsilon > 0$ and every $\theta \in \Theta$, a consistent sequence W_n will satisfy

$$\lim_{n \rightarrow \infty} P_\theta(|W_n - \theta| \geq \epsilon) = 0.$$

Definition 10.1.1 should be compared to Definition 5.5.1, the definition of convergence in probability. Definition 10.1.1 says that a consistent sequence of estimators converges in probability to the parameter θ it is estimating.

To summarize: a consistent estimator will produce an estimate extremely close to the true value(s) of the statistical object as the sample size become extremely large. "How close?" and "For what sample size?" are questions left to the discussion on variance below. An inconsistent estimator will converge to the wrong value or no finite value at all.

283 **6.5 Identifiability and Identification**

284 Identification is critical for statistical estimation. For a given statistical model, identifiability asks the question of
 285 whether it is even possible to create an estimator that could recover the underlying parameters of interest. In basic
 286 statistical procedures, identifiability often is not explicitly discussed since it is easy to prove, but it always exists
 287 in the background. For more advanced statistical procedures such as those that OSPEAD will use, identification
 288 must be discussed explicitly. [Allman et al., 2009] state,

289 General formulations of the identification problem were made by several authors, and pioneering works
 290 may be found in [27, 28]. The study of identifiability proceeds from a hypothetical exact knowledge of
 291 the distribution of observed variables and asks whether one may, in principle, recover the parameters.
 292 Thus identification problems are not problems of statistical inference in a strict sense. However, since
 293 nonidentifiable parameters cannot be consistently estimated, identifiability is a prerequisite of statistical
 294 parameter inference.

295 Here the authors are using “consistently estimated” in its technical sense defined above. In essence, statistical
 296 estimation is mostly hopeless without identifiability. They continue:

297 In the following, we are interested in models defined by a family $\mathcal{M}(\Theta) = \{\mathbb{P}_\theta, \theta \in \Theta\}$ of probability
 298 distributions on some space Ω , with parameter space Θ (not necessarily finite dimensional). The classical
 299 definition of identifiability, which we will refer to as *strict identifiability*, requires that for any two different
 300 values $\theta \neq \theta'$ in Θ , the corresponding probability distributions \mathbb{P}_θ and $\mathbb{P}_{\theta'}$ are different.

301 In other words, different parameter values *must* produce different data if a model is to be identifiable. [Casella & Berger, 2002]
 302 give another definition (pp. 523):

303 **Definition 11.2.2** A parameter θ for a family of distributions $\{f(x|\theta) : \theta \in \Theta\}$ is *identifiable* if distinct
 304 values of θ correspond to distinct pdfs or pmfs. That is, if $\theta \neq \theta'$, then $f(x|\theta)$ is not the same function
 305 of x as $f(x|\theta')$.

306 Identifiability is a property of the model, not of an estimator or estimation procedure. However, if a
 307 model is not identifiable, then there is difficulty in doing inference. For example, if $f(x|\theta) = f(x|\theta')$
 308 then observations from both distributions look exactly the same and we would have no way of knowing
 309 whether the true value of the parameter was θ or θ' . In particular, both θ and θ' would give the likelihood
 310 function the same value.

311 Realize that problems with identifiability can usually be solved by redefining a model.

312 “Redefining a model” can mean many things. There is “partial identification”, which involves narrowing down the
 313 possible values of θ to a set or a function of parameters. We will set aside partial identification. The redefinition of
 314 a model that is more relevant to our setting is to impose additional constraints, assumptions, or “structure” on the
 315 model, based on aspects of the data that are known with certainty or high probability. Fortunately, our problem
 316 — identification of Monero’s real spend age distribution — has enough structure for identification.

317 In another setting that may be familiar to the reader, ordinary least squares (OLS) regression is not identified
 318 when there is linear dependence between the columns of the “independent” variables, i.e. the design matrix \mathbf{X} . A
 319 typical cause of this would be perfect correlation of a pair of variables. In this case, the OLS estimator would give
 320 an infinite number of solutions to any particular parameter set $\boldsymbol{\theta}$ when minimizing the sum of squared errors. The
 321 \mathbf{X} matrix does not have full rank, so the inversion procedure $(\mathbf{X}'\mathbf{X})^{-1}$ in the $(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} = \boldsymbol{\theta}$ formula for OLS
 322 would be invalid. $\mathbf{X}'\mathbf{X}$ is not invertible in this case. Put differently, the sum of squared errors function in this case
 323 would achieve its minimum for an infinite number of values of $\boldsymbol{\theta}$, causing the function surface to have a “valley” of
 324 infinite length and fixed depth.

325 6.6 Variance of an Estimator

326 The expression for the variance of a random variable X is:

$$\text{Var}(X) = E \left[(X - E[X])^2 \right]$$

327 where $E[X]$ is the expectation, i.e. theoretical mean, of X .

328 An estimator itself has a variance since it is a function of random data from a sample of the population. In a
 329 frequentist paradigm — rather than a Bayesian paradigm — which OSPEAD orients itself within, an estimator's
 330 *estimand*, i.e. the parameter(s) to be estimated, are considered to be fixed in Nature with no variability in itself.
 331 Generally in practical applications an estimator's standard error, i.e. the square root of the variance of the estimator,
 332 is used since confidence intervals can be directly constructed from a standard error when an estimator is normally
 333 distributed.

334 For most estimators, only the asymptotic variance is known. [Casella & Berger, 2002] give a definition:

335 In calculating an asymptotic variance, we are, perhaps, tempted to proceed as follows. Given an esti-
 336 mator T_n based on a sample of size n , we calculate the finite-sample variance $\text{Var } T_n$, and then evaluate
 337 $\lim_{n \rightarrow \infty} k_n \text{Var } T_n$, where k_n is some normalizing constant. (Note that in many cases $\text{Var } T_n \rightarrow 0$ as
 338 $n \rightarrow \infty$, so we need a factor k_n to force it to a limit.)

339 **Definition 10.1.7** For an estimator T_n , if $\lim_{n \rightarrow \infty} k_n \text{Var } T_n = \tau^2 < \infty$, where $\{k_n\}$ is a sequence of
 340 constants, then τ^2 is called the *limiting variance* or *limit of the variances*.

341 [Some discussion of the drawbacks of using limiting variance.]

342 **Definition 10.1.9** For an estimator T_n , suppose that $k_n(T_n - \tau(\theta)) \rightarrow n(0, \sigma^2)$ in distribution [i.e. when
 343 scaled by k_n , the difference between the estimator and the estimand converges to a normal distribution
 344 with mean 0 and variance σ^2]. The parameter σ^2 is called the *asymptotic variance* or *variance of the
 345 limit distribution* of T_n .

346 A very large number of estimators are asymptotically normally distributed. Many estimators are essentially a type
 347 of sum. The sum of an infinite number of independent random variables converges to a normal distribution due
 348 to the Central Limit Theorem. However, it cannot be simply assumed that an estimator will be asymptotically
 349 normal especially when it is complex or unusual. In the vast majority of parametric settings, the k_n normalizing
 350 constant series is \sqrt{n} . In this case, such estimators are said to be " \sqrt{n} -consistent". In essence, the "speed" or "rate"
 351 of convergence for these estimators is \sqrt{n} . When attempting to nonparametrically estimate an entire distribution,
 352 the rate of convergence is typically slower than \sqrt{n} . This fact guides later recommendations about appropriate
 353 sample size.

354 6.7 Bias of an Estimator

355 If variance is a measurement of the precision of an estimator, then bias indicates its accuracy. [Casella & Berger, 2002]
 356 explain:

357 **Definition 7.3.2** The *bias* of a point estimator W of a parameter θ is the difference between the
 358 expected value W and θ ; that is, $\text{Bias}_\theta W = E_\theta W - \theta$. An estimator whose bias is identically (in θ) equal
 359 to 0 is called *unbiased* and satisfies $E_\theta W = \theta$ for all θ .

360 To be clear, many, many consistent estimators have some bias. With a large sample, some bias may not be
 361 concerning since bias generally shrinks as sample size increases. In informal discussions, "bias" is sometimes used
 362 when "inconsistency" is meant. For example, "omitted variable bias" actually refers to inconsistency; with an infinite
 363 sample size, omitted variable bias will cause an estimator to mis-estimate the true parameter.

With nonparametric density estimation such as kernel density estimation, some bias is generally unavoidable unless variance is allowed to balloon such as with an “undersmoothing” theoretical condition. This is just a particular instance of the classic variance-bias tradeoff in statistics.

7 Finite Mixture Models

Finite mixture models are the centerpiece of my estimation strategy. [McLachlan et al., 2019] describe finite mixtures:

The probability density function, or probability mass function in the discrete case, of a finite mixture distribution of a p -dimensional random vector \mathbf{Y} , takes the form

$$f(\mathbf{y}) = \sum_{i=1}^g \pi_i f_i(\mathbf{y}), \quad (2)$$

where the mixing proportions π_i are nonnegative and sum to one and where the $f_i(\mathbf{y})$ are the component densities. We refer to the $f_i(\mathbf{y}; \boldsymbol{\theta}_i)$ as densities, since even if the vector \mathbf{Y} is discrete, we can still view the $f_i(\mathbf{y})$ as densities by the adoption of counting measure. Typically, the component densities are taken to be known up to a vector $\boldsymbol{\theta}_i$ of parameters. In this case, we can write the mixture density as

$$f(\mathbf{y}; \boldsymbol{\Psi}) = \sum_{i=1}^g \pi_i f_i(\mathbf{y}; \boldsymbol{\theta}_i), \quad (3)$$

where $\boldsymbol{\Psi} = (\boldsymbol{\xi}^\top, \pi_1, \dots, \pi_{g-1})^\top$ denotes the vector of unknown parameters and where $\boldsymbol{\xi}_i$ consists of the elements of the $\boldsymbol{\theta}_i$ known a priori to be distinct. Here the superscript \top denotes transposition. In many applications, the component densities $f_i(\mathbf{y}; \boldsymbol{\theta}_i)$ are taken to belong to the same parametric family, for example, the multivariate normal.

Our setting is univariate since the focus is solely on age as a variable. Therefore, “ p -dimensional random vector \mathbf{Y} ” would be a one-dimensional random variable Y .

In general, identification of the mixing proportions π_i and component densities $f_i(y)$ is quite difficult without more structure imposed on the model. This should not be surprising. If the $f_i(y)$ components are fully nonparametric and their respective supports overlap fully, there is no clear way to use a set of single empirical draws from $f(y)$ to recover π_i and $f_i(y)$. In existing work, researchers have suggested that estimating the real spend age distribution of Monero may be impossible. [Ronge et al., 2021] remarks, “Due to the anonymous nature of anonymous cryptocurrencies, it is (supposedly) infeasible to learn the real signer distributions.” [Ni et al., 2021] states, “adversaries cannot guess the probability distribution of consumed tokens.” The conclusion of [Aeeneh et al., 2021] claims,

We showed that if the distribution used for selecting decoy TXOs are known, then a malicious party can significantly weaken the privacy features of CryptoNote by tracing back the transactions. However, our studies on Monero, a CryptoNote-based blockchain, show that users are using different source codes each having a different algorithms for mixin selection. Fortunately, since there is no distribution that a strong majority of users are using for mixin selection, an attacker cannot derive the distribution of truly spent transactions and arrange the second attack on Monero.

The general direction of the intuition of existing Monero research is correct: statistical identification of nonparametric mixture models can be exceedingly difficult. But ultimately their conclusion is incorrect. Identification is possible by leveraging the structure of Monero’s ring signatures.

Each given Monero ring is constructed from a set of decoys independently selected from a particular distribution by the user’s wallet software. In the statistical literature, such a structure is known as “repeated

measure(ment)s” [Crowder & Hand, 1990]. Sometimes the concept appears under the labels “exchangeable sequences”, “grouped observations”, “clustering with instance-level (must-link) constraints”, “mixture of product distributions”, “multivariate mixtures under the conditionally i.i.d. assumption” “latent-class models” and “latent structure” ([Allman et al., 2009], [Benaglia et al., 2009b], [Ritchie et al., 2020], [Wei & Nguyen, 2022]). A ring’s repeated measurements structure provides the critical ingredient for identification and then consistent estimation of the real spend age distribution.[Bonhomme et al., 2016] explains:

Mixture models are most often parametrically specified, and estimation may be done by maximum likelihood or indirect inference in a frequentist setting, or via Markov chain Monte Carlo techniques in a Bayesian approach.

There is a growing literature on non-parametric identification of finite mixtures. Univariate mixtures are generally not identified non-parametrically. (Additional restrictions may lead to identification. Bordes et al. (2006) and Hunter et al. (2007) studied two- and three-component mixtures of symmetric location families. Kitamura (2004) and Henry et al. (2013) studied models with conditioning variables.) In contrast, data on multiple measurements can represent a powerful source of identification. Hettmansperger and Thomas (2000), Hall and Zhou (2003) and Allman et al. (2009) provided identification results on component distributions and mixing proportions under the assumption that measurements are conditionally independent and the number of components is known. Hu (2008) and Kasahara and Shimotsu (2009) established related results in econometrics.

8 Modeling the Real Spend and Decoy Distributions as a Mixture

I will proceed with the notation of [Bonhomme et al., 2016] and then modify their basic model. Let x be a random variable with support $\mathcal{X} \subseteq \mathbb{R}$. For Monero, the support is block age going back to the first RingCT outputs. Let repeated independent and identically distributed measurements x_1, x_2, \dots, x_M of x have joint cumulative distribution function (CDF)

$$F(x_1, x_2, \dots, x_M) = \sum_{k=1}^K \omega_k \left(\prod_{m=1}^M F_k(x_m) \right) \quad (4)$$

where K is the number of decoy selection algorithms (DSA), ω_k is the proportion of Monero transactions using the k th DSA, M is the number of decoys, and F_k is the CDF of the k th DSA. Of course, each ring also contains the real spend. Let x_S be the random variable representing the age of the real spend. Redefine F_k as $F_{k,D}$ as the CDF of the k th DSA. Then the age of ring members have joint CDF

$$F(x_1, x_2, \dots, x_M, x_S) = \sum_{k=1}^K \omega_k \left(F_{k,S}(x_S) \times \prod_{m=1}^M F_{k,D}(x_m) \right) \quad (5)$$

where $F_{k,S}$ is the real spend age distribution within rings that are constructed by the k th DSA. The $F_{k,S}(x_S)$ can be multiplied by $\prod_{m=1}^M F_{k,D}(x_m)$ since the decoys are selected independently from the real spend. Here I allow for the possibility — indeed, likelihood — that $F_{k,S}$ is distinct for each DSA. This would occur if certain types of users were more likely to use wallet software with certain DSAs. Spending habits of light wallet users are likely different from users using the Monero GUI (i.e. standard) DSA. If services or centralized exchanges use custom DSAs, their true spend age distribution would also likely be distinct.

In taking the next step there is a small theoretical problem. I aim to transform formula (5) back into the form of (4). This requires that $\prod_{m=1}^{M+1} F_k(x_m) \stackrel{?}{=} F_{k,S}(x_S) \times \prod_{m=1}^M F_{k,D}(x_m)$. If $F_k(x_m) = \frac{1}{M+1} F_{k,S}(x_S) + \frac{M}{M+1} F_{k,D}(x_m)$ then a draw from F_k can be modeled as first drawing a Bernoulli random variable with probability of success $p = \frac{1}{M+1}$ ([McLachlan et al., 2019]). If the Bernoulli draw results in a success then draw from $F_{k,S}$. Otherwise, draw from $F_{k,D}$.

The $\prod_{m=1}^{M+1} F_k(x_m) \stackrel{?}{=} F_{k,S}(x_S) \times \prod_{m=1}^M F_{k,D}(x_m)$ equality supposes that for every $M + 1$ draws from F_k there will be exactly one draw from $F_{k,S}$ and M draws from $F_{k,D}$. In the $\prod_{m=1}^{M+1} F_k(x_m)$ expression there could also be zero draws of $F_{k,S}$, two draws of $F_{k,S}$, three, and so forth. Therefore, $F_{k,S}(x_S) \times \prod_{m=1}^M F_{k,D}(x_m)$ is more “exact”, “strict”, or “restricted” than $\prod_{m=1}^{M+1} F_k(x_m)$. The number of draws from $F_{k,S}$ can be modeled as a binomial distribution with number of draws $n = M + 1$ and success probability $p = \frac{1}{M+1}$ ([McLachlan et al., 2019]). When $M = 15$, i.e. rings of size 16, there will be exactly a single draw from $F_{k,S}$ with 38 percent probability. The mean, median, and mode of $\text{Binomial}(M + 1, \frac{1}{M+1})$ is 1 for all values of $M \geq 0$. I do not think that this theoretical problem causes a significant practical problem. Therefore, I will proceed as if $\prod_{m=1}^{M+1} F_k(x_m) = F_{k,S}(x_S) \times \prod_{m=1}^M F_{k,D}(x_m)$.

The model that OSPEAD seeks to estimate is:

$$F(x_1, x_2, \dots, x_M, x_S) = \sum_{k=1}^K \omega_k \left(\prod_{m \in \{1, 2, \dots, M, S\}} F_k(x_m) \right) \quad (6)$$

$$\text{s.t. } F_k(x_m) = \frac{1}{M+1} F_{k,S}(x_S) + \frac{M}{M+1} F_{k,D}(x_m)$$

The form of the model suggests a two-step procedure: first estimate $F(x_1, x_2, \dots, x_M, x_S)$ and second estimate the decomposed form of $F_k(x_m)$. To estimate $F(x_1, x_2, \dots, x_M, x_S)$, a good or practical guess of K , the number of DSAs being used to construct Monero transactions, is necessary. If K is under-estimated then likely some minority DSAs will be grouped together as the same k th DSA. Such a miscellaneous grouping is not necessary to identify the real spend age distribution of the bulk of Monero users, as I will argue in a moment. The first estimation step will produce good estimates of ω_k and $F_k(x_m)$ for each k .

For the second step in which $F_k(x_m)$ is estimated, the mixing proportions $\frac{1}{M+1}$ and $\frac{M}{M+1}$ are known since the number of decoys M is known. For identification and estimation of $F_{k,S}$ for a particular k it is necessary to have some knowledge of $F_{k,D}$ for the corresponding k . Fortunately, for several of the $F_{k,D}$ with large ω_k , i.e. the F_k that are most prevalent in the empirical on-chain data, it is possible to obtain an exact form of $F_{k,D}$ since they are described in open source code of Monero wallet software implementations. For the $F_{k,D}$ DSAs where there is not available open source code, my recommendation is to ignore them since identification of $F_{k,S}$ would be on very shaky ground.

The “standard” DSA in `wallet2` that the GUI wallet and many third-party wallets use may be the only relevant $F_{k,D}$. This is a judgment call that the panel should give input on. Since the purpose of OSPEAD is to replace the standard DSA, its protection can really only extend to users who are using the standard DSA. The objective is to set $F_{k,D}(x_m)$ as close to $F_{k,S}(x_S)$ as possible for any particular k .

An alternative choice is to consider the real spend age distribution in aggregate as much as possible for some subset of K DSAs that we have an open source reference for. As of now I am aware of only one other open source DSA that differs from the standard one: MyMonero.³ If the real spend age distribution of users who use `wallet2` and MyMonero is used to create the OSPEAD DSA, then the new OSPEAD DSA would be:

$$\tilde{F}_{OSPEAD,D}(x_m) = \frac{\omega_{wallet2} F_{wallet2,S}(x_m) + \omega_{MyMonero} F_{MyMonero,S}(x_m)}{\omega_{wallet2} + \omega_{MyMonero}} \quad (7)$$

Note that when $F_{wallet2,S}(x_m) \neq F_{MyMonero,S}(x_m)$ this version (7) of the OSPEAD DSA $\tilde{F}_{OSPEAD,D}$ would not offer best protection for `wallet2` users since the best protection occurs when $F_{OSPEAD,D}(x_m) \approx F_{wallet2,S}(x_m)$.

³See “Catalogue of Monero decoy selection algorithms” <https://github.com/monero-project/research-lab/issues/99>

470 9 Two-Step Estimation

471 My proposed estimator of the real spend age distribution is a type of two-step estimator. These estimators are used
 472 widely in econometrics.[Murphy & Topel, 1985] explain the mechanics of these models and their challenges:

473 The estimation of models that contain unobservable, though estimable, variables is now common in
 474 several areas of applied econometrics...In these and other econometric studies, unobserved variables
 475 are either replaced by their predicted values from an auxiliary statistical model when estimating the
 476 relationship of interest, or they are estimated jointly with that model. This article provides simple but
 477 theoretically correct procedures for statistical inference in this class of econometric models when joint
 478 procedures are infeasible or inappropriate.

479 There are two main approaches to estimating this type of model. The first, which we will refer to as the
 480 two-step (T-S) estimator, will be our main interest in this article. The T-S procedure simply replaces the
 481 unobserved components with their estimated or predicted values from the auxiliary model. Critically,
 482 in most applications these values are then treated as if they are known for purposes of estimation and
 483 inference in the second-step model, which is usually the model of interest. It is well known that T-S
 484 procedures yield consistent estimates of second-stage parameters under fairly general conditions. It is
 485 also well known that the second-step estimated standard errors and related test statistics based on these
 486 procedures are incorrect. In some cases this feature of the T-S estimator is acknowledged (e.g., Topel
 487 1982. Blanchard 1983), but more commonly the problem is ignored, and no correction of second-step
 488 test statistics is attempted. (Some authors who use the T-S procedure write as if the problem were
 489 merely one of efficiency [i.e. lower variance]. This conclusion is false, as our subsequent analysis will
 490 show.)

491 Following the suggestion of Liederman (1979), the usual alternative to the T-S procedure is to estimate
 492 the first- and second-step models via some joint method, such as full information maximum likelihood
 493 (FIML). Under appropriate assumptions, FIML procedures yield efficient estimators and asymptoti-
 494 cally correct estimates of standard errors. Unfortunately, in many situations FIML estimation is both
 495 computationally complex and costly to implement, facts that account in part for its limited use.

496 To summarize: two-step estimators are typically consistent, but naive computation of the standard errors of the
 497 second-stage parameters would be incorrect since the effect of sampling error in the first stage would be ignored.
 498 Joint estimation of the two models will provide correct standard errors, although joint estimation comes with its
 499 own drawbacks, assuming it is even possible to develop such a joint estimator.

500 Correct standard errors (i.e. correct confidence intervals or measures of an estimator's precision generally) are
 501 vital for scientific hypothesis testing. For OSPEAD, the important thing is that the two-step estimator is consistent.
 502 Obtaining estimates of the precision of the estimator is secondary. As long as the variance of the estimator is not
 503 uselessly enormous, using the estimator is preferable to using the log-gamma estimate from the pre-RingCT era
 504 based on the analysis of [Möser et al., 2018]. Correct standard errors would be of greatest interest to an attacker
 505 aiming to de-anonymize Monero users by leveraging an estimated real spend age distribution as input to a statistical
 506 attack. The attacker may want to use hypothesis testing to determine the false positive rate of such an attack.
 507 [Liu et al., 2019] provides a framework for statistical privacy based on hypothesis testing. I will return to the
 508 question of estimating standard errors in **Section 12 Confidence Interval Estimation**.

509 10 First Step: Bonhomme-Jochmans-Robin Estimator

510 The choice of estimator for the first step, i.e. the estimator of ω_k and $F_k(x_m)$ of equation (6), is probably the
 511 single most important decision to be made. In every other step there is a way to judge the correctness of the
 512 estimates in one way or another. For the first step there is no "ground truth." If we had a reliable ground truth, as
 513 [Möser et al., 2018] did, there would hardly be a need to perform statistics in the first place.

514 There are several methods to consistently estimate (4). The asymptotic variance of most of the existing methods
 515 has not yet been determined. The consistency of several methods has not even been rigorously established. They

516 also tend to be very computationally expensive. One exception is the method proposed in [Bonhomme et al., 2016].
 517 They explain the state of the field:

518 To estimate multivariate mixtures non-parametrically, computational procedures akin to the [Expectation
 519 Maximization] algorithm (Dempster et al., 1977) have recently been introduced by Benaglia et al.
 520 (2009) and Levine et al. (2011). These approaches are applicable more generally than are the earlier
 521 proposals of Hettmansperger and Thomas (2000) and Elmore et al. (2004), and of Hall and Zhou (2003)
 522 and Hall et al. (2005). Although simulation evidence suggests that these estimators work well in finite
 523 samples, their statistical properties are currently unknown. Chauveau et al. (2014) have provided an
 524 account of these developments.

525 [Benaglia et al., 2009b] describe the earlier “cutpoint” methods:

526 Hettmansperger and Thomas (2000); Cruz-Medina, Hettmansperger, and Thomas (2004); and Elmore,
 527 Hettmansperger, and Thomas (2004) treat the case [that is equivalent to Monero’s setting]....

528 The authors named above have developed an estimation method for the conditionally i.i.d. model. This
 529 method, the *cutpoint approach*, discretizes the continuous measurements by replacing each r -dimensional
 530 observation, say $\mathbf{X}_i = (x_{i1}, \dots, x_{ir})$, by the p -dimensional multinomial vector (n_1, \dots, n_p) , where $p \geq 2$
 531 is chosen by the experimenter along with a set of cutpoints $-\infty = c_0 < c_1 < \dots < c_p = \infty$, so that for
 532 $a = 1, \dots, p$,

$$n_a = \sum_{k=1}^r I\{c_{a-1} < x_{ik} \leq c_a\}.$$

533 Note that the multinomial distribution is guaranteed by the conditional i.i.d. assumption, and the
 534 multinomial probability of the a th category is equal to $\theta_a \equiv P(c_{a-1} < X_{ik} \leq c_a)$.

535 The cutpoint approach is completely general in the sense that it can be applied to any number of
 536 components m and any number of repeated measures r , just as long as $r \geq 2m - 1$, a condition
 537 that guarantees identifiability (Elmore and Wang 2003). However, some information is lost in the
 538 discretization step, and for this reason it becomes difficult to obtain density estimates of the component
 539 densities.

540 The $r \geq 2m - 1$ identifiability requirement for cutpoint methods is probably not very restrictive for our problem.
 541 For rings of size 11, the maximum number K of allowable density components F_k to be estimated would be 6.
 542 With ring size 16, it would be 8. The main problem with cutpoint methods for us, as [Benaglia et al., 2009b] say,
 543 is that the method destroys some information that is necessary to precisely estimate the density components F_k .
 544 The primary objective in this first step is to estimate F_k so it can be passed to the second step. Therefore, cutpoint
 545 methods are probably a dead end for us.

546 After cutpoint methods, Expectation Maximization (EM) methods were developed. EM algorithms are often
 547 computationally expensive ([O’Hagan et al., 2012]). Some of the papers mentioned above attempt to support the
 548 conjecture that their method is consistent through Monte Carlo simulations. For example, [Benaglia et al., 2009a]
 549 say,

550 Finally, we include a simulation study whose purpose is to explore the possible rate of convergence for the
 551 algorithm, for fixed m [number of distribution components] and r [number of repeated measurements],
 552 as the sample size n [corresponds to number of rings in Monero] tends to infinity. Note that the plots
 553 here do not constitute a proof of the asymptotic rate of convergence, nor even of consistency, yet they
 554 are interesting nonetheless because they suggest that such a theoretical result is possible.

555 Monte Carlo simulations can be helpful in ruling out obvious inconsistency of an estimator, but at best they occupy
 556 a distant second place compared to rigorous mathematical proofs. The main issue is that Monte Carlo simulation
 557 results are depended upon the particular distributions and parameters chosen for the simulations. Other data

samples may have different properties, resulting in possible inconsistency. Mathematical proofs generally cover all possible cases. It is much safer to go with an estimator that has a rigorous mathematical proof of consistency.

The Bonhomme-Jochmans-Robin estimator described in [Bonhomme et al., 2016] is the complete package: An estimator for an identifiable model that is mathematical proven to be consistent and asymptotically normal, with closed-form expressions for its asymptotic variance and bias. The estimation algorithm is reasonably computationally efficient and the paper authors released a MATLAB implementation. For space reasons I will not fully describe the estimator; the reader is referred to [Bonhomme et al., 2016] itself for full details. I will discuss a few of its features.

10.1 Linear Independence Assumption

Like any somewhat complicated estimator, the Bonhomme-Jochmans-Robin estimator requires that certain assumptions be met. Most of the required assumption are regularity conditions that are easily satisfied in practical applications. The single assumption that is more strict than a simple regularity condition is their Assumption 1:

Assumption 1 (rank). The matrix B has maximal column rank.

Assumption 1 is similar to the identification condition of [Allman et al., 2009], theorem 8, which require component distributions to be linearly independent. Indeed, with $\chi_i(x_m) = \mathbf{1}\{x_m \leq v_i\}$ for a set of chosen values v_1, \dots, v_I we have $\mathbb{E}[\chi_i(x_m)] = F_k(v_i)$ and assumption 1 demands linear independence of the component distributions on the grid v_1, \dots, v_I . Note that assumption 1 is testable. Indeed,

$$A \equiv \mathbb{E}[A(x_m)] = B\Omega B'.$$

Hence, for given I , assumption 1 is equivalent to A having rank K , which is a testable restriction. See the on-line supplementary material for details.

Theorem 8 of [Allman et al., 2009] requires that no component distribution can be completely reconstructed by a linear combination of the other component distributions. This would require reconstruction at every point x in the support of such a component. Note that this meaning of linear independence is different from the meaning discussed in **Section 6.5 Identifiability and Identification** with OLS estimation. As argued in [Mbakop, 2017] and [Kwon & Mbakop, 2021], this linear independence assumption of Theorem 8 of [Allman et al., 2009] is generally easily satisfied when the component distributions are continuous.

The story is a bit different with Assumption 1 of the Bonhomme-Jochmans-Robin estimator because the component distributions are evaluated at a discrete number of points even if the distributions are continuous, which could make the assumption harder to satisfy. Using a larger number of points, i.e. setting I to be larger, would probably increase the likelihood that this assumption is satisfied. On the bright side, this Assumption 1 is testable. To give context, it is always very good when a statistical assumption is testable with the sample. When an important assumption is not testable by data, we are left wondering if our estimate is valid without any way to directly test it. OSPEAD will include a statistical test of this Assumption 1 condition.

10.2 Permutation of Triples

The core of the Bonhomme-Jochmans-Robin Estimator involves using every permutation (*not* combination) of ordered triples (i.e. 3-tuples) from each set of repeated measures. Certain characteristics of the estimator flow from the triple permutations. First, the minimum number of repeated measurements necessary for identification and estimation is 3. That presents no problems for us since ring size is now 16, up from the pre-hardfork 11. As will be discussed in detail in **Section 12 Confidence Interval Estimation**, the variance of the estimator shrinks in rough relation to n choose k , i.e. the binomial coefficient $\binom{n}{k} = \frac{n!}{k!(n-k)!}$, where n is the ring size and $k = 3$.

The triple permutations also influence the computational expense of the estimator. By taking triples rather than dealing with the entire set of repeated measures, the Bonhomme-Jochmans-Robin estimator skirts the curse of dimensionality that afflicts other estimators of its class. The main computational bottleneck requires computation and summation of $N \times I^3 \times \frac{(M+1)!}{(M+1-3)!}$ terms, where N is the number of “rings”, M is the number of decoys to maintain the same notation as above (And thus $M + 1$ is ring size), and I is the number of “evaluation points” of the distributions. I have verified this quantity through code profiling of the paper’s MATLAB code. The choice of I is important for precision of the estimator. Larger I increases precision. Table 1 of [Bonhomme et al., 2016] gives the Monte Carlo bias and standard errors of the estimates of the means of the component distributions with $I = 5$ and $I = 10$. With greater I there is a moderate improvement in precision. Computational expense is the limitation on I because expense increases with the cube of I . Therefore, it would be best to convert the MATLAB implementation into a C++ implementation, which I discuss in **Section 18 OSPEAD Requests to Monero C++ Developer(s)**, so that I can be made as large as possible.

10.3 Determining the Number of Distribution Components K

Historically, determining the number of distribution components K for fitting finite mixture models has relied on heuristics with questionable foundations, especially in the nonparametric setting [McLachlan & Peel, 2000]. Usually, K must be specified at the very beginning of the estimation procedure and therefore OSPEAD would be incomplete without a specification of how to choose K . In our setting, the modeling choice or estimate of K is equivalent to the choice or estimate of the number of decoy selection algorithms being used in the wild on the blockchain. Most papers on estimating mixture models consider the choice of K to be outside of their scope, and [Bonhomme et al., 2016] is no exception. Fortunately, in the last decade a few papers have made progress.

[Kasahara & Shimotsu, 2014] analyze the case of multivariate mixtures where the multiple variables “are independently (but not necessarily identically) distributed within each component.” Our repeated measures is a special case of this multivariate setting: the variables are both independently and identically distributed within each component (or at least identically distributed in the sense of Equation (6)). Therefore, the [Kasahara & Shimotsu, 2014] method would also work in our stricter setting. [Kasahara & Shimotsu, 2014] give an identification result when there are at least two variables (in our setting, at least two ring members) as well as a procedure to consistently estimate a lower bound on the number of distribution components.

[Kwon & Mbakop, 2021] improves on [Kasahara & Shimotsu, 2014] by creating a consistent estimator of the exact number of components (not just a lower bound). From their Monte Carlo simulations it appears that their estimator is not very precise when K is high but N (number of rings) is low. In our case we have both K and N high, so their estimator should correctly estimate the number of components. Therefore, as of now I plan to use their procedure to estimate K despite the fact that there are two complications. The first is that [Kwon & Mbakop, 2021] do not make their code available. Either I could re-implement the method or we could ask them if they are willing to share the code. The second complication is that the method requires computing each element of two $N \times N$ matrices, solving for their matrix square root and then computing the singular value decomposition of their product (see Algorithm 1 of [Kwon & Mbakop, 2021]). Given that our target N is about 400,000, [Kwon & Mbakop, 2021] could get computationally expensive. Taking a random subsample may be necessary.

10.4 Comparison with Ritchie-Vandermeulen-Scott Estimator

I will now compare the Bonhomme-Jochmans-Robin estimator to an estimator described in [Ritchie et al., 2020], which is based on theoretical results from [Vandermeulen & Scott, 2019]. At one point I considered the [Ritchie et al., 2020] estimator to be the leading candidate for a potential estimator for the first step, but on deeper analysis it is clear to me that it is inferior to the Bonhomme-Jochmans-Robin estimator for our particular problem.

Overall, [Ritchie et al., 2020] approaches the problem from more of a machine learning angle. (Bonhomme, Jochmans, and Robin are econometricians.) They show identifiability and consistency of their estimator, but do not give an expression for its variance, unlike [Bonhomme et al., 2016]. They are particularly concerned with classification of observations into the component distributions after estimation and compare their estimator with existing clustering methods in the vein of machine learning.

[Ritchie et al., 2020] and [Vandermeulen & Scott, 2019] make some theoretical achievements such as showing that a mixture distribution with repeated measurements is identifiable even if the number of repeated measurements is only two as long as an additional “joint irreducibility” condition is met. [Bonhomme et al., 2016] requires at least three. This achievement is not relevant for our problem since our number of repeated measurements is 11 and 16. The method of [Ritchie et al., 2020] also can be used with multiple variables, unlike that of [Bonhomme et al., 2016]. This is not relevant to us since the age of spent transaction outputs is univariate.

As an aside, the [Ritchie et al., 2020] describe why conventional clustering methods cannot identify nonparametric mixture models with repeated measurements:

Mixture models are often utilized to solve the clustering problem. Parametric mixture models, such as GMMs, are able to capture overlapping clusters. Most clustering algorithms, however, such as k -means [11, 12], DBSCAN [13], and spectral clustering [14, 15], assume clusters are non-overlapping and hence fail when clusters overlap.

By “overlapping clusters” here they mean settings where the component distributions F_k have partial or completely overlapping support (such the setting of Monero’s ring signature age distribution). See [Vankadara et al., 2021] for more discussion of this point.

The computational expense of [Ritchie et al., 2020] basically torpedoes its use in our setting, unfortunately. The Python implementation code and the example data all used two repeated measures. It took me some time to realize why: the computational expense of their estimator explodes as the number of repeated measurements increases. According to my reading of the paper and the included Python code, the number of terms that need to be summed is approximately $N \times I^{M+1}$ where N is the number of “rings”, M is the number of decoys, and I is the number of “evaluation points” of the distributions as before. The meaning of evaluation points is different in [Ritchie et al., 2020], but it serves a similar roles as in [Bonhomme et al., 2016]. A minimum acceptable value of I is probably about 10. If $N = 100,000$, $I = 10$, and $M = 15$, then the number of terms to sum would be $100,000 \times 10^{16} = 10^{21}$, which is computationally infeasible. With Bonhomme-Jochmans-Robin estimator, the required number of terms to sum would be only $100,000 \times 10^3 \times \frac{16!}{13!} = 3.36 \times 10^{11}$. (With ring size 128, the estimator would require summation of $100,000 \times 10^3 \times \frac{128!}{125!} = 2.03 \times 10^{14}$ terms, which seems computationally feasible with an optimized implementation.)

One way to reduce the computational expense of the [Ritchie et al., 2020] estimator is to take a random subsample of each ring, e.g. only select 4 ring members from each ring. In theory there should be no theoretical problem with doing this. The estimator should still be consistent since the sample was selected in an independent way. Such an approach is probably inferior to using the Bonhomme-Jochmans-Robin estimator with the full sample. It would be hard to justify the subsampling given that we do not even know the variance of the [Ritchie et al., 2020] estimator. Why we should accept excluding data when the Bonhomme-Jochmans-Robin estimator uses all observations?

10.5 Options for the First Step Estimator

My recommendation is to use the Bonhomme-Jochmans-Robin Estimator for the first step. Possible alternatives are:

1. The Ritchie-Vandermeulen-Scott Estimator. It is likely inferior to the Bonhomme-Jochmans-Robin Estimator since its variance is unknown and its computational expense is large.

682 2. The estimator described in [Benaglia et al., 2009a]. The `npEM` function in the `mixtools` R package implements
 683 this estimator. It is likely inferior to the Bonhomme-Jochmans-Robin Estimator since its consistency has not
 684 yet been proven.

685 3. The estimator described in [Levine et al., 2011]. The `npMSL` function in the `mixtools` R package implements
 686 this estimator. It is likely inferior to the Bonhomme-Jochmans-Robin Estimator since its consistency has not
 687 yet been proven.

688 Another option is to use the [Benaglia et al., 2009a] and/or [Levine et al., 2011] EM estimators as “checks” on
 689 the Ritchie-Vandermeulen-Scott Estimator. I do not recommend this because it is a kind of tautology. If those
 690 EM estimators arrive at similar estimates as the Ritchie-Vandermeulen-Scott Estimator, then we conclude that
 691 the Ritchie-Vandermeulen-Scott Estimator gives a good estimate. If the EM estimators arrive at a dramatically
 692 different estimate (i.e. well outside of the confidence intervals — this can be made formal with a hypothesis test and
 693 joint bootstrapping), then what should we conclude? It seems obvious to trust the estimate of a consistent estimator
 694 rather than the estimate of estimators with unknown consistency properties. We use the Ritchie-Vandermeulen-
 695 Scott Estimator anyway.

696 The number of distribution components K must also be decided. I recommend the estimator described in
 697 [Kwon & Mbakop, 2021]. Possible alternatives are:

- 698 1. Use the [Kasahara & Shimotsu, 2014] estimator to estimate a lower bound on K . Then choose some K greater
 699 than or equal to that estimated value.
- 700 2. Use another method with poorer theoretical foundations such as sequential hypothesis testing (SHT), Akaike
 701 Information Criteria (AIC), Bayesian Information Criteria (BIC), and Hannan-Quinn Information Criteria
 702 (HQ), which both [Kasahara & Shimotsu, 2014] and [Kwon & Mbakop, 2021] compared their proposed esti-
 703 mators against.
- 704 3. Select K based on “expert knowledge” of the probable number of decoy selection algorithms being used in the
 705 wild.

706 11 Second Step: Patra-Sen Inversion Estimator

707 Once each $F_k(x_m)$ from model (6) is estimated, the real spend age distributions $F_{k,S}(x_S)$ of the each component
 708 of interest must be estimated. [Patra & Sen, 2016] describe a model and estimator that can be used to estimate
 709 $F_{k,S}(x_S)$. They begin their paper by setting up their main model of interest:

710 Consider a mixture model with two components, i.e.

$$F(x) = \alpha F_s(x) + (1 - \alpha) F_b(x) \quad (8)$$

711 where the cumulative distribution function (CDF) F_b is known, but the mixing proportion $\alpha \in [0, 1]$
 712 and the CDF F_s ($\neq F_b$) are unknown. Given a random sample from F , we wish to estimate (non-
 713 parametrically) F_s and the parameter α

714 In this paper we provide a methodology to estimate α and F_s (non-parametrically), without assuming
 715 any constraint on the form of F_s .

716 The paper briefly discusses the model that corresponds to our problem: when α — the proportion of ring members
 717 that are real spends — is known:

718 Suppose that we observe an independent and identically distributed sample X_1, X_2, \dots, X_n from F as
 719 in model (8). If $\alpha \in (0, 1]$ were known, a naive estimator of F_s would be

$$\hat{F}_{s,n}^\alpha = \frac{\mathbb{F}_n - (1 - \alpha)F_b}{\alpha} \quad (9)$$

720 where \mathbb{F}_n is the empirical CDF of the observed sample, i.e. $\mathbb{F}_n(x) = \sum_{i=1}^n \mathbf{1}\{X_i \leq x\}/n$. Although
 721 this estimator is consistent, it does not satisfy the basic requirement of a CDF: $\hat{F}_{s,n}^\alpha$ need not be non-
 722 decreasing or lie between 0 and 1. This naive estimator can be improved by imposing the known shape
 723 constraint of monotonicity. This can be accomplished by minimizing

$$\int \left\{ W(x) - \hat{F}_{s,n}^\alpha(x) \right\}^2 d\mathbb{F}_n(x) \equiv \frac{1}{n} \sum_{i=1}^n \left\{ W(X_i) - \hat{F}_{s,n}^\alpha(X_i) \right\}^2 \quad (10)$$

724 over all CDFs W . Let $\check{F}_{s,n}^\alpha(x)$ be a CDF that minimizes expression (10).

725 Note that $\mathbf{1}\{x\}$ is the indicator function.⁴ Patra and Sen then describe a numerical algorithm to minimize (10) with
 726 computational complexity $O(n)$. Equation (9) has been referred to as an “inversion estimator” [Milhaud et al., 2021].
 727 It is a type of plug-in estimator because estimated quantities (here, only \mathbb{F}_n needs to be estimated since α and F_b are
 728 known) are plugged into the equation to determine the estimate. Their equation (9) is very similar to equation (2)
 729 in my HackerOne submission. [Patra & Sen, 2016] used a CDF form. I used a PDF form. The definitions of α and
 730 $(1 - \alpha)$ are switched in the two papers. I also recognized the issue of constructing a valid probability distribution
 731 and chose a less sophisticated correction:

732 Note that due to the long thin tails of the distributions and the noisy nature of the empirical data
 733 and jberman’s simulated data, it can be the case that the calculated value of $f_S(x)$ [the PDF of the
 734 real spend age distribution] can be negative for some values of x , which violates an axiom of probability
 735 theory. In those cases, the negative value of $f_S(x)$ was replaced by $f_M(x)$ [the PDF of the decoy selection
 736 algorithm], since that would be the most conservative approach.

737 Patra and Sen’s (9) is also similar to equation (2) in [Aeeneh et al., 2021], which develops the MAP Decoder attack
 738 on Monero, although in that paper it is described as a probability relationship rather than as a plug-in estimator
 739 per se.

740 We are quite fortunate that the structure of ring signatures guarantees that α is known with certainty rather
 741 than being a parameter that must be estimated. The value of α is not just known in the population, i.e. *in*
 742 *expectation*, but is known to be exactly $\frac{1}{M+1}$ in the sample. [Patra & Sen, 2016] note, “When α is unknown, the
 743 problem is considerably more difficult; in fact, it is non-identifiable.” This mixture model is sometimes referred to
 744 as a contamination or admixture model [Milhaud et al., 2021].

745 [Patra & Sen, 2016] does not derive the variance of $\check{F}_{s,n}^\alpha$. There are a few papers with similar, more general
 746 models (such as treating α as unknown) that have some variance estimator. [Arias-Castro & Jiang, 2021] use a
 747 bootstrap procedure, for example.

748 11.1 Options for the Second Step Estimator

749 The equation (8) has four unknowns to start with: $F(x)$, α , $F_s(x)$, and $F_b(x)$. When two of them (α and $F_b(x)$)
 750 are known with certainty and one of them ($F(x)$) can be estimated, we have one linear equation in one unknown.
 751 Therefore, it seems that the Patra-Sen Inversion Estimator is very suitable. I am not aware of an alternative
 752 estimator. There is not an obvious way to improve on the estimator, except perhaps in the rectification step of
 753 equation (10). The rectification of $\hat{F}_{s,n}^\alpha$ into a valid CDF is not just a matter of theoretical tidiness. $\hat{F}_{s,n}^\alpha$ needs

⁴https://en.wikipedia.org/wiki/Indicator_function

754 to be passed to yet another step to fit a parametric distribution over an estimate of a valid PDF. Therefore, some
 755 rectified estimate such as $\check{F}_{s,n}^\alpha(x)$ is necessary. In saying that the rectification step could possibly be improved, I do
 756 not mean that my ad hoc rectification that I performed in my HackerOne submission is an improvement. Rather,
 757 it is possible that there exists a published rectification method that is superior to the Patra-Sen method for our
 758 particular setting. I am not aware of any. The options are thus:

- 759 1. Use the Patra-Sen Inversion Estimator, including the rectification in Equation (10).
 760 2. Use the Patra-Sen Inversion Estimator, but use the ad hoc rectification that I used in my HackerOne submis-
 761 sion.

762 I recommend option (1).

763 I will state here an alternative to the two-step estimation procedure I have recommended. We could just ignore
 764 the existence of other decoy selection algorithms and apply the Patra-Sen Inversion Estimator directly to the data
 765 without using the first step. This is roughly equivalent to what I did in my HackerOne submission. Obviously, I do
 766 not recommend this, but I wanted to raise it as an option. The performance of such an estimator would worsen with
 767 the ring size increase from 11 to 16 because there would be even more “noise” in the raw data from the nonstandard
 768 decoy selection algorithms compared to the signal of real spends.

769 12 Confidence Interval Estimation

770 As I stated before, for our current problem the estimation of confidence intervals is secondary to the main goal
 771 of consistently estimating the real spend age distribution. It is still useful to estimate confidence intervals since
 772 excessively large confidence intervals might suggest that sample size should increase — e.g. increase the time window
 773 from one week to one month.

774 Given time constraints and the technical difficulty of deriving a closed-form expression for the two-step estimator
 775 outlined above, bootstrapping is the obvious method for estimating confidence intervals.

776 An interesting object of investigation is how the width of confidence intervals respond to a change in ring size.
 777 ArticMine has taken a particular interest in this question. A larger ring size gives the Bonhomme-Jochmans-Robin
 778 estimator more information to estimate each component F_k in the first step. Therefore, raising ring size will shrink
 779 confidence intervals in the first step. In the second step, we might expect the confidence interval on the estimate $F_{k,s}$
 780 to expand as the signal-to-noise ratio falls when ring size increases. These two effects move in opposite directions,
 781 so the direction of their combined effect is a question of their relative magnitudes.

782 12.1 Variance of the First Step Estimator

783 [Bonhomme et al., 2016] give the following closed-form expression for the asymptotic distribution of their estimator
 784 of the component densities:

785 *Proposition 2.* Let assumptions 1–3 and 5 and 6 hold. Then $|\hat{f}_k(x) - f_k(x)| = o_P(1)$ and

$$\sqrt{(Nh)} \left\{ \hat{f}_k(x) - f_k(x) \right\} \xrightarrow{L} \mathcal{N}(\sqrt{c}\mu_f, \mathcal{V}_f),$$

786 as $N \rightarrow \infty$ where $c \geq 0$ is a finite constant,

$$\mu_f \equiv \frac{1}{2} f_k''(x) \int_{-\infty}^{\infty} u^2 \kappa(u) du,$$

$$\mathcal{V}_f \equiv M \left\{ \frac{(M-3)!}{M!} \right\}^2 q_k(x) f(x) \int_{-\infty}^{\infty} \kappa(u)^2 du,$$

787 provided that $Nh \rightarrow \infty$ and $Nh^5 \rightarrow c$.

788 In proposition 2 we allow the bandwidth to vanish at the optimal rate of $N^{-1/5}$.

789 Proposition 2 says that the difference between the true density and the estimator converges to a normal distribution
 790 with mean (i.e. bias) $\sqrt{c}\mu_f$ and variance \mathcal{V}_f when the kernel bandwidth is h . As I said in **Section 6.6 Variance**
 791 **of an Estimator**, nonparametric density estimators have a slower convergence rate than parametric estimators.
 792 And some bias is generally unavoidable. $f_k''(x)$ in μ_f is the second derivative of the k th component's PDF at point
 793 x . The $\kappa(u)$ is the chosen kernel function, which is usually the Normal kernel.

794 The meaning of $q_k(x)$ in the definition of \mathcal{V}_f is available in [Bonhomme et al., 2016] and not too important for
 795 the following discussion. The relevant fact is that $q_k(x)$ does not depend on M . The point I want to make is that
 796 for a fixed value of $q_k(x)f(x) \int_{-\infty}^{\infty} \kappa(u)^2 du$, the estimator's asymptotic variance \mathcal{V}_f is proportional to $M \left\{ \frac{(M-3)!}{M!} \right\}^2$.
 797 Since the estimator's asymptotic distribution is Normal, it is easy to compute the effect of an increased ring size on
 798 the confidence interval when our sample size is very large.

799 A 90% confidence interval (a reasonable interval for our current problem) would be $[\hat{f}_k(x) - 1.64\sqrt{\mathcal{V}_f}, \hat{f}_k(x) + 1.64\sqrt{\mathcal{V}_f}]$.

800 In other words, the width of the confidence interval is proportional to the square root of the variance of the estima-
 801 tor, i.e. its standard error. Therefore, to evaluate the effect of ring size on the confidence interval of estimator $\hat{f}_k(x)$,
 802 the ratio of the two $\sqrt{\mathcal{V}_f}$ can be calculated. In other words, $\sqrt{M \left\{ \frac{(M-3)!}{M!} \right\}^2} = \sqrt{M} \frac{(M-3)!}{M!}$ should be calculated
 803 for two different M values. Then the expected ratio of the width of the confidence intervals when ring size is 16
 804 compared to ring size 11 is

$$\left(\sqrt{16} \frac{(16-3)!}{16!} \right) / \left(\sqrt{11} \frac{(11-3)!}{11!} \right) \approx 0.36$$

805 Therefore, we may expect that the confidence intervals of $\hat{f}_k(x)$ when ring size is 16 would only be 36% the
 806 width of confidence intervals when ring size is 11, i.e. confidence interval width would fall by 64 percent. In the
 807 broader context, the Flourine Fermi hard fork thus contributed to greatly improved precision *for the first step of*
 808 *the estimator*. If Monero's ring size increases to 128 as is probable with the Seraphis upgrade, then the ratio of the
 809 confidence intervals for ring size 16 vs 128 would be

$$\left(\sqrt{128} \frac{(128-3)!}{128!} \right) / \left(\sqrt{16} \frac{(16-3)!}{16!} \right) \approx 0.005$$

810 So the confidence interval would be less than half a percent in width with ring size 128 compared to 16. In other
 811 words, ring size of 128 (without binning) would produce an extremely precise estimate of $f_k(x)$. The computational
 812 expense of the estimator would rise, however.

813 12.2 Variance of the Second Step Estimator

814 This is only half of the story, of course. The variance of the Patra-Sen Inversion Estimator must be considered as
 815 well. As I stated in **Section 11 Second Step: Patra-Sen Inversion Estimator**, I am unaware of any closed-
 816 form expression for the asymptotic variance of the Patra-Sen Inversion Estimator. In cases such as these, Monte
 817 Carlo simulation usually gives a good approximation of the variance of the estimator, at least for the parameters
 818 chosen for the simulation. The `patra-sen-monte-carlo.R` code for a simulation is included in this zip file.

819 For this simulation I set the real spend age distribution to be negative binomial with mean parameter 100
 820 and dispersion parameter 3. Then the decoy distribution was negative binomial with mean parameter 150 and
 821 dispersion parameter 2. 100,000 rings were constructed for each of 1,000 Monte Carlo simulations, with ring size
 822 11, 16, and 128. I applied the Patra-Sen Inversion Estimator with rectification $\check{F}_{s,n}^{\alpha}(x)$ in each simulation (plus a
 823 kernel smoothing step) and obtained a 90 percent confidence interval at each point in the support x by calculating
 824 the 5th and 95th percentile of the simulated estimates. I chose confidence intervals instead of standard errors as

the metric of comparison since it is unknown at this point in time whether $\check{F}_{s,n}^\alpha(x)$ is asymptotically normal — and therefore it is unknown whether the width of the confidence intervals are strictly proportional to the standard errors. I plotted these confidence intervals and split the support to show detail in Figure 1.

The first thing to notice in these plots is that the confidence intervals when ring size is 16 are only slightly larger than when ring size is 11. I calculated the median of the width of all confidence intervals when $x \in \{20, 21, \dots, 200\}$ for ring size 11, 16, and 128 to generate a summary statistic for comparison. According to this median, the confidence interval for ring size 16 was only 14 percent larger than for ring size 11. However, precision suffered quite a bit when ring size rose to 128: the confidence interval with ring size 128 was 212 percent wider than with ring size 16.

Another thing to notice is the estimator's apparent bias in the support range 1 to 20. As I have mentioned before, some bias is generally unavoidable with nonparametric kernel density estimators. Looking carefully at the bias here, it is notable that the estimator is *downward* biased when the decoy distribution is greater than the real spend age distribution. My intuition tells me that the bias probably is not being caused by the density of decoy distribution; if it were, I would expect an upward — not downward — bias in this range of the support. The bias may be due to the typical bias with kernel smoothing or some effect of the imposing the monotonicity constraint on the CDF as in equation (10). In general, bias correction procedures are possible, but in practice they are often avoided since they usually increase the variance of estimators too much.

Despite the analysis above, we still cannot say with certainty whether a rise in ring size from 11 to 16 increases or decreases the precision of my proposed two-step estimator. This is for two reasons. First, the confidence interval of the first step's estimator may shrink by 64 percent and the confidence interval of the second step's estimator may expand by only 14 percent, but we cannot know which effect dominates without knowing the initial confidence interval widths for both estimators at ring size 11. If the confidence interval of the first step estimator is very small to begin with and the confidence interval of the second step estimator is very large, then 14 percent of a large number could be much bigger than 64 percent of a small one.

The second reason is the general complication of variance estimation in two-step models discussed in [Murphy & Topel, 1985] and related work. Unfortunately, the answer to this overall question will have to wait for the full bootstrap or Monte Carlo estimations. To be clear, bootstrapping and Monte Carlo estimates for two-step estimators must do the (re)sampling at the first step and then carry that same sample to the second step estimate, then resample again at the first step for the next iteration.

Bootstrap-based confidence intervals involve a special challenge with finite mixture models. [Bonhomme et al., 2016] explains:

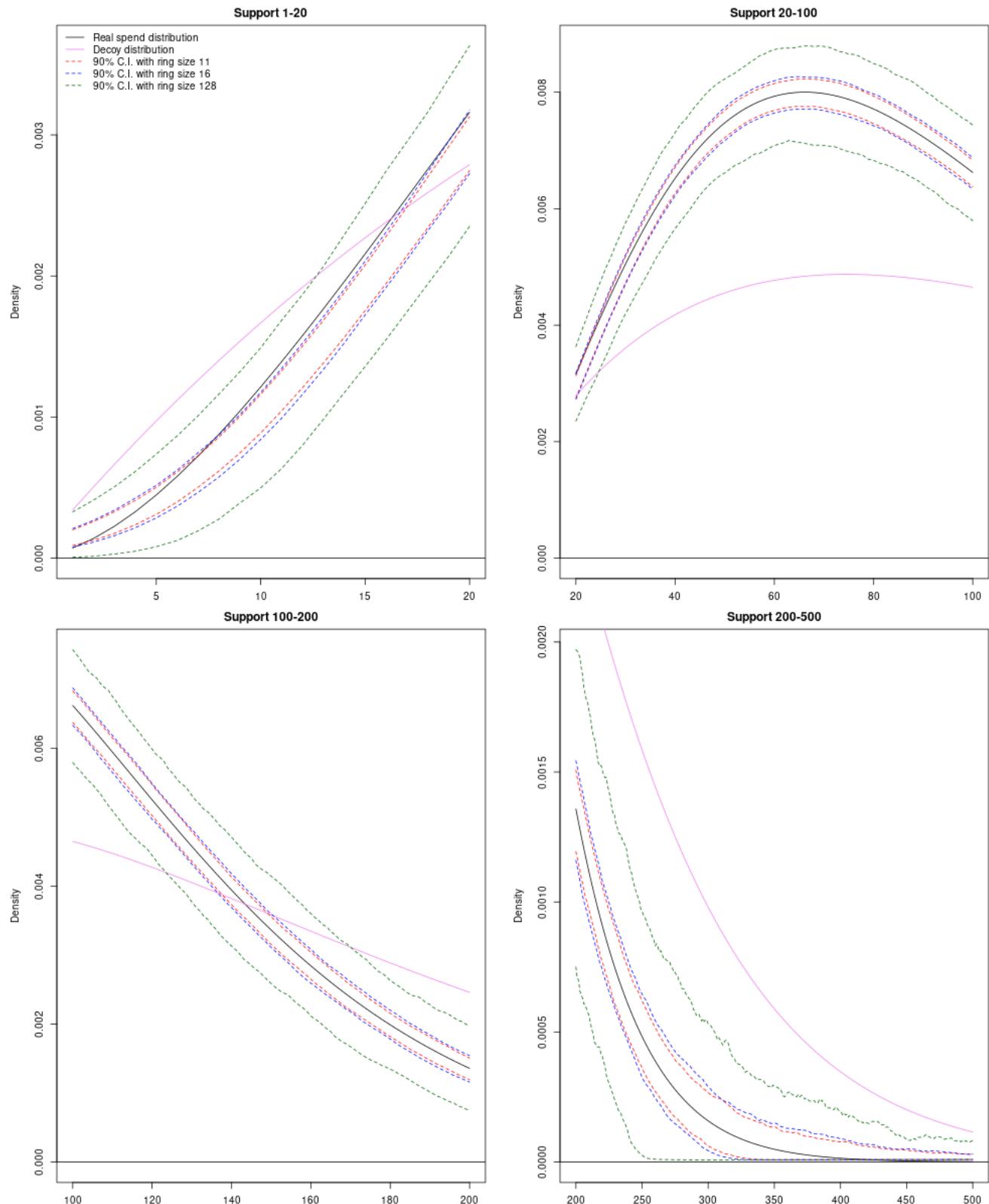
We note that, although it must be taken care of in our simulations, label swapping does not cause any complications for estimation and inference based on our approach. Label swapping does present an important challenge for inference methods based on resampling algorithms such as the bootstrap or jackknife, and on Markov chain Monte Carlo procedures (see, for example, Stephens (2000)).

We need to be able to say that the k th F_k estimated distribution in the first step is associated with the k th decoy selection algorithm for estimation in the second step. It is perhaps not too difficult to do this when a human can examine a single estimate. Bootstrapping involves many such estimates; I would aim for 1,000 for OSPEAD. Therefore, some automated way of labeling is recommended. [Bonhomme et al., 2016] describe their own method to avoid label-swapping in their Monte Carlo simulations:

To deal with label swapping in our simulation study we proceed as follows. In each replication we first estimate the mean of each component. We then label the estimated component with the smallest estimated mean as the first mixture component. We carefully checked for label swapping and found no mix-up.

The mean of distributions is not very meaningful nor specific in our setting. Here I tentatively develop two methods that may work alone or combined to perform automated labeling. The first is to simply assume that the mixing

Figure 1: Patra-Sen Monte Carlo Confidence Intervals by Ring Size



proportions ω_k would remain stable in each bootstrap iteration. My hypothesis is that most transactions are being constructed with the standard `wallet2` DSA, followed by the MyMonero DSA as the second most prevalent. We will have to see what the data tells us to confirm or invalidate that hypothesis.

The second possible approach to automated labeling involves enforcement of one of the probability axioms: that probability cannot be negative. The idea is similar to the imposition of monotonicity of the CDF in equation (10) from the Patra-Sen Inversion Estimator. Once a specified DSA distribution is subtracted from the mixture distribution F_k , the remaining portion is the estimate of the real spend age distribution $\hat{F}_{k,S}$ (which is the label we care about). Of course, due to sampling error some portion of the corresponding estimated real spend age distribution $\hat{F}_{k,S}$ may appear as negative; this suggest a statistical test should be developed, or at least a heuristic. As ring size increases, this type of automated labeling should get more accurate since there is less flexibility in the thickness of the real spend age distribution $F_{k,S}$ as a proportion of the component distribution F_k .

12.3 Variance of the Estimator of the Number of Distribution Components K

The estimator for the number of distribution components K is also subject to sampling error. [Kwon & Mbakop, 2021] do not give a way to compute their estimator's variance or confidence intervals. Bootstrapping could be a reasonable strategy because it is usually valid even when the parameter space is discrete [Newton, 1996]. The parameter space of K is the natural numbers. However, my recommendation would be to ignore the variance of \hat{K} and to use the (weekly) estimate of K as its true value because:

1. K is discrete and probably less than 10. Therefore, the practical variance of \hat{K} is probably negligible given our sample sizes.
2. A proper bootstrap procedure for \hat{K} creates special problems for the downstream estimators. Say a bootstrap estimate of K produces varied estimates of K , which anyway would be the whole point of computing a bootstrap iteration: to estimate the variability of \hat{K} . Then some of the downstream bootstrap estimates of the Bonhomme-Jochmans-Robin Estimator and the Patra-Sen Inversion Estimator would have a varied number of estimated components. Then combining those estimates for bootstrap confidence intervals would result in some “missing” and “extra” components, greatly complicating valid statistical inference.
3. The estimator described in [Kwon & Mbakop, 2021] is computationally expensive. If bootstrapped, the procedure would have to be done about 1,000 times for each week of data rather than a single time for each week.

12.4 Options for Confidence Interval Estimation

My recommendation is to perform about 1,000 nonparametric bootstrap iterations, starting each iteration from the first step and then passing the results of the first step to the second step in each iteration. A nonparametric bootstrap, also known as a case resampling bootstrap, requires the fewest assumptions about our model and estimator compared to other bootstrap subtypes. Setting the number of iterations to 1,000 should provide enough information to calculate confidence intervals with good coverage.⁵ If the computation is very expensive, then we could settle for 100 iterations. Computation of a closed-form asymptotic confidence interval is not really an option due both to the two-step nature of the estimator (as [Murphy & Topel, 1985] warn) and because there is not yet a closed-form expression for the variance of the Patra-Sen Inversion Estimator.

Another option would be to not estimate any confidence intervals and just expect that the consistency properties of the estimator will give us a good estimate. As a statistician, having no metric of precision at all would make me nervous. Granted, the sample size is quite large.

⁵“Coverage” used with this meaning: https://en.wikipedia.org/wiki/Coverage_probability

910 13 Sampling Strategy and Data Features that Support Estimation

911 Earlier I said that we have no “ground truth” about the different distribution components F_k . Even with no ground
 912 truth, the data has certain features that can serve as a check on the estimation.

913 The plan is to take each week from September 1st, 2021 to October 31, 2022 as
 914 a sample. Monero version 0.17.2.3, released September 1, 2021 was the last time
 915 that the standard decoy selection algorithm was changed. Starting the sample in
 916 September 2021 would also avoid the apparent July-August 2021 flood incident that
 917 distorted the ring member age distribution ([Krawiec-Thayer et al., 2021]).

918 There is a clear weekly cycle in the on-chain transaction volume data. Taking each
 919 week as the estimation unit would avoid the need to deal with the cycles explicitly.
 920 The average number of rings (i.e. transaction inputs) per week since September 2021
 921 has been about 375,000, according to the `monero-blockchain-stats` command line
 922 utility. This sample size is probably large enough for acceptable precision of my
 923 proposed estimator. Two-step estimators tend to have larger variance. The fact that
 924 there is a nonparametric kernel estimator involved also raises variance. If 375,000
 925 observations are not enough, then we can take samples 4 weeks at a time, making the
 926 sample size 1.5 million rings. It is crucial to break the sample into time units rather
 927 than pooling the entire sample together because the real spend age distribution is
 928 expected to change over time.

929 13.1 Leveraging the Hard Fork

930 I want to have at least two months of post-hard fork data for a few different reasons.
 931 In my OSPEAD CCS proposal I wrote⁶:

932 The upcoming hard fork, which does not yet have a fixed date, will include
 933 an increase in the ring size. The discontinuity that the hard fork creates
 934 can be leveraged to better understand how ring signatures work in practice
 935 [sic] on the Monero blockchain. Therefore, some of the research work will
 936 occur after the hard fork.

937 First, it is likely that the larger ring size will shrink the confidence intervals of the estimates, as I discussed in
 938 **Section 12 Confidence Interval Estimation**. Second, the hard fork caused user and software disruption that
 939 could generate “cleaner” estimates. Table 1 lists the daily on-chain transaction volume around the time of the hard
 940 fork. The day before the August 13 hard fork, transaction volume rose to about double normal levels. The day
 941 after the hard fork, transaction volume plummeted by about 2/3rds.

942 A large number of wallet implementations did not update in time for the hard fork. Until they were updated,
 943 no transactions constructed from those wallets would have appeared on the blockchain. Some of these wallets are
 944 known to use non-standard decoy selection algorithms, such as MyMonero, Edge, and Exodus.⁷ The decoy selection
 945 algorithms of some of the other wallets is unknown. Table 2 contains the approximate dates that some of these
 946 wallets could begin constructing valid transactions again. During the period of time that these wallets were out of
 947 commission, the proportion of transactions on the blockchain with the standard decoy selection algorithm, i.e. the
 948 mixing proportion $\omega_{\text{wallet}2}$ in Equation (6), would have been much higher. Therefore, it would likely be easier to
 949 estimate $F_{\text{wallet}2}$ without the distracting other components F_k during this immediate post-hardfork period.

Table 1: Transactions Per Day Around the Hard Fork

Date	Transactions
2022-08-06	21,315
2022-08-07	23,297
2022-08-08	29,083
2022-08-09	27,178
2022-08-10	32,995
2022-08-11	25,566
2022-08-12	61,305
2022-08-13	47,536
2022-08-14	9,288
2022-08-15	18,289
2022-08-16	20,555
2022-08-17	27,398
2022-08-18	29,639
2022-08-19	40,914
2022-08-20	23,031
2022-08-21	20,702

Source:

`monero-blockchain-stats`

⁶<https://ccs.getmonero.org/proposals/Rucknium-OSPEAD-Fortifying-Monero-Against-Statistical-Attack.html>

⁷<https://github.com/monero-project/research-lab/issues/99>

Table 2: Delayed Hardfork-compatible Wallets as of August 30, 2022

Fix date	Wallet	Source
August 19	WooKey	https://github.com/WooKeyWallet/monero-wallet-android-app/releases/tag/v2.2.1
August 25	Exodus	https://twitter.com/exodus_io/status/1562918301181034496
August 27	Edge	https://twitter.com/EdgeWallet/status/1563584457361149952
August 29	MyMonero	https://twitter.com/MyMonero/status/1564149853478760448
Not fixed	Atomic	https://www.reddit.com/r/atomicwallet/comments/x0wljr/xmr_balance_showing_zero_i_sent_a_transfer_a_few/
Not fixed	Coinomi	https://www.reddit.com/r/COINOMI/comments/wrqvw7/monero_node_is_down_need_to_export_the_mnemonic/
Not fixed	Guarda	https://libera.monerologs.net/monero-dev/20220831#c143060
Not fixed	Zelcore	https://twitter.com/CactusHashAZ/status/1564020934482178048#m

950 13.2 MyMonero Fee Fingerprinting

951 jberman discovered that the MyMonero wallet used a method to calculate transaction fees that is different from
952 `wallet2`'s method.⁸ jberman concluded, “Thus, this difference from `wallet2` should be fingerprintable on chain.”
953 MyMonero also uses a distinct decoy selection algorithm.⁹ The fee fingerprint can be used in two ways. First, the
954 proportion of MyMonero transactions $\omega_{MyMonero}$ in Equation (6) can be calculated exactly for each week rather
955 than estimated with the Bonhomme-Jochmans-Robin estimator. This can serve as a “ground truth” check on the
956 estimated $\hat{\omega}_{MyMonero}$ that would hopefully reveal any problems with the Bonhomme-Jochmans-Robin estimator in
957 this setting.

958 The second way that the MyMonero fee fingerprinting could be useful is in partitioning the sample. If the
959 MyMonero component of the mixture distribution is completely excluded from the first step estimate, then a certain
960 source of noise is removed. We could exclude these MyMonero transactions and be done with it if we wanted to
961 construct the OSPEAD decoy selection algorithm based only on the real spend age distribution of users of the
962 `wallet2` DSA. If we wanted to use the combined `wallet2` and MyMonero real spend age distribution as envisaged
963 in equation (7), then we could still exclude MyMonero from the main estimation, but then take the MyMonero-only
964 sample and apply the Patra-Sen Inversion Estimator to it directly, skipping the first step Bonhomme-Jochmans-
965 Robin estimator.

966 13.3 Nonstandard Wallets are Nonstandard in Multiple Ways

967 The sample partitioning idea can be applied with more wallets than just MyMonero. Wallets with a nonstandard
968 DSA are also more likely to have nonstandard fee calculation, `tx_extra` contents, and possibly other transaction
969 uniformity defects. Through querying databases provided by neptune I have tentatively discovered some `tx_extra`
970 formats associated with nonstandard DSAs.

971 Once labeled as such, those transactions can be excluded from the estimation because they contribute only
972 noise. With the estimation approach I have outlined above, we can only recover the real spend age distribution
973 from transactions that use a known decoy selection algorithm. If a wallet implementation is using an unknown non-
974 standard DSA and distinguishes its transactions with other nonstandrad features, its transactions can be removed
975 from the sample before estimation.

976 13.4 `wallet2` Integer Truncation Bug

977 Between Monero versions v0.14.1.0 and v0.17.2.2, `wallet2`'s DSA would integer-truncate a value that dictated the
978 spread of the decoy distribution. This triggered a discontinuity in the distribution at the point in time that the
979 value would be truncated to the next integer. The triggering is apparent at approximate block height 2,381,000 in

⁸<https://github.com/mymonero/mymonero-core-cpp/pull/36>

⁹<https://www.getmonero.org/2021/09/20/post-mortem-of-decoy-selection-bugs.html>

980 June 2021 in some of the plots of [Krawiec-Thayer et al., 2021]. I do not plan to investigate the trigger point since
 981 the reward is probably not worth the effort. If we wanted to investigate it, would could get an approximate number
 982 for the proportion of transaction constructed with `wallet2` around the time of the truncation trigger.

983 14 Future Work to Improve the Two-Step Estimator

984 A general note about estimators: If a single valid (i.e. consistent) estimator of a particular estimand exists, then
 985 multiple valid estimators exist almost always. In general what distinguishes the estimators may be the assumptions
 986 required for their validity and their variance and bias. My proposed two-step estimator is probably not the lowest-
 987 variance estimator of the estimand. The estimator “leaves information on the table” in a number of ways. In
 988 general, when more true information (do not impose false assumptions if avoidable!) about the structure of the
 989 statistical problem is considered by the estimator, its estimates become more precise. For example, when more true
 990 over-identifying restrictions are imposed by Generalized Method of Moments estimators, variance usually falls. I
 991 give three ways that the precision of the estimator can be improved by leveraging more information.

992 First of all, my proposed estimator does not use the fact that all rings in the same transaction almost certainly use
 993 the same decoy selection algorithm and therefore all belong to the same F_k component. (Note, however, that rings
 994 can be cached when a transaction is first constructed but not broadcast, freezing the ring in time.)¹⁰ Every estimator
 995 that I have reviewed that attempts to estimate mixture models with repeated measurements nonparametrically,
 996 i.e. that attempts to estimate the model of Equation (4), has an equal number of observations for repeated
 997 measurements. If we could take into account the hierarchical transaction structure of the data, then the number of
 998 repeated measures $M + 1$ could vary according to the number of rings in a given transaction. As seen in **Section**
 999 **12 Confidence Interval Estimation**, a larger $M + 1$ increases the precision of the Bonhomme-Jochmans-Robin
 1000 estimator.

1001 The second way is to convert the two-step estimator into a single-step joint estimator in which a single objective
 1002 function is optimized to produce estimates in both steps simultaneously. Among maximum likelihood estimators,
 1003 such an estimator is known as “Full Information Maximum Likelihood” (FIML) [Murphy & Topel, 1985]. Creating
 1004 such an estimator would likely require a great deal of effort and skill in theoretical statistics.

1005 The third way may or may not be mutually exclusive with the second. In the first step, we do not exploit the
 1006 fact that we know quite a bit about several of the F_k components — specifically, that the component distributions
 1007 have an admixture/contamination structure with a known parametric form of the largest sub-component (i.e. the
 1008 DSA) for several of the components. There may be a way to include this information from the second step into the
 1009 first step estimator.

1010 There are a few recent papers on using Bayesian methods to model “mixture of mixtures” or “hierarchical mixture
 1011 models”([Malsiner-Walli et al., 2017], [Argiento et al., 2020], [Frühwirth-Schnatter et al., 2021]).

1012 I consider these possible improvements to the estimator to be out of scope for this initial OSPEAD CCS. The
 1013 gain is not worth the effort at this time.

1014 15 Rucknium Ratio Attack/MAP Decoder/Discriminant Analysis

1015 An inebriated Homer Simpson once said, “To alcohol! The cause of, and solution to, all of life’s problems.” Similar
 1016 to alcohol, an estimate of the real spend age distribution is the cause of, and solution to, many of Monero’s statistical
 1017 attack problems.

1018 The potency of statistical attacks is a critical overall factor influencing the disclosure decision in **Section 17**
 1019 **Disclosure Considerations**. [Ronge et al., 2021] give the following metrics of statistical attack potency:

¹⁰<https://libera.monerologs.net/monero-dev/20220622#c111774>

1020 **Definition 3.1.** *The guessing probability of X is*

$$\text{Guess}(X) := \max_X \Pr[X = X].$$

1021 This gives an upper bound on the probability than an (unbounded) adversary can “guess” the value of
1022 the random variable X correctly.

1023 **Definition 3.2.** *The (average) conditional guessing probability of X given Y is defined as*

$$\text{Guess}(X|Y) := \sum_Y \Pr[Y = Y] \max_X \Pr[X = X | Y = Y].$$

1024 This gives an upper bound on the probability that an (unbounded) adversary “guesses” the value of the
1025 random variable X correctly when given a sample of Y

1026 **Theorem 4.4.** *Let \mathcal{A} be any computationally unbounded adversary, who inputs a ring $\Pi(S)$ (where Π
1027 [the ring sampler, i.e. Decoy Selection Algorithm,] is possibly subverted by \mathcal{A}) and some leakage $\Lambda(S)$,
1028 where S is sampled from the distribution \mathcal{S} (possibly influenced or specified by \mathcal{A}), and outputs a guess
1029 S' . The probability of \mathcal{A} correctly guessing the signer S , i.e., $S' = S$, is upper bounded by*

$$\text{Guess}(S|\Pi(S), \Lambda(S)) = 2^{-\alpha(S, \Pi, \Lambda)}. \quad (11)$$

1030 The right hand side of Equation (11) is not important for this discussion. $\Lambda(S)$, which represents side-channel
1031 information like view keys, logs of centralized exchanges, EAE/EABE attacks, etc., is out of scope. Therefore, we
1032 will deal with the simpler $\text{Guess}(S|\Pi(S))$. Possible active attacks on the DSA by adversary \mathcal{A} like flooding/black
1033 marble attacks are also out of scope.

1034 In my HackerOne submission I performed a Monte Carlo simulation to approximate the probability of success
1035 of the “Rucknium Ratio Attack” when the decoy selection algorithm diverges from the real spend age distribution,
1036 where I found $E[\widehat{\text{Guess}}_{RRD}(S|\Pi(S))] \approx 0.35$. The Rucknium Ratio Attack guesses that the most likely real spend
1037 is the ring member, compared to all other ring members in a ring, whose age x is at the maximum ratio $\frac{f_S(x)}{f_D(x)}$,
1038 i.e. the ratio of the real spend distribution to the decoy distribution. Unknown to me at the time, a few months
1039 prior [Aeeneh et al., 2021] had published an essentially identical attack they called the “Maximum A Posteriori
1040 (MAP) Decoder”. They did not attempt to estimate the empirical potency of the attack because they did not
1041 think obtaining an estimate of the real spend age distribution $f_S(x)$ was possible. I will refer to this attack as
1042 the MAP Decoder since it is more descriptive of the technique and [Aeeneh et al., 2021] were first to publish. The
1043 Maximum A Posteriori Decoder is related to the Maximum A Posteriori (MAP) estimator. The MAP estimator
1044 find the statistical mode (i.e. the most frequent value) of the posterior distribution for a particular distribution
1045 [Bassett & Deride, 2019]. The MAP Decoder finds the most frequent value for the ratio of two distributions.

1046 In their Equation 13, [Aeeneh et al., 2021] state the closed-form expression for the error probability of their
1047 attack. I will translate it into our notation. M is the number of decoys. Let $R = \{r_1, r_2, \dots, r_{M+1}\}$ be the set
1048 of ring members that have corresponding block height $B = \{b_1, b_2, \dots, b_{M+1}\}$. Without loss of generality, assume
1049 that the real spend is r_1 . Let J be the number of blocks that have RingCT outputs. In other words, J is the
1050 maximum age of a possible ring member, in units of block height. Then when the real spend r_1 is at block height
1051 b_1 , the probability that the MAP Decoder correctly guesses the real spend is:

$$\left(\sum_{j=1}^J f_D(j) \mathbf{1} \left\{ \frac{f_S(j)}{f_D(j)} < \frac{f_S(b_1)}{f_D(b_1)} \right\} \right)^M \quad (12)$$

1052 The $\mathbf{1} \left\{ \frac{f_S(j)}{f_D(j)} < \frac{f_S(b_1)}{f_D(b_1)} \right\}$ indicator function term checks if a potential decoy drawn from the j th block would have
1053 a lower $\frac{f_S(x)}{f_D(x)}$ value than the real spend r_1 . If the potential decoy does not have a lower value than the real spend,

then the real spend would not be guessed by the MAP Decoder. The $f_D(j)$ weights each $\mathbf{1} \left\{ \frac{f_S(j)}{f_D(j)} < \frac{f_S(b_1)}{f_D(b_1)} \right\}$ term by the probability that the decoy selection algorithm would choose a decoy from that j th block height. Then the $\sum_{j=1}^J$ operator sums each weighted indicator function for every possible J blocks that a decoy could be chosen from. Then for the whole expression we need to calculate the joint probability for every decoy selected in the ring. We can take the M th power of the expression for M decoys since the decoy selection algorithm (or at least that of **wallet2**) chooses decoys in an independent and identically distributed manner.

Expression 12 gives the probability that the MAP Decoder correctly guesses the real spend for a particular real spend r_1 at block height b_1 . To calculate the expectation of $\text{Guess}_{MAPD}(S|\Pi(S))$ (i.e. the average) over all real spends, we weight Expression 12 by the probability that a user spends from each block height age $f_S(x)$ and then sum:

$$E[\text{Guess}_{MAPD}(S|\Pi(S))] = \sum_{b_1=1}^J f_S(b_1) \left(\sum_{j=1}^J f_D(j) \mathbf{1} \left\{ \frac{f_S(j)}{f_D(j)} < \frac{f_S(b_1)}{f_D(b_1)} \right\} \right)^M \quad (13)$$

Equation (13) apparently requires summation of J^2 elements. Given that J is in excess of one million, computation of (13) could be computationally expensive. However, by intelligently ordering the quantities to be compared, the computational expense can be reduced by several orders of magnitude. In **map-decoder-attack-potency.R** I perform the calculations with a naive implementation and an efficient one. With estimates of $f_S(x)$ and $f_D(x)$ from my HackerOne submission, Equation (13) gives $E[\text{Guess}_{MAPD}(S|\Pi(S))] \approx 35.43693\%$, which is within about one hundredth of a percentage point of 35.4512%, the value I found from my original Monte Carlo simulation. When ring size is 16, $E[\text{Guess}_{MAPD}(S|\Pi(S))] \approx 30.39\%$ with my HackerOne estimates of $f_S(x)$ and $f_D(x)$.

It is beneficial to have the Equation (13) closed-form expression of $E[\text{Guess}_{MAPD}(S|\Pi(S))]$ with an efficient implementation rather than needing to resort to Monte Carlo simulation. Minimizing $E[\text{Guess}_{MAPD}(S|\Pi(S))]$ is one of several criteria that I propose for determining the OSPEAD decoy selection algorithm later in **Section 21.6 Maximise Resistance to an Attack**. For a numerical minimization algorithm it is much better to have function evaluation computation time measured in seconds rather than minutes. There are also particular difficulties (though surmountable) with having the objective function be stochastic as would be the case with $E[\text{Guess}_{MAPD}(S|\Pi(S))]$ estimated with Monte Carlo simulation.

My HackerOne estimate of $E[\text{Guess}_{MAPD}(S|\Pi(S))]$ was based on the false simplifying assumption that the **wallet2** decoy selection algorithm was the only one being used on the Monero blockchain, using a simplified version of the Patra-Sen Inversion Estimator to obtain an estimate of $f_S(x)$. My proposed two-step estimator will produce a consistent estimate of $f_S(x)$ under the assumption that multiple decoy selection algorithms are being used. A more accurate estimate of $E[\text{Guess}_{MAPD}(S|\Pi(S))]$ could be produced from this improved estimate of $f_S(x)$. I am not sure if this more accurate estimate of $E[\text{Guess}_{MAPD}(S|\Pi(S))]$ would be greater than or less than the original HackerOne estimate. There are two factors that cut in opposite directions.

First, the original estimate would have conflated decoys from nonstandard decoy selection algorithms with real spends, potentially biasing the estimate of $E[\text{Guess}_{MAPD}(S|\Pi(S))]$ upward. The magnitude of this source of bias is dependent on the divergence between the **wallet2** DSA and the nonstandard DSAs, plus the proportion of on-chain transactions that use the nonstandard DSAs.

On the other hand, the original estimate of $E[\text{Guess}_{MAPD}(S|\Pi(S))]$ may have been biased downward because it did not take into account the variability of $f_S(x)$ over time. It gave an “average” over the time period under consideration. Assume for a moment that the **wallet2** DSA performed adequately for the average $f_s(x)$ over time. If $f_s(x)$ was highly variable from week to week, $E[\text{Guess}_{MAPD}(S|\Pi(S))]$ could still be high since the MAP Decoder could exploit the extremes of each week.

We will not have an improved estimate of $E[\text{Guess}_{MAPD}(S|\Pi(S))]$ until we have the improved estimate of $f_S(x)$.

1095 15.1 Attack Optimality

1096 [Aeeneh et al., 2021] claims in their Proposition 4 that the MAP Decoder attack is optimal in the sense that it
 1097 minimizes the (guessing) error probability. If there is no attack more powerful, then it may make sense to
 1098 construct the OSPEAD DSA by minimizing the MAP Decoder effectiveness as I suggest in **Section 21.6 Maximize**
 1099 **Resistance to an Attack.** It is not clear to me whether maximizing resistance to the most powerful attack also
 1100 maximizes resistance to a less powerful attack. For example, the “guess newest heuristic” of [Kumar et al., 2017] and
 1101 [Möser et al., 2018] is a less powerful attack that exploits *part* of the information about the discrepancy between
 1102 $f_S(s)$ and $f_D(x)$.

1103 In formal terms: Let $\Pi_{MAPD}(S)$ be the DSA that is constructed to minimize $E[\text{Guess}_{MAPD}(S|\Pi(S))]$ (with the
 1104 constraint that $\Pi_{MAPD}(S)$ is static and parametric). Let $\Pi_{Alt}(S)$ be some alternative DSA that does not minimize
 1105 $E[\text{Guess}_{MAPD}(S|\Pi(S))]$. Let $\text{Guess}_{MAPD_{inferior}}$ be the guessing probability of an attack $MAPD_{inferior}$ that is
 1106 inferior to the Map Decoder. Does there exist a $\Pi_{Alt}(S)$ and $MAPD_{inferior}$ such that

$$E[\text{Guess}_{MAPD}(S|\Pi_{MAPD}(S))] < E[\text{Guess}_{MAPD}(S|\Pi_{Alt}(S))]$$

1107 (which is true by Proposition 4 of [Aeeneh et al., 2021]), yet also

$$E[\text{Guess}_{MAPD_{inferior}}(S|\Pi_{MAPD}(S))] > E[\text{Guess}_{MAPD_{inferior}}(S|\Pi_{Alt}(S))]? \\$$

1108 If the answer is in the affirmative and an adversary can implement $MAPD_{inferior}$ but not $MAPD$, then it
 1109 could be better to adopt something other than $\Pi_{MAPD}(S)$ as the new DSA. In such case, it might be better to use
 1110 one of the other criteria I develop in **Section 21 Criteria for Best Fit**. That said, in the following discussion I
 1111 lay out an argument that the MAP Decoder attack may have deep optimality properties.

1112 My search of the statistics literature suggested that the MAP Decoder/Rucknium Ratio Attack is just a specific
 1113 instance of a very old problem and class of solutions: Discriminant Analysis, a type of supervised classification
 1114 ([McLachlan, 1992] [Fraley & Raftery, 2002]). As early as 1951, the year that [Fix & Hodges, 1989] was originally
 1115 written under a U.S. Air Force contract, this problem was considered to be, “in a sense, completely solved.” In 1989
 1116 [Fix & Hodges, 1989] was finally published in a journal due to its importance. After consistent estimation of the
 1117 real spending distribution, I take our problem as falling into the category of “subproblem (i)” below:

1118 The discrimination problem (two population case) may be defined as follows: a random variable Z , of
 1119 observed value z , is distributed over some space (say, p -dimensional) either according to distribution F ,
 1120 or according to distribution G . The problem is to decide, on the basis of z , which of the two distributions
 1121 Z has.

1122 The problem may be classified in various ways into subproblems. One pertinent method of classification
 1123 is according to the amount of information assumed to be available about F and G . We may distinguish
 1124 three stages:

- 1125 (i) F and G are completely known;
- 1126 (ii) F and G are known except for the values of one or more parameters;
- 1127 (iii) F and G are completely unknown, except possibly for assumption about the existence of
 densities, etc.

1129 Subproblem (i) has been, in a sense, completely solved. The solution is implicit in the Neyman-Pearson
 1130 lemma (Neyman & Pearson, 1936), and was made explicit by Welch (1939). We may without loss of
 1131 generality assume the existence of density functions, say f and g , corresponding to F and G , since F
 1132 and G are absolutely continuous with respect to $F + G$. If f and g are known, the discrimination
 1133 should depend only on $f(z)/g(z)$. An appropriate (positive) constant c is chosen, and the following rule
 1134 is observed:

1135 if $f(z)/g(z) > c$, we decide in favor of F ;
 1136 if $f(z)/g(z) < c$, we decide in favor of G ;
 1137 if $f(z)/g(z) = c$, the decision may be made in an arbitrary manner.

1138 These procedures are known to have optimum properties with regard to control of probability of mis-
 1139 classification (probability of wrong decision). We shall refer to this as the 'likelihood ratio procedure',
 1140 and denote it by $L(c)$

1141 The choice of c depends on considerations related to the relative importance of the two possible errors:
 1142 saying Z is distributed according to G when in fact it is distributed to F , and conversely. Two choice
 1143 of c have been widely advocated:

- 1144 (a) take $c = 1$;
 1145 (b) choose c so that the two probabilities of error are equal.

1146 Choice (a) has been called 'logical'; choice (b) yields the minimax procedure. In this paper we shall
 1147 not concern ourselves with the choice of c , but shall assume that a given positive c is a datum of the
 1148 problem.

1149 "If f and g are known, the discrimination should depend only on $f(z)/g(z)$." This 70-year-old principle is used in
 1150 the MAP Decoder attack. The main difference between MAP Decoder and this passage from [Fix & Hodges, 1989]
 1151 is that the inequality criterion $f(z)/g(z) \geq c$ is not governed by a constant c . In the MAP Decoder attack, c is
 1152 variable, determined by the ratio $f(z)/g(z)$ for every other ring member in the ring, i.e. $c = \max_i\{f(z_i)/g(z_i)\}$.

1153 [McLachlan et al., 2019] give a more general and modern explanation of the principle, placing it within the
 1154 framework of finite mixture models:

1155 The mixture model give by Equation (3) can be used to provide a model-based approach to clustering
 1156 of the observed data in \mathbf{y}_{obs} by conceptualizing that $\mathbf{y}_1, \dots, \mathbf{y}_n$ come from a mixture in proportions
 1157 π_1, \dots, π_g of g groups G_1, \dots, G_g in which \mathbf{Y}_j has density $f_1(\mathbf{y}_j; \boldsymbol{\theta}_1), \dots, f_g(\mathbf{y}_j; \boldsymbol{\theta}_g)$, respectively. This
 1158 is irrespective of whether these groups do externally exist.

1159 For clustering purposes, each component in the mixture model given by Equation (3) is usually taken
 1160 to correspond to a cluster. The posterior probability that the j th observation \mathbf{y}_j arose from group G_i
 1161 is given by the fitted posterior probability,

$$\tau_i(\mathbf{y}_j; \hat{\boldsymbol{\Psi}}) = \hat{\pi}_i f_i(\mathbf{y}_j; \hat{\boldsymbol{\theta}}_i) / f_i(\mathbf{y}_j; \hat{\boldsymbol{\Psi}}) \quad (i = 1, \dots, g; j = 1, \dots, n), \quad (14)$$

1162 where $\hat{\boldsymbol{\Psi}}$ denotes an estimate of $\boldsymbol{\Psi}$.

1163 A probabilistic clustering of the data $\mathbf{y}_1, \dots, \mathbf{y}_n$ into g clusters can be obtained in terms of the fitted
 1164 posterior probabilities of component membership $\tau_i(\mathbf{y}_j; \hat{\boldsymbol{\Psi}}) (i = 1, \dots, g)$.

1165 An outright partitioning of the observations into g nonoverlapping clusters C_1, \dots, C_g is effectuated
 1166 by assigning each \mathbf{y}_j to the group G_i to which it has the highest estimated posterior probability of
 1167 belonging. That is, the i th cluster C_i contains those observations \mathbf{y}_j with $\hat{z}_{ij} = (\hat{\mathbf{z}}_j)_i = 1$, where

$$\begin{aligned} \hat{z}_{ij} &= 1, && \text{if } i = \arg \max_h \hat{\tau}_h(\mathbf{y}_j; \hat{\boldsymbol{\Psi}}), \\ &= 0, && \text{otherwise.} \end{aligned} \quad (15)$$

1168 As the notation implies, \hat{z}_{ij} can be viewed as an estimate of z_{ij} which, under the assumption that the
 1169 observations come from a mixture of g groups G_1, \dots, G_g , is defined to be one or zero according to
 1170 whether \mathbf{y}_j did or did not arise from $G_i (i = 1, \dots, g)$.

1171 The above rule for assigning the unclassified data points \mathbf{y}_j to the g groups corresponds to the Bayes
 1172 rule of allocation in the supervised classified case with known parameter vector $\boldsymbol{\Psi}$ ([McLachlan, 1992],
 1173 section 1.3).

1174 The $\hat{\pi}_i$ in equation (14) is the estimated mixing proportion of the i th density component. If we were to apply the
 1175 assignment rule of (15) to a ring to classify ring members into the real spend and decoy distribution, the decoy
 1176 distribution would always be heavily weighted since $\pi_D = \frac{M}{M+1}$, where M is the number of decoys. Thus, in many
 1177 cases all ring members would be classified as decoys if this particular Bayes rule of allocation were to be used.

1178 The general idea of this Bayes rule is packaged in a piece of advice given to physicians about diagnosis of disease:
 1179 “When you hear hoofbeats in the night, look for horses — not zebras.”¹¹ In other words, when you have some limited
 1180 symptom data that is consistent with two diseases, i.e. $f_1(\mathbf{y}_j; \hat{\theta}_1) = f_2(\mathbf{y}_j; \hat{\theta}_2)$ for symptom data \mathbf{y}_j and diseases 1
 1181 and 2, choose the diagnosis of the disease that is more prevalent in the population, i.e. choose the disease with the
 1182 larger mixing proportion π_i .

1183 Our problem is different in that we can use information from all $M+1$ observations in a set of repeated measures
 1184 — all members of a ring — for classification. The general term for this setting is the “compound decision problem”
 1185 [Copas, 1974]. To extend the metaphor, we know that there is exactly one zebra in the corral (there is exactly one
 1186 real spend in a Monero ring). To re-imagine the MAP Decoder in the framework of the (15) assignment rule, we
 1187 can instead make an assignment (a guess) based on:

$$\begin{aligned}\hat{z}_{ij} &= 1, && \text{if } j = \arg \max_k \hat{\pi}_i(\mathbf{y}_k; \hat{\Psi}), \\ &= 0, && \text{otherwise.}\end{aligned}$$

1188 In other words, apply the decision rule over the rows (i th dimension where i denotes the distribution component)
 1189 of the \hat{z}_{ij} matrix rather than over the columns (j th dimension where j denotes the observation). Thus, we consider
 1190 each distribution (in our problem the two distributions: real and decoy) and then assign the most likely observation
 1191 (ring member) to that distribution — we can ignore the decoy assignment since we are interested in the assignment
 1192 of the real spend — rather than considering each observation and assigning the most likely distribution to it.

1193 With a quick search I have not found a paper that covers our exact case of a compound decision problem in
 1194 which the proportion of the sample belonging to each component is known exactly. Diving further into this corner
 1195 of statistics may uncover more properties of the MAP Decoder such as how to maximize its error rate when one of
 1196 the distributions, i.e. the decoy distribution, can be partially manipulated.

1197 To conclude, it is likely that the MAP Decoder attack is the most dangerous statistical attack that exploits the
 1198 timestamp information on the Monero blockchain. Therefore, I am personally leaning toward using the minimization
 1199 of $E[\text{Guess}_{MAPD}(\mathbf{S}|\Pi(\mathbf{S}))]$ as the criteria to select the OSPEAD decoy selection algorithm.

1200 16 Postestimation Classification of On-Chain Rings into DSAs

1201 Given estimates of the mixture model components F_k , is it possible to probabilistically classify rings on the Monero
 1202 blockchain according to their associated decoy selection algorithm? I believe it may be possible. There are two
 1203 main reasons we might be interested in doing this.

1204 The first is to better understand the degree of defects in transaction uniformity on the blockchain — the
 1205 “anonymity puddles.” Transaction uniformity defects themselves can be used as ammunition in statistical attacks
 1206 against Monero user privacy. The second reason is to evaluate the safety of implementing an OSPEAD-derived
 1207 DSA in a new release of the Monero `wallet2` software at a time other than a hard fork. It could be possible to
 1208 probabilistically distinguish between transactions constructed with the new DSA and the old DSA, which could
 1209 endanger user privacy. Without a hard fork, users who use the GUI or CLI wallet and wallet developers who use
 1210 `wallet2` to construct transactions are not required to upgrade to the new version. The question is whether the
 1211 cure would be worse than the disease:

¹¹<https://quoteinvestigator.com/2017/11/26/zebras/>

$$\Pr(MAPD_{\text{wallet2_DSA}}) \stackrel{?}{\leq} \Pr(MAPD_{\text{OSPEAD_DSA}} \cup \text{Classifier}_{\text{wallet2_DSA}/\text{OSPEAD_DSA}}) \quad (16)$$

where $MAPD_{\text{wallet2_DSA}}$ is the event that the real spend is guessed by the MAP Decoder attack with the existing DSA, $MAPD_{\text{OSPEAD_DSA}}$ is the same with the OSPEAD-derived DSA, and $\text{Classifier}_{\text{wallet2_DSA}/\text{OSPEAD_DSA}}$ is the event that a DSA-based classifier correctly guesses the real spend by distinguishing between transactions using the existing DSA and the OSPEAD-derived DSA. Further data analysis is required because the direction of the inequality in (16) must be resolved empirically rather than theoretically.

Inequality (16) is not very formal. It does not consider the “adoption curve” over time of users updating to a new `wallet2` version. A more complete decision criteria would consider the risk over time and an educated guess about the adoption curve. An educated guess could be based on an estimate of the rate of adoption of the bug fixes to the decoy selection algorithm written by jberman in 2021, using the Bonhomme-Jochmans-Robin estimator.

There are a few candidate techniques for probabilistic classification of rings using their repeated measures structure. [Bagui et al., 2006] describe the state of the literature as of 2006:

Various nonparametric procedures are available in situations where independent repeated measurements are available on the entity to be classified; they have been considered by Das Gupta [7], Hudimoto [15], Kanazawa [18], Gupta and Kim [13], Lin [17], Govindarajulu and Gupta [12], and Haack and Govindarajulu [14] among others. DasGupta [7] proposed an allocation rule using the Wilcoxon statistic. Govindarajulu and Gupta [12] considered an allocation rule based on linear rank statistics for the s -group problem....

Bagui [2] considered the problem of classification of m independent multiple observations under the mixture sampling scheme which is common in prospective studies and diagnostic situations. Bagui et al. [3] proposed the allocation of multiple observations using a sub-sample approach under mixture sampling. Bagui and Mehra [4] used a rank nearest neighbor (RNN) classification rule to classify univariate independent multiple observations between two classes. The present article deals with the classification of m independent (multivariate) multiple observations under a separate (or conditional) sampling scheme.

I consider exploration of ring classification to be out of scope of this OSPEAD CCS. It would require much additional research. Future work could cover this ground.

17 Disclosure Considerations

Given the potential threat to retroactive user privacy of an estimate of the real spend age distribution combined with the MAP Decoder attack, there is a question of whether to disclose the estimator. A sophisticated adversary who pays close attention likely would already be aware of the MAP Decoder attack since it was published in [Aeeneh et al., 2021]. All that is missing is an estimate of the real spend age distribution. The decoy selection algorithm of course would be plainly visible in Monero’s source code, so it is a matter of disclosing how the parameters of the algorithm was chosen.

The process to decide on disclosure is not very formal at this point. It is not clear what exact body of individuals plays the role of deciding to disclose or not. After I completed my initial investigations in August/September 2021 I submitted my findings through HackerOne in compliance with Monero’s Vulnerability Response Process. After discussion, the VRP manager deferred to my judgment about whether to publicly disclose immediately. I chose not to disclose because the risk to user privacy was still not completely understood and we had no fix ready. It would make sense for the revoew panel to give input on about the advisability of disclosure.

As I stated in Section 15, we will not have an improved estimate of the potency of the MAP Decoder attack until we have the improved estimate of the real spend distribution Therefore, we do not need to make a final call at this point in time. Below I discuss some pros and cons of disclosure.

1254 **17.1 Pros**

1255 **17.1.1 Greater peer review of OSPEAD**

1256 Disclosure would enable wider peer review of my proposed solution. Disclosure of this document and my HackerOne
 1257 submission would potentially put more eyes on the OSPEAD estimator. Furthermore, if we choose disclosure I am
 1258 willing to prepare a manuscript for submission to a privacy-focused journal for a traditional peer review process.
 1259 We could then benefit from suggested improvements to OSPEAD.

1260 **17.1.2 More attention to this research area**

1261 Many papers such as [Egger et al., 2022][Ronge et al., 2021], and [Yu et al., 2019] have suggested a “partitioning”
 1262 (i.e. a single bin) DSA, partially since they did not think it was feasible to estimate the real spend age distribution.
 1263 Disclosure of OSPEAD could revitalize research on mimicking DSAs.

1264 Nonparametric mixture models are an active area of statistical research. With disclosure, we might be able to
 1265 catch the interest of theoretical statisticians to help us refine and strengthen the OSPEAD approach.

1266 An estimate of the real spend age distribution may also be useful to potential future research on churning.

1267 **17.1.3 User trust**

1268 Community trust in the security of the code and the privacy of its protocol is essential. Many community members
 1269 might not accept concealment of any part of the process of how certain designed decisions were determined. When
 1270 I submitted my CCS proposal, my suggestion that parts of OSPEAD be remain a secret was very controversial.¹²

1271 Community members’ concern is not without some basis in reality. In the past, malicious parties have used
 1272 “suggested” cryptographic parameters to effectively create backdoors, such as the suspected NSA backdoor in
 1273 Dual_EC_DRBG.¹³

1274 On the other hand, there is precedent for publication of certain security parameters without explicit justification
 1275 because the justification would reveal attacks, such as with the discovery of differential cryptanalysis:¹⁴

1276 The discovery of differential cryptanalysis is generally attributed to Eli Biham and Adi Shamir in the late
 1277 1980s, who published a number of attacks against various block ciphers and hash functions, including a
 1278 theoretical weakness in the Data Encryption Standard (DES). It was noted by Biham and Shamir that
 1279 DES was surprisingly resistant to differential cryptanalysis but small modifications to the algorithm
 1280 would make it much more susceptible.

1281 In 1994, a member of the original IBM DES team, Don Coppersmith, published a paper stating that
 1282 differential cryptanalysis was known to IBM as early as 1974, and that defending against differential
 1283 cryptanalysis had been a design goal. According to author Steven Levy, IBM had discovered differential
 1284 cryptanalysis on its own, and the NSA was apparently well aware of the technique. IBM kept some
 1285 secrets, as Coppersmith explains: “After discussions with NSA, it was decided that disclosure of the
 1286 design considerations would reveal the technique of differential cryptanalysis, a powerful technique that
 1287 could be used against many ciphers. This in turn would weaken the competitive advantage the United
 1288 States enjoyed over other countries in the field of cryptography.”

1289 **17.1.4 Probabilistic attacks may be irrelevant for certain threat models**

1290 For many user threat models, probabilistic/statistical attacks based on the divergence of the real spend age dis-
 1291 tribution from the decoy selection algorithm may not be relevant. Plausible deniability is still likely intact even
 1292 with statistical attacks. Due to Monero’s anonymous nature, we cannot really know the threat models that users
 1293 confront, but the existence of statistical attacks might not be a threat to most users’ privacy.

¹²https://www.reddit.com/r/Monero/comments/py8ub3/ccs_proposal_ospead_fortifying_monero_against/

¹³https://en.wikipedia.org/wiki/Dual_EC_DRBG

¹⁴https://en.wikipedia.org/wiki/Differential_cryptanalysis

[Deuber et al., 2022] attempt to formalize the reliability of certain classes of de-anonymization attacks and heuristic analysis of cryptocurrency transactions, including Bitcoin, Zcash, and Monero. They rate different techniques by relevance to law enforcement: “A high relevance means that the attack is very likely actually employed by law enforcement; medium means that it is not known whether the attack is used but we see potential use-cases; low means that we do not see a potential use-case because the exploited weakness has been fixed or the attack would require too much effort.” For example, they rate the “multi-output assumption” against Monero, which is statistical, as having “medium” relevance. They do not assess the relevance of decoy timing attacks since they consider the guess newest attack rectified with the implementation of the gamma distribution.

17.1.5 Monero adversaries may already be using these techniques

It is possible that chain analysis companies already have a good estimate of the real spend age distribution and are aware of the MAP Decoder attack. In such case, there would be little to gain from withholding the information. In 2020, the U.S. Internal Revenue Service awarded a contract to Chainalysis and Integra to develop Monero and BTC Lightning Network tracing tools.¹⁵ Ciphertrace claims to have Monero tracing tools.¹⁶ On the other hand, an Elliptic representative claimed that no chain analysis companies can trace Monero.¹⁷

17.1.6 An improved DSA mimics the real spend age distribution

If an adversary is aware of the MAP Decoder attack and is observing public discussion of the OSPEAD research, then they may assume that the new proposed public DSA closely mimics the real spend age distribution, which could enable the MAP Decoder attack. The main ambiguity that an adversary would deal with would be the question of how closely the new DSA mimics the real spend age distribution, the time period over which the distribution was fit, and how stable the distribution is over time. In other words, an adversary would not necessarily know how well a MAP Decoder attack would perform on past and future time periods of blockchain data.

17.2 Cons

17.2.1 Retrospective analysis

A principal challenge of cryptocurrency privacy, compared to privacy in other types of protocols, is that historical data exists forever publicly. Therefore, passive retroactive analysis is always a potential threat, regardless of whether any particular techniques existed at the time of the target transactions.

In a recent case, Swedish-Russian national Roman Sterlingov was arrested in the Los Angeles airport in 2021 on charges of running a BTC mixing service.¹⁸ A significant part of the evidence supporting the arrest is based on chain analysis of BTC transactions that occurred in 2011.

17.2.2 Probabilistic attacks may be relevant for extreme threat models

Certain users may be in environments where probabilistic analysis does not rise to the standard of proof for using legal investigative tools, filing charges, and/or conviction in a court of law. On the other hand, in other environments potential adversaries of Monero users may not feel themselves constrained by the rule of law. One-party states, dictatorships, absolute monarchies, and criminal gangs come to mind. Probabilistic analysis also opens the door to false positives, false accusations, and mis-targeted robberies. Even users in jurisdictions that ostensibly respect the

¹⁵<https://news.bitcoin.com/chainalysis-and-integra-win-1-25-million-irs-contract-to-break-monero/>

¹⁶<https://ciphertrace.com/heads-up-monero-v15-hard-fork-update/>

¹⁷<https://monero.observer/elliptic-liat-shetret-no-blockchain-analytics-company-can-track-monero/>

¹⁸<https://www.wired.com/story/bitcoin-fog-roman-sterlingov-blockchain-analysis/>

rule of law may be falsely targeted. There is a long history of the use of unreliable pseudoscience by law enforcement in such jurisdictions [Lilienfeld & Landfield, 2008].

17.2.3 The potency of probability-weighted transaction graph attacks is unknown

The unpublished MRL Research Bulletin #0011 raised the possibility of overlaying a set of probability weights on the edges of the Monero transaction graph to strengthen graph-based attacks on anonymity. I consider evaluating the risk of graph-based attacks to be out of scope of OSPEAD at this time. I personally do not have a strong grasp of graph theory.

Graph-based attacks may present a different risk from the “single hop” risk that my overall analysis focuses on. As such, it could present a blind spot whose danger we may fully realize only in the future. A few recent papers on non-probabilistic graph attacks make me less nervous about this blind spot, however. [Ye et al., 2020] executed the chain reaction attack of [Kumar et al., 2017] on newer data. See their Section 3.2.1. They found that the share of transactions that were fully or partially traceable using this attack fell to almost zero when RingCT was introduced in January 2017. [Vijayakumaran, 2021] performed the Dulmage-Mendelsohn (DM) decomposition on the Monero transaction graph, which “is optimal in the sense that it eliminates every output-key image association which is incompatible with the transaction history.” [Vijayakumaran, 2021] found “For pre-RingCT outputs in Monero, the DM decomposition technique performs better than existing techniques. For RingCT outputs in Monero, the DM decomposition technique has the same performance as existing techniques, with only five out of approximately 29 million outputs being identified as spent.”

[Egger et al., 2022] approach the question from a theoretical perspective. They attempt to put a bound on the necessary size of rings to defeat certain graph attacks when the decoy selection algorithm is a partitioning one rather than a mimicking one. I discussed my interpretation of the paper in the July 6, 2022 MRL meeting.¹⁹ They argued that to defeat global graph-based attacks in a transaction graph with $|U|$ transaction outputs (and a partitioning DSA), the ring size should be at least $\ln(2 \cdot |U|) + \sqrt{2 \ln(2 \cdot |U|)}$. The current $|U|$ of Monero of 75 million would require ring size 25 in this framework. From the preceding set of empirical and theoretical results we may tentatively hope that adding probability weights to graph attacks might not significantly harm Monero user privacy. Note that isthmus asserts “anyways, unfortunately Egger et al’s optimistic takeaway does not translate to the heuristics framework described in the doc, since NRL’s work uses transaction metadata to establish priors for partitioning the graph by fingerprint, whereas the partitioning samplers in their work are oblivious to that information”.²⁰

18 OSPEAD Requests to Monero C++ Developer(s)

Several tasks that would greatly help OSPEAD can only be completed by someone who can read and write C++ code, which I cannot. mj-xmr is willing and able to complete these tasks. The tasks are here in decreasing order of priority:

1. Converting the `wallet2` decoy selection algorithm into a closed-form probability density function, i.e. $f_{wallet2,D}$. OSPEAD cannot be done at all without this step. mj-xmr and I have made progress on this task already.²¹
2. Write a fast C++ implementation of the Bonhomme-Jochmans-Robin estimator. With the current MAT-LAB/Octave implementation, bootstrapping and Monte Carlo simulations are not computationally feasible. Therefore, we cannot really know the precision of the two-step estimator without a faster implementation. In

¹⁹<https://libera.monerologs.net/monero-research-lab/20220706#c117351>

²⁰<https://libera.monerologs.net/monero-research-lab/20220817#c137450>

²¹<https://github.com/mj-xmr/monero-mrl-mj/tree/master/decoy>

1367 Octave testing with a single CPU core, with $M = 10$ decoys, $N = 1000$ rings, $K = 3$ distinct decoy selection
 1368 algorithms, and $I = 10$ evaluation points, the estimation completed in about 24 hours. Assuming that the
 1369 estimation time scales linearly with N , for a single week of data a single CPU core would take about a year
 1370 of wall clock time to complete the estimation.

- 1371 3. Determine a method to fingerprint MyMonero transactions by their fees. According to jberman, “simplest/quickest way would be to run their fee algorithm at various times to get a range of what fees their
 1372 txs would be. hardest way would be to definitively calculate the plausible fees it could have chosen for each
 1373 tx. they can likely be pinpointed with effort”²² As I stated in **Section 13.2 MyMonero Fee Fingerprinting**, fingerprinted MyMonero transactions could help with the estimations.
- 1374
 1375 4. Converting the MyMonero decoy selection algorithm into a closed-form probability density function, i.e.
 1376 $f_{MyMonero,D}$. This is necessary for OSPEAD if we want to have a decoy selection algorithm that combines
 1377 the real spend age distribution of `wallet2` users and MyMonero users as in Equation (7). It would also
 1378 help overcome the “label swapping” problem discussed in **Section 12.2 Variance of the Second Step
 1379 Estimator**.
- 1380
 1381 5. A fast implementation of a forecast evaluator and simulator. mj-xmr has developed `tsqsim` in C++, which
 1382 compares very favorably to existing methods coded in slower languages, but certain modifications are needed
 1383 such as dealing with multivariate inputs. ²³

²²<https://libera.monerologs.net/monero-research-lab/20220615#c109279>

²³https://www.reddit.com/r/Monero/comments/tl5w1b/tsqsim_benchmark_new_tool_designed_for_monero/

₁₃₈₄ Part II

₁₃₈₅ Disclosed Portion

₁₃₈₆ 19 Plain English Explanation of the Problem

₁₃₈₇ 19.1 Timing Metadata in Monero Transactions

₁₃₈₈ Through great effort, Monero researchers and software developers have been able to conceal most information about
₁₃₈₉ transactions even though all Monero nodes share the blockchain database itself. However, transactions still contain
₁₃₉₀ some information that distinguishes them from one another.

₁₃₉₁ One key piece of information embedded in each Monero transaction is the time that it was confirmed in a mined
₁₃₉₂ block. Timestamps are an inherent property of blockchain-based cryptocurrencies. The timestamp ensures that
₁₃₉₃ coins are not double-spent.

₁₃₉₄ Each Monero transaction creates at least two *outputs*, which are amounts of Monero that can be spent in a
₁₃₉₅ subsequent transaction. When a user that has received an output wishes to spend it, the wallet software he or she
₁₃₉₆ uses will include that output in the set of outputs within the ring signature. Monero's current ring size is 16, so
₁₃₉₇ there will be 15 *decoy* outputs listed in the ring and 1 *real spend* output.

Table 3: Ring Members of Transaction 8bc3b73e48881d48cd69f1bf82a06b6a7b94bfd0dbf85edac535dab9270a305b

Output Public Key	Block Number	Block Timestamp
033a3d4b817bfaba6472fdb0dc6b00ecfe20983cea5e608fe5e679b0d23b	2666408	2022-07-13 13:27:42
866c4cf39af3e8785cbda21a9cbaa120424df76d5d05518d3ec888b94acf4c2b	2683943	2022-08-06 20:03:35
c8c3f8660776fad9dee0dd9bd6a3e02f52ba5e4b4785419db7a4131f1d1aa1ff	2692667	2022-08-19 02:10:18
ca2f709c371c21a8ab33cdab4b593c24dd52b0ca9e4835ddcbcf17f0b38e8c7	2696371	2022-08-24 06:06:50
1f11eb618ad7bb9a26a436af1190532c3744d0c0c3b335b6d3fd65b0a9911585	2697585	2022-08-25 21:24:57
6cd9e83b9f75e44a00310aec8935954a23bcdfa34c4fb5aa5e177212fd45ce95	2698303	2022-08-26 22:06:16
9ab87cffbc8eb76a0affbc9d004967ba1ca1ab008ebe68444d687a5d5a1cb145	2698869	2022-08-27 15:36:55
52ae53b9b67316c932db550c386fa7092952d65f9a9d52e261dc0f6ccba78c0	2699496	2022-08-28 12:52:01
19357e641d525fd001945d87829045e3ac6ad829ca92c90b729f152e7dc86b4a	2699840	2022-08-29 00:34:55
2bcdcb2b0b656cefefec925cbe122ecc0bb3d83bf61dfa0ac5a6c96a586ea0560	2699843	2022-08-29 00:39:49
56672f13802cca5363cd5def25c9af4be5ad2e625401d6381568902dca5b0b35	2700102	2022-08-29 09:31:16
bbae1b9ca63e017d030e295a5c9f92ed13152c82a32fdc10458912b9e0fc4e1a	2700241	2022-08-29 14:12:24
cc94142537e87fb36713c84c2bf8465bde03f442d155391025e078bac2154cef	2700606	2022-08-30 02:04:45
dda4c8f4c51b387e03b5ce97b187a5ccb4fb03ea622c62de5afb698f0351215	2700755	2022-08-30 06:59:12
846717c3b56dcf2c87a92b71890fdde13c72081d8d7e40ddbcba2e82d866c063	2700908	2022-08-30 12:10:55
6241dc9af536fd23584c11b6a9c84c9fd692fc5dbc9326741db6dde3f05832e2	2700918	2022-08-30 12:33:33

Source: <https://xmrchain.net/tx/8bc3b73e48881d48cd69f1bf82a06b6a7b94bfd0dbf85edac535dab9270a305b>

₁₃₉₈ The timestamp of all of those outputs in a ring is available to anyone running a full Monero node or even
₁₃₉₉ someone who knows how to use a web-based Monero block explorer. Table 3 contains an example of a ring with
₁₄₀₀ the time stamps of each ring member. Timing information is not the only data that is available to external
₁₄₀₁ observers of the Monero blockchain, but the timing data presents special challenges. Over a year, there will be
₁₄₀₂ about $30 \cdot 24 \cdot 365 = 262,800$ different timestamps (each timestamp corresponding to a unique block on the Monero
₁₄₀₃ blockchain). Since the number of timestamp values is so large compared to the number of decoys, timing information
₁₄₀₄ can help an anti-privacy adversary possibly narrow down which ring member is the real spend.

₁₄₀₅ Other types of data on the blockchain cannot take so many values in practice. The number of outputs of a
₁₄₀₆ transaction are limited to 16. Currently the maximum number of inputs of a transaction is about 146, but in practice

1407 the vast majority of transactions have fewer than 5 inputs.²⁴ Therefore, the amount of information contained in
 1408 timestamp data dwarfs the information from the number of inputs or outputs of a transaction. Transaction fee
 1409 and the `tx_extra` field in theory could have many distinguishing values. In practice, the “official” reference wallet
 1410 software allows users to set one of only 4 fee priority levels and restricts the format of data in the `tx_extra` field.
 1411 The difference between the timing data and other types of observable data in a Monero transaction is similar to
 1412 the difference between data on a person’s birth date and race. Birth date reveals much more information than race
 1413 since it can take so many values. Race can only be a pretty limited set of values. For an in-depth discussion of
 1414 metadata available to Monero blockchain observers, see [Krawiec-Thayer et al., 2021].

1415 19.2 Decoys Must Be Credible

1416 The core challenge of decoys is intuitive: Like decoys in all contexts, a decoy only serves its purpose if — to observers
 1417 — it looks like the real thing. Unfortunately, previous versions of Monero did a poor job of selecting decoys that
 1418 looked like the real thing. According to one estimate, 80 percent of Monero transactions prior to February 2017
 1419 could be traced simply by guessing that the youngest ring member was the real spend because selected decoys
 1420 tended to be much older than the real spends ([Möser et al., 2018]). Monero’s decoy selection algorithm has been
 1421 changed in recent years to correct flaws found by existing research, but there is still a lot of room for improvement.

1422 To explain the problem I will use the extreme case of what can happen if a given real spend has zero decoys
 1423 nearby in the timestamp distribution. Monero’s actual decoy selection algorithm does not suffer from this extreme
 1424 flaw, but it serves its purpose as an introduction to the issue.

Figure 2: Adequate Decoy Selection Algorithm

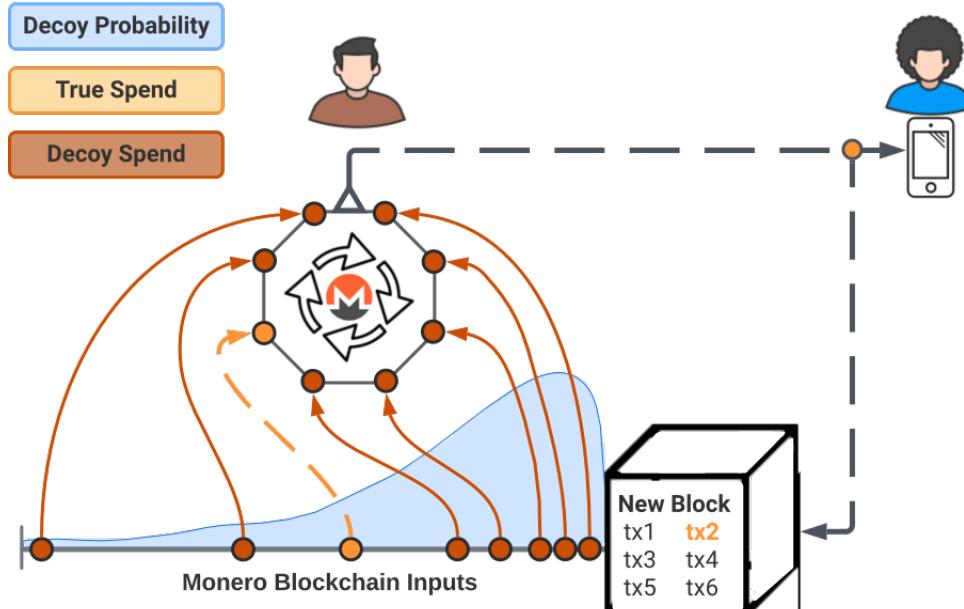


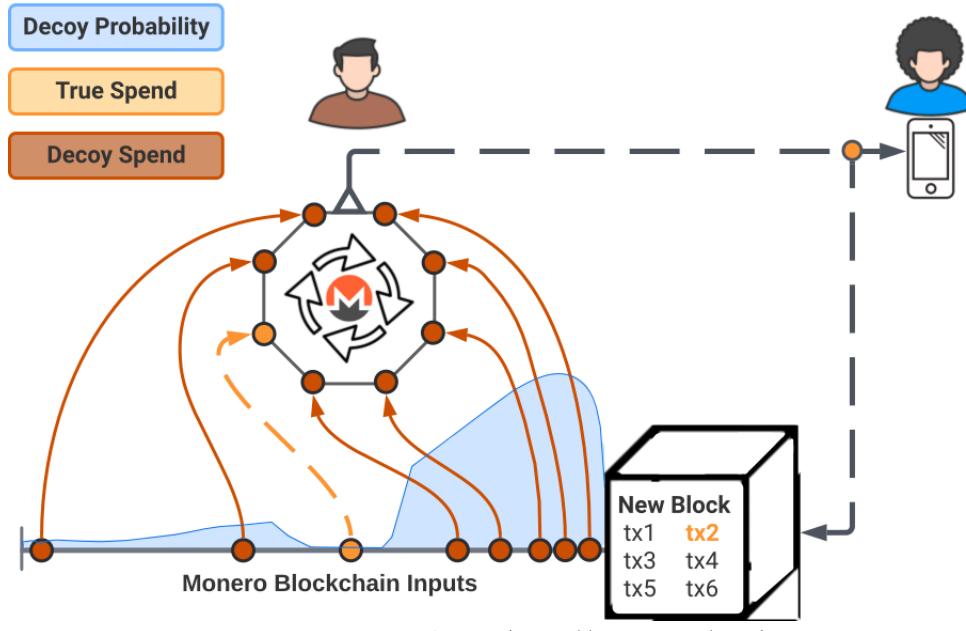
Image designed by ACK-J (<https://github.com/ack-j>)

1425 The height of the blue shape in Figure 2 represents the probability that the decoy selection algorithm selects
 1426 outputs as decoys from certain blockchain time stamps. Younger outputs are more likely to be selected as decoys
 1427 since younger outputs are also more likely to be the real spends. Since there is a high probability that a decoy (red
 1428 circle) could be selected from the same time interval as the real spend (orange), it is difficult for an anti-privacy
 1429 adversary to deduce that the orange circle is the real spend.

²⁴See discussion: <https://libera.monerologs.net/monero-dev/20220526>

1430 Figure 3 shows what would happen if the decoy selection algorithm were severely defective. There is a portion
 1431 of the age distribution that the decoy selection algorithm will never select from. The height of the blue shape is
 1432 zero at that interval. In this hypothetical, the fact that no decoy would be selected in that interval would be known
 1433 by an anti-privacy adversary since the decoy selection algorithm is written into Monero's public open source code,
 1434 at least for Monero wallets that are open source. Therefore, there would be a 100 percent probability that the ring
 1435 member selected in that interval was the real spend. This particular transaction would be 100 percent traceable.
 1436 To repeat, this is the extreme case that does not happen in practice. The current decoy selection algorithm is not
 1437 severely defective like this hypothetical example, but it has significant shortcomings.

Figure 3: Defective Decoy Selection Algorithm



1438 The ideal decoy selection algorithm would select decoys from every point on the age distribution in exact
 1439 proportion to the real spends. For a ring size of 16, this ideal algorithm would provide 15 decoys of "camouflage"
 1440 on top of every real spend across the entire age distribution. In ideal circumstances, such an algorithm would
 1441 not give any hint about the real spend to an anti-privacy adversary. The best an adversary could do is random
 1442 guessing about the real spend, which would only achieve a success rate of $1/16 = 6.25\%$. The "defective" algorithm
 1443 described above provided zero decoys in a certain age interval, which would lead to full traceability of transactions
 1444 that spend outputs that were a certain age. Monero's current decoy selection algorithm lies somewhere between the
 1445 two extremes of ideal and severely defective. The goal of OSPEAD is to move the decoy selection algorithm much
 1446 closer to the ideal shape.

1447 20 The Solution

1448 The goal of OSPEAD is to use a special type of mathematical function called a parametric probability density
 1449 function to match the real spend age distribution as closely as possible. The close match will provide about 15
 1450 decoys for every real spend along nearly every interval of the age distribution. The match cannot be perfect, however.
 1451 Parametric probability density functions can only serve as an approximation of complex real human behavior.

1452 There are only a limited number of decoys available per ring. Therefore, if more decoy "camouflage" is moved
 1453 from section A to section B on the age distribution, users who happen to spend outputs of an age that falls into

section A will experience a reduction in privacy; users spending from section B will experience an increase in privacy. Therefore with the imperfect parametric probability distribution functions there is an unavoidable trade-off. The terms of the trade-off can be made precise mathematically, but the correct choice of trade-off is a judgment call.

Given that the “correct” trade-off is a judgment call, OSPEAD involves a procedure to obtain the best decoy selection algorithm under several different sets of criteria. Once these candidate algorithms are determined, a single best algorithm will be selected though a judgment call. In the next section I describe the proposed sets of criteria: privacy impoverishment, economic welfare, inequality minimization, worst-case-scenario minimization, Maximum Likelihood Estimation, and maximize resistance to a specific attack.

21 Criteria for Best Fit

At first glance constructing $f_D(x)$, the probability density function (PDF) of the decoy selection algorithm, appears to be a standard problem: Extract the real spend age distribution $f_S(x)$ “data” through some method and then fit some distribution to it in order to form $f_D(x)$. Typically in applied statistical work, “fitting some distribution” would be done via nonparametric means or Maximum Likelihood Estimation (MLE) for a parametric approach. Since we are setting aside nonparametric methods in the short term, MLE may see like the way forward.

I would argue instead that what we confront here is not actually a standard problem, despite appearances. An MLE approach, which is the approach taken in [Möser et al., 2018], may make sense if what we are interested in is conducting hypothesis tests on the parameters of parametric distribution families, a typical scientific exercise. In fact, MLE is quite good at that task, assuming certain crucial assumptions are met, since asymptotically it is usually guaranteed to achieve minimum variance among unbiased estimators.

What we are trying to do here with construction of $f_D(x)$ has little to do with conducting hypothesis tests, however. Most statistical methods are interested in separating statistical signal from noise. What we are trying to do is quite the opposite: merge the signal (real spends) and noise (decoys) so that even a determined and sophisticated statistician cannot distinguish them. We are not trying to minimize variance in fitting $f_D(x)$ to $f_S(x)$. Rather, we are trying to minimize risk to user privacy. We are trying to minimize traceability of Monero transactions. This goal requires us to be more creative in our approach.

In general, our goal is to “cover” every age interval with 15 decoys for every real spend that comes from that interval. A user who spends from a point of the age distribution that is not covered by at least 15 decoys experiences a privacy deficit. A user who spends from a point of the age distribution that is covered by more than 15 decoys has a surplus of decoys. To approach the ideal of all users being protected by a single-hop anonymity set of 16, we can move these decoys from the surplus at particular age x values to the x values with deficits. I will build OSPEAD around the idea of minimizing this privacy deficit for all x , defined as $h(x) = \max\{0, f_S(x) - f_D(x)\}$. I also include a symmetric option: $h_{sym}(x) = |f_S(x) - f_D(x)|$ as a possible minimization objective.

To justify my proposed methodology, I will quote extensively from the second edition of *Statistical Inference* by George Casella & Roger L. Berger ([Casella & Berger, 2002]). In the book’s preface, the authors write, “The purpose of this book is to build theoretical statistics (as different from mathematical statistics) from the first principles of probability theory....The book is intended for first-year graduate students majoring in statistics or in a field where a statistics concentration is desirable.” As far as I can tell, it is widely used for that purpose, at least within economics. For example, the first course in the econometrics sequence of MIT’s economics doctoral program uses the book as its primary textbook.²⁵ Therefore, in my view the book is authoritative enough to lean on in justifying my proposed approach.

I will begin by quoting from pages 348–350:

²⁵<https://ocw.mit.edu/courses/economics/14-381-statistical-method-in-economics-fall-2018/syllabus/>
Other economics doctoral programs likely do as well, but we cannot know for sure since many of them keep their syllabi hidden behind student login screens.

1495 *7.3.4 Loss Function Optimality*

1496 Our evaluations of point estimators have been based on their mean squared error performance. Mean
 1497 squared error is a special case of a function called a *loss function*. The study of the performance, and
 1498 the optimality, of estimators evaluated through loss functions is a branch of *decision theory*.

1499 After the data $\mathbf{X} = \mathbf{x}$ are observed, where $\mathbf{X} \sim f(\mathbf{x}|\theta)$, $\theta \in \Theta$, a decision regarding θ is made. The set of
 1500 allowable decisions is the *action space*, denoted by \mathcal{A} . Often in point estimation problems \mathcal{A} is equal to
 1501 Θ , the parameter space, but this will change in other problems (such as hypothesis testing—see Section
 1502 8.3.5).

1503 The loss function in a point estimation problem reflects the fact that if an action a is close to θ , then
 1504 the decision a is reasonable and little loss is incurred. If a is far from θ , then a large loss is incurred.
 1505 The loss function is a nonnegative function that generally increases as the distance between a and θ
 1506 increases. If θ is real-valued, two commonly used loss functions are

$$\text{absolute error loss, } L(\theta, a) = |a - \theta|,$$

1507 and

$$\text{squared error loss, } L(\theta, a) = (a - \theta)^2.$$

1508 Both of these loss functions increase as the distance between θ and a increases, with minimum value
 1509 $L(\theta, \theta) = 0$. That is, the loss is minimum if the action is correct. Squared error loss gives relatively more
 1510 penalty for large discrepancies, and absolute error loss gives relatively more penalty for small discrepancies.
 1511 A variation of squared error loss, one that penalizes overestimation more than underestimation,
 1512 is

$$L(\theta, a) = \begin{cases} (a - \theta)^2 & \text{if } a < \theta \\ 10(a - \theta)^2 & \text{if } a \geq \theta. \end{cases} \quad (17)$$

1513 A loss that penalizes errors in estimation more if θ is near 0 than if $|\theta|$ is large, a relative squared error
 1514 loss, is

$$L(\theta, a) = \frac{(a - \theta)^2}{|\theta| + 1}.$$

1515 Notice that both of these last variations of squared error loss could have been based instead on ab-
 1516 solute error loss. **In general, the experimenter must consider the consequences of various**
 1517 **errors in estimation for different values of θ and specify a loss function that reflects these**
 1518 **consequences.** [my emphasis]

1519 In a loss function or *decision theoretic* analysis, the quality of an estimator is quantified in its *risk*
 1520 function; that is, for an estimator $\delta(\mathbf{x})$ of θ , the risk function, a function of θ , is

$$R(\theta, \delta) = E_\theta L(\theta, \delta(\mathbf{X})). \quad (18)$$

1521 At a given θ , the risk function is the average loss that will be incurred if the estimator $\delta(\mathbf{x})$ is used.

1522 Since the true value of θ is unknown, we would like to use an estimator that has a small value of $R(\theta, \delta)$
 1523 for all values of θ . This would mean that, regardless of the true value of θ , the estimator will have a
 1524 small expected loss. If the qualities of two different estimators, δ_1 and δ_2 , are to be compared, then they
 1525 will be compared by comparing their risk functions, $R(\theta, \delta_1)$ and $R(\theta, \delta_2)$. If $R(\theta, \delta_1) < R(\theta, \delta_2)$ for all
 1526 $\theta \in \Theta$, then δ_1 is the preferred estimator because δ_1 performs better for all θ . More typically, the two
 1527 risk functions will cross. Then the judgment as to which estimator is better may not be so clear-cut.

1528 The risk function for an estimator δ is the expected loss, as defined in (18). For squared error loss, the
 1529 risk function is a familiar quantity, the mean squared error (MSE) that was used in Section 7.3.1. There
 1530 the MSE of an estimator was defined as $\text{MSE}(\theta) = E_\theta(\delta(\mathbf{X}) - \theta)^2$, which is just $E_\theta L(\theta, \delta(\mathbf{X})) = R(\theta, \delta)$
 1531 if $L(\theta, a) = (a - \theta)^2$. As in Chapter 7 we have that, for squared error loss,

$$R(\theta, \delta) = \text{Var}_\theta \delta(\mathbf{X}) + (\text{E}_\theta \delta(\mathbf{X}) - \theta)^2 = \text{Var}_\theta \delta(\mathbf{X}) + (\text{Bias}_\theta \delta(\mathbf{X}))^2.$$

This risk function for squared error loss clearly indicates that a good estimator should have both a small variance and a small bias. A decision theoretic analysis would judge how well an estimator succeeded in simultaneously minimizing these two quantities.

It would be an atypical decision theoretic analysis in which the set \mathcal{D} of allowable estimators was restricted to the set of unbiased estimators, as was done in Section 7.3.2. [my emphasis] Then, minimizing the risk would just be minimizing the variance. A decision theoretic analysis would be more comprehensive in that both the variance and bias are in the risk and will be considered simultaneously. An estimator would be judged good if it had a small, but probably nonzero, bias combined with a small variance.

The text then goes into several specific examples of risk analysis and then briefly explores risk within a Bayesian framework. What does Section 7.3.2 (page 334), referenced above, say?

7.3.2 Best Unbiased Estimators

As noted in the previous section, a comparison of estimators based on MSE considerations may not yield a clear favorite. Indeed, there is no one “best MSE” estimator. Many find this troublesome or annoying, and rather than doing MSE comparisons of candidate estimators, they would rather have a “recommended” one.

The reason that there is no one “best MSE” estimator is that the class of all estimators is too large a class. (For example, the estimator $\hat{\theta} = 17$ cannot be beaten in MSE at $\theta = 17$ but is a terrible estimator otherwise.) One way to make the problem of finding a “best” estimator tractable is to limit the class of estimators. A popular way of restricting the class of estimators, the one we consider in this section, is to consider only unbiased estimators.

If W_1 and W_2 are both unbiased estimators of a parameter θ , that is, $\text{E}_\theta W_1 = \text{E}_\theta W_2 = \theta$, then their mean squared error are equal to their variances, so we should choose the estimator with the smaller variance. If we can find an unbiased estimator with uniformly smallest variance—a best unbiased estimator—then our task is done.

[Casella & Berger, 2002] then go on to give a formal definition of best unbiased estimator — also known as uniform minimum variance unbiased estimator (UMVUE) — in Definition 7.3.7 and then give the Cramér–Rao Lower Bound (CRLB) in Theorem 7.3.9 and Corollary 7.3.10. Corollary 7.3.15 is also useful. The gist of the discussion is that certain unbiased estimators of a parameter θ (the text uses the notation $\tau(\theta)$, with τ being some continuous function of θ , in order to be more general) of a distribution $f(\mathbf{x}|\theta)$ may achieve the lowest variance possible in this setting: the Cramér–Rao Lower Bound (CRLB).

Now I will move on to examining Maximum Likelihood Estimation (MLE) in this entire context. First I will reference the definition of asymptotic variance from [Casella & Berger, 2002], page 471. Let k_n be some normalizing constant.

Definition 10.1.9 For an estimator T_n , suppose that $k_n(T_n - \tau(\theta)) \rightarrow n(0, \sigma^2)$ in distribution. The parameter σ^2 is called the *asymptotic variance* or *variance of the limit distribution* of T_n .

I also need their definition of asymptotic efficiency to make my point:

Definition 10.1.11 A sequence of estimators W_n is asymptotically efficient for a parameter $\tau(\theta)$ if $\sqrt{n}[W_n - \tau(\theta)] \rightarrow n[0, v(\theta)]$ in distribution and

$$v(\theta) = \frac{[\tau'(\theta)]^2}{\text{E}_\theta \left(\left(\frac{\partial}{\partial \theta} \log f(X|\theta) \right)^2 \right)};$$

that is, the asymptotic variance of W_n achieves the Cramér–Rao Lower Bound.

1572 Finally, we come to one of the primary reasons MLE is so popular in applied statistical work:

1573 **Theorem 10.1.12 (Asymptotic efficiency of MLEs)** Let X_1, X_2, \dots , be iid $f(x|\theta)$, let $\hat{\theta}$ denote
1574 the MLE of θ , and let $\tau(\theta)$ be a continuous function of θ . Under the regularity conditions in Miscellanea
1575 10.6.2 on $f(x|\theta)$ and, hence $L(\theta, \mathbf{x})$ [this expression, the likelihood function, is defined in Theorem 10.1.6
1576 as $L(\theta|\mathbf{x}) = \prod_{i=1}^n f(x_i|\theta)$],

$$\sqrt{n} [\tau(\hat{\theta}) - \tau(\theta)] \rightarrow n [0, v(\theta)],$$

1577 where $v(\theta)$ is the Cramér–Rao Lower Bound. That is, $\tau(\hat{\theta})$ is a consistent and asymptotically efficient
1578 estimator of $\tau(\theta)$.

1579 Later, at the top of page 477, [Casella & Berger, 2002] make this comment:

1580 Since the MLE is typically asymptotically efficient, another estimator cannot hope to beat its asymptotic
1581 variance. However, other estimators may have other desirable properties (ease of calculation, **robust-**
1582 **ness to underlying assumptions** [my emphasis]) that make them desirable. In such situations, the
1583 efficiency of the MLE becomes important in calibrating what we are giving up if we use an alternative
1584 estimator.

1585

1586 Let us break down what [Casella & Berger, 2002] are saying in these passages, and how they relate to OSPEAD.
1587 They set up a framework for decision-making based on statistical analysis. Loss functions, decision theory, and risk
1588 functions are the main elements of this framework. They then discuss some of the most common loss functions,
1589 Mean Squared Error (MSE) being one of them. I note in passing that they discuss a “variation” on MSE in equation
1590 (17) that bears some resemblance to my privacy deficit notion, defined as $h(x) = \max \{0, f_S(x) - f_D(x)\}$, in that
1591 it is asymmetric in penalization of under- and over-estimation of some quantity.

1592 Minimizing MSE for estimates of a particular parameter θ is a fair goal for statistical analyses whose purpose
1593 is to conduct traditional hypothesis tests regarding the true value of θ in line with the Popperian falsification
1594 paradigm of science. Lower MSE generally would lead to higher statistical power and therefore greater ability to
1595 avoid Type II error in hypothesis tests. [Casella & Berger, 2002] write, “In general, the experimenter must consider
1596 the consequences of various errors in estimation for different values of θ and specify a loss function that reflects
1597 these consequences.” By minimizing bias and variance explicitly, MSE — as a loss function — performs well in
1598 avoiding the Type I and Type II error consequences in null hypothesis testing and therefore MSE often makes sense
1599 to use as a loss function in scientific hypothesis testing.

1600 [Casella & Berger, 2002] go further and observe, “It would be an atypical decision theoretic analysis in which
1601 the set \mathcal{D} of allowable estimators was restricted to the set of unbiased estimators, as was done in Section 7.3.2.”
1602 As I highlighted, MLE typically has lowest variance, asymptotically, among estimators with zero bias. Therefore,
1603 restricting ourselves to MLE in tackling the problem before us — construction of the PDF for the decoy selection
1604 algorithm — would likely seem ill-advised to [Casella & Berger, 2002], since the class of unbiased estimators might
1605 not be suitable.

1606 Note also that maximum likelihood estimators are point estimators for the particular θ parameters of PDFs. In
1607 theory, the estimated $\hat{\theta}$ will yield the corresponding theoretical distributions that fit the target empirical distributions
1608 well, although only under the assumption that the empirical distribution being fitted exactly equals the chosen
1609 theoretical PDF. There is no guarantee that variance — or bias for that matter — for the distribution itself will
1610 be small when using MLE to fit the wrong theoretical PDF to an empirical distribution. Given the fact that we
1611 are interested in the whole distributions themselves, i.e. $f_S(x)$ and $f_D(x)$, and not just parameters of parametric

1613 distributions, the justification for using MLE in this setting is even weaker. In just a moment I will convert the
 1614 problem of fitting distributions into a parametric one, so that the true problem is better illustrated.

1615

1616

1617 Until now I have dealt with $f_S(x)$ and $f_D(x)$ as if they were probability density functions, i.e. as if the domain
 1618 of the functions were continuous. However, in reality they are probability mass functions, i.e. the domain of the
 1619 functions are discrete. The real spend age distributions are only meaningful in terms of the discrete blocks on the
 1620 blockchain, since miners can arrange valid transactions within a block in any order they choose. We can leverage
 1621 this fact for the following analysis.

1622 Define a set of parameters θ as follows:

$$\theta = \{\theta_1, \theta_2, \dots, \theta_N\} = \{f_S(1), f_S(2), \dots, f_S(N)\}$$

1623 So θ_1 is the value of $f_S(x)$ when $x = 1$, i.e. the first-available block that an output can be spent in. θ_N refers to
 1624 the first block where RingCT outputs appeared. My definition here is not completely rigorous, since the blockchain
 1625 is lengthening constantly, and therefore the number of elements of θ is increasing by the hour.

1626 What we wish to estimate is θ — every element of it. Or, more specifically, we wish to construct some $f_D(x)$
 1627 such that $f_D(x_i)$ is “close” — in some sense to be defined shortly — to each θ_i , for every i . With a nonparametric
 1628 approach, we may be able to tackle each element θ_i individually, more or less. With a parametric approach, we
 1629 cannot hope to do so. Therefore, we must establish some overall metric, or metrics, of “success”. In other words,
 1630 we must define and justify a set \mathcal{L} of loss functions. This is our first task. Our second task is to define a set \mathcal{D} of
 1631 allowable estimators.

1632 In equation (18), [Casella & Berger, 2002] defined the risk function as $R(\theta, \delta) = E_\theta L(\theta, \delta(\mathbf{X}))$. For the time
 1633 being, we will assume that θ is deterministic and therefore the risk function equals the loss function. The θ is
 1634 actually stochastic — it changes over time as user spending patterns change over time — which will be dealt with
 1635 in **Section 24 Dynamic Risk and Forecasting**.

1636 There are some similarities between the Differential Privacy framework and Monero’s ring signature privacy
 1637 model. Differential Privacy seeks a balance between the desire to perform statistical analysis on private data and
 1638 the need to protect individuals from discovery of their private information. Monero seeks to maximize privacy
 1639 of users, but is constrained by reasonable limits on ring size. I looked through the Differential Privacy literature
 1640 for some criteria for fitting a decoy selection algorithm. In general there are utility-based criteria for resolving
 1641 the balance between statistical analysis and user privacy for each user, but there is not yet a way to resolve our
 1642 problem([Hsu et al., 2014]). Our problem is different in that a particular ring size is fixed, and then trade-offs
 1643 between users need to be determined.

1644 In “Research Roadmap for an Overhaul of Monero’s Mixin Selection Algorithm”, which I submitted to Monero’s
 1645 Vulnerability Response Process in 2021, I wrote:

1646 How do we construct $f_M(x)$ [the decoy selection algorithm] for best overall privacy if we restrict ourselves
 1647 to parametric distributions, and therefore cannot achieve $f_S(x) \approx f_M(x)$ for all values of x ? Well, it
 1648 depend [sic] on how we define “best”. Below are six approaches. I have the mathematical definitions of
 1649 these worked out in my head, but they are not written here:

1650

- 1651 1) Privacy impoverishment
- 1652 2) Economic welfare
- 1653 3) Inequality minimization
- 1654 4) Worst-case-scenario minimization
- 1655 5) Maximum Likelihood Estimation

1656 6) Maximize resistance to a specific attack

1657
1658 Now that a set of objective functions have been defined, one could imagine optimizing these objective
1659 functions in a numerical optimization procedure where each parametric distribution family with support
1660 of $[0, \infty)$ is permuted over 1-6 above.

1661 These objective functions are the loss functions in the framework of [Casella & Berger, 2002]. Many of them are
1662 forms of minimum divergence estimators ([Maji et al., 2019]). Just two more bits of housekeeping before I can
1663 mathematically define these 1-6 loss functions. First, to be consistent with the notation of [Casella & Berger, 2002],
1664 define $f_D(x_i) \equiv a_i$.

1665 Second, we must think about how to give weight to each θ_i . Do we give equal weight to every θ_i ? That would
1666 mean that outputs years old would be given the same importance as outputs just a few minutes or hours old.
1667 Weighting the loss function by the value of θ_i seems more reasonable since that would, in effect, give each spent
1668 output equal importance. However, we are trying to protect the privacy of people, not outputs.

1669 Say that User A re-spends outputs frequently, and so has real spends that are in the thick portion of $f_S(x)$.
1670 Say that User B re-spends outputs infrequently, and therefore has real spends that are in the thin portion of $f_S(x)$.
1671 Furthermore, User A would be able to generate more transactions than User B simply because of rapidly re-spending
1672 outputs. Therefore, only weighting the loss function by the value of θ_i would (maybe unfairly) give more importance
1673 to User A than User B, just due to the way that the mathematics work out. The way forward is not exactly clear.
1674 I think it could make sense to try several intermediate weighting schemes.

1675 First, my recommendation is to only give weight to elements of $\boldsymbol{\theta}$ when there is a privacy deficit, i.e. when
1676 $h(x) = \max \{0, f_S(x) - f_D(x)\}$ is nonzero. Later we will use $\mathbf{1}\{x\}$, the indicator function, for this part of the
1677 weighting scheme, or the $\min\{x\}$ operator, depending on the context.²⁶ For comparison purposes, I will also
1678 include the symmetric counterpart $h_{sym}(x) = |f_S(x) - f_D(x)|$ that seeks to avoid privacy surpluses as much as it
1679 seeks to avoid privacy deficits. To allow for multiple intermediate weighting schemes, define this weight function:

$$w(\theta_i, \lambda) = \lambda\theta_i + (1 - \lambda)\frac{1}{N} \quad (19)$$

1680 Thus, when $\lambda = 1$, $w(\theta_i, \lambda)$ is fully weighted by the value of θ_i . When $\lambda = 0$, $w(\theta_i, \lambda)$ gives equal weight to each
1681 θ_i . Tentatively, let $\boldsymbol{\lambda} = \{0, 0.5, 0.9, 0.95, 0.99, 0.999, 0.9999, 1\}$.

1682 Now we are ready to define the set \mathcal{L} of loss functions.

1683 21.1 Privacy impoverishment

1684 To me, the privacy deficit formulation is somewhat reminiscent of a poverty line. Below a certain defined threshold
1685 individuals are considered to be impoverished. In the case of a standard poverty indicator, there is some poverty
1686 line z defined by a researcher or government entity. In the case of Monero user privacy, the corresponding “poverty
1687 line” would be $f_S(x)$, which is different for every x .

1688 The Foster–Greer–Thorbecke (FGT) indices are a family of widely-used poverty indicators. They are defined as

$$FGT_\alpha = \frac{1}{N} \sum_{i=1}^H \left(\frac{z - y_i}{z} \right)^\alpha$$

1689 where N is the number of people in the population under study, H is the number of people within that population
1690 who are below the poverty line, z is the poverty line, y_i is the income (or consumption) of each individual i who is
1691 under the poverty line, and α is a parameter that controls weighting.

²⁶https://en.wikipedia.org/wiki/Indicator_function

If α is high, then people far below the poverty line are given much more weight than people just barely below the poverty line. When $\alpha = 0$, FGT_0 is simply the poverty headcount. When $\alpha = 1$, FGT_1 is the poverty gap index. When $\alpha = 2$, FGT_2 weights each person's poverty gap by its square, and for $\alpha = 3$, by its cube, and so forth.

A loss function analogue can be defined:

$$L_{FGT_\alpha}(\boldsymbol{\theta}, \mathbf{a}, \lambda) = \frac{1}{N} \sum_{i=1}^N \mathbf{1}\{a_i < \theta_i\} w(\theta_i, \lambda) \left(\frac{\theta_i - a_i}{\theta_i} \right)^\alpha \quad (20)$$

The $\mathbf{1}\{a_i < \theta_i\}$ indicator function ensures that a_i 's are only counted when there is some privacy deficit. The symmetric counterpart $L_{FGT_{\alpha, sym}}(\boldsymbol{\theta}, \mathbf{a}, \lambda)$ would omit the $\mathbf{1}\{a_i < \theta_i\}$ term and replace $\left(\frac{\theta_i - a_i}{\theta_i} \right)^\alpha$ with $\left(\left| \frac{\theta_i - a_i}{\theta_i} \right| \right)^\alpha$. Using $\alpha = 0$ probably doesn't make sense for our purposes. Using $\alpha = \{1, 2, 3\}$ is reasonable.

21.2 Economic welfare

A full discussion of the meaning and theory of economic welfare is outside of the scope of this document. The basic idea is that people are theorized to have preferences that can be approximated via a mathematical function. These approximating functions are called utility functions. When people's preferences are satisfied, they obtain utility. When the preferences are satisfied to a higher degree, they obtain higher utility. In essence, when an individual's utility is higher, their happiness is higher. Within the framework of neoclassical economics, individuals strive to maximize their utility, subject to the constraints they encounter in the world. The analysis of aggregate population-level utility is the realm of welfare economics.

We may assume that Monero users prefer to have privacy and therefore some "privacy utility function" in this specific context could be defined. In practice, the loss function derived from this type of analysis will look similar to the $L_{FGT_\alpha}(\boldsymbol{\theta}, \mathbf{a}, \lambda)$ privacy impoverishment loss function, but with a different interpretation.

There are dozens of utility functions to choose from. One of the more appropriate ones in the present context is the Constant Relative Risk Aversion (CRRA) utility function:

$$u_{CRRA_\eta}(y) = \begin{cases} \frac{y^{1-\eta}}{1-\eta} & \eta \geq 0, \eta \neq 1 \\ \ln(y) & \eta = 1 \end{cases} \quad (21)$$

where η is the coefficient of constant relative risk aversion. The CRRA utility function makes sense to use in this context because (1) At a basic level, it takes a single argument, unlike other classes of utility functions that deal with multiple goods and services; (2) It explicitly deals with risk; (3) It has an adjustable parameter η that we can use to explore the sensitivity of the results; (4) For $\eta \geq 1$, the CRRA utility function approaches $-\infty$ as y approaches zero.

The (4) characteristic serves as an important advantage compared to the privacy impoverishment framework since in theory the $L_{FGT_\alpha}(\boldsymbol{\theta}, \mathbf{a}, \lambda)$ loss function would allow $f_D(x)$ to have zero mass at some x values — and therefore any ring members having an age corresponding to those x values would be clearly identifiable as real spends. In addition, a numerical optimization process that used the CRRA utility function as its basis would avoid at all costs $f_D(x) = 0$ for any and all x values as long as η is chosen so that $\eta \geq 1$.

In a welfare economics framework, generally some discussion of the Pareto weights and Pareto efficiency would be appropriate. However, that type of discussion has little practical effect on the analysis here and it would generally only be of interest to economists, so I will elide it here.

Now we are ready to define an economic welfare loss function:

$$L_{Welfare_\eta}(\boldsymbol{\theta}, \mathbf{a}, \lambda) = (-1) \frac{1}{N} \sum_{i=1}^N w(\theta_i, \lambda) \left(u_{CRRA_\eta} \left(\min \left\{ \frac{a_i}{\theta_i}, 1 \right\} \right) \right) \quad (22)$$

The (-1) scalar is to make this a *loss* function to be minimized. As is typical, the CRRA utility function increases with its argument, and privacy increases as a_i becomes closer to θ_i , i.e. as $f_D(i)$ becomes closer to $f_S(i)$. We use the $\min \{x\}$ operator here. When $a_i < \theta_i$, the utility function is applied to $\frac{a_i}{\theta_i}$. When $a_i \geq \theta_i$, the utility function is applied to one, as if $a_i = \theta_i$, since users spending outputs from this block i suffer no privacy deficit as I have defined it. The symmetric counterpart $L_{Welfare_{\eta, sym}}(\boldsymbol{\theta}, \mathbf{a}, \lambda)$ would replace $\min \left\{ \frac{a_i}{\theta_i}, 1 \right\}$ with $\frac{a_i}{\theta_i}$.

The choice for η is somewhat open-ended. Perhaps five values should be tested, to get some sense of the sensitivity of the results to the choice of η . One of the values should be $\eta = 1$ so that log utility is used. It is less clear what the remaining four values should be. Some experimentally-determined values for η for the utility of income are available in the literature, but how individuals' utility functions for income and privacy relate to one another is unclear. Tentatively, let us set $\boldsymbol{\eta}$ to $\boldsymbol{\eta} = \{0.5, 1, 2, 5\}$.

21.3 Inequality minimization

Another possible loss function is an inequality metric. The privacy impoverishment and economic welfare frameworks dealt with the absolute privacy of each user. We may also want to consider a framework that explicitly compares users to each other, in terms of the privacy that a $f_D(x)$ provides. A loss function that attempts to minimize inequality of privacy among users may be appropriate.

One of the most widely-used inequality metrics is the Gini coefficient (or Gini index). The Gini coefficient has many attractive theoretical properties that I will not recite here. Its main drawback is that its value is not very interpretable for laypeople. (The most intuitive explanation involves the Lorenz curve.) Given that we are already far into the realm of difficulty with interpretation of these loss functions, the low interpretability of the Gini coefficient should not stop us from using it.

One formulation of the Gini coefficient is:

$$G(\mathbf{x}) = \frac{\sum_{i=1}^n \sum_{j=1}^n |x_i - x_j|}{2 \sum_{i=1}^n \sum_{j=1}^n x_j} = \frac{\sum_{i=1}^n \sum_{j=1}^n |x_i - x_j|}{2n^2 \bar{x}}$$

This formulation is somewhat easy to interpret, although it is computationally expensive as it requires about $(n^2 - n)/2$ arithmetic operations. It does not contemplate weights, so it must be modified. [Creedy, 2015] provides some guidance.

First, the formula suggested by Creedy requires that the weights sum up to n , i.e. $\sum_{i=1}^n w_i = n$. The weights that we will use, $w(\theta_i, \lambda)$ could easily be normalized to ensure $\sum_{i=1}^n w(\theta_i, \lambda) = n$.

In his equation (17), Creedy suggests the following weighted formula:

$$G(\mathbf{x}) = \frac{\sum_{i=1}^n \sum_{j=i+1}^n w_i w_j |x_i - x_j|}{\sum_{i=1}^n w_i \sum_{i=1}^n w_i x_i}$$

If we let x_i be $\frac{a_i}{\theta_i}$ as with the economic welfare loss function, the corresponding loss function for our problem would then be:

$$L_{Gini}(\boldsymbol{\theta}, \mathbf{a}, \lambda) = \frac{\sum_{i=1}^n \sum_{j=i+1}^n w(\theta_i, \lambda) \cdot w(\theta_j, a_j, \lambda) \cdot \left| \min \left\{ \frac{a_i}{\theta_i}, 1 \right\} - \min \left\{ \frac{a_j}{\theta_j}, 1 \right\} \right|}{\sum_{i=1}^n w(\theta_i, \lambda) \sum_{i=1}^n w(\theta_i, \lambda) \cdot \min \left\{ \frac{a_i}{\theta_i}, 1 \right\}} \quad (23)$$

$$\text{subject to } \sum_{i=1}^n w(\theta_i, \lambda) = n$$

1756 Note that the formula given in equation (17) in [Creedy, 2015] is just one possible way to deal with weighting
1757 to compute a Gini coefficient. The reason that there are multiple ways to do weighting is that weighting attempts
1758 to, in essence, interpolate parts of the Lorenz curve. How to do that interpolation is up for debate. In his equation
1759 (12), [Creedy, 2015] also gives a weighted Gini formula that is less computationally expensive. However, it would
1760 require conversion of $w(\theta_i, \lambda)$ to integers, which is feasible but adds another layer of complication. The symmetric
1761 counterpart $L_{Gini,sym}(\boldsymbol{\theta}, \mathbf{a}, \lambda)$ would replace $\min \left\{ \frac{a_i}{\theta_i}, 1 \right\}$ with $\frac{a_i}{\theta_i}$.

1762 21.4 Worst-case-scenario minimization

1763 Worst-case-scenario minimization is inspired by Appendix F: Minimum Untraceability in [Möser et al., 2018]. They
1764 define Ge_{min} , the minimum possible guessing entropy of the untraceability of a transaction input:

$$Ge_{min} = \frac{\frac{1}{2}m(m+1)}{\frac{r_{max}}{r_{min}} + m} \quad (24)$$

1765 where m is the number of ring members and

$$r_{max} = \max_{\forall x} \left(\frac{f_S(x)}{f_D(x)} \right), r_{min} = \min_{\forall x} \left(\frac{f_S(x)}{f_D(x)} \right)$$

1766 Since m is a constant, and optimization procedures are insensitive to constants, for the purposes of constructing
1767 a loss function, Ge_{min} can be simplified to

$$Ge_{min(optimizer)} = \frac{r_{max}}{r_{min}} \quad (25)$$

1768 The logic to construct $Ge_{min(optimizer)}$ is as follows:

$$\begin{aligned} Ge_{min(optimizer)} &= \frac{\frac{1}{2}m(m+1)}{\frac{r_{max}}{r_{min}} + m} \\ &\propto \frac{1}{\frac{r_{max}}{r_{min}} + m} && \text{scaling by } \left(\frac{1}{2}m(m+1) \right)^{-1} \text{ is a monotonic transformation} \\ &\cong \frac{r_{max}}{r_{min}} + m && g(x) = \frac{1}{x} \text{ is a strictly decreasing monotonic transformation} \\ &\cong \frac{r_{max}}{r_{min}} && \text{subtracting } m \text{ is a monotonic transformation} \end{aligned}$$

1769 (I am using the \cong symbol loosely here.) Therefore, a $f_D(x)$ that minimizes (25) will also maximize (24)

1770

1771 Putting these ideas into the loss function notation that I have been using, i.e.

$$r_{max} = \max_{\forall i} \left(\frac{\theta_i}{a_i} \right), r_{min} = \min_{\forall i} \left(\frac{\theta_i}{a_i} \right)$$

1772 gives us

$$L_{Worst\ case}(\boldsymbol{\theta}, \mathbf{a}) = \frac{\max_{i \in \{1, \dots, N\}} \left(\frac{\theta_i}{a_i} \right)}{\min_{i \in \{1, \dots, N\}} \left(\frac{\theta_i}{a_i} \right)} \quad (26)$$

1773 It is not clear to me how weighting by $w(\theta_i, \lambda)$ might work here since the numerator and denominator are both
1774 single numbers. For the time being I will leave it out.

1775 Although I include $L_{Worst\ case}$ in the proposed set \mathcal{L} of loss functions, I am skeptical of its usefulness. It is
1776 not clear to me that we should care about the minimum θ_i and a_i among literally hundreds of thousands of them.
1777 Maybe the tail should not wag the dog. However, it could serve as a useful comparison to other approaches, so I
1778 include it in \mathcal{L} . Another issue is that I suspect that $L_{Worst\ case}$ will not be a well-behaved function for the purposes
1779 of numerical optimization. Given the min and max operators and how they are used here, I foresee that $L_{Worst\ case}$
1780 as a function will have some significant discontinuities with respect to i .

1781 21.5 Maximum Likelihood Estimation

1782 As I have stated, I do not favor a MLE approach here. However, including it in the set of loss functions \mathcal{L} might
1783 serve as a useful comparison.

1784 Let $\boldsymbol{\beta} = \{\beta_1, \dots, \beta_k\}$ be the set of parameters of some parametric probability density function $g(x|\boldsymbol{\beta})$. Then the
1785 likelihood function for some sample \mathbf{x} is

$$L(\boldsymbol{\beta}|\mathbf{x}) = \prod_{i=1}^n g(x_i|\boldsymbol{\beta})$$

1786 To put MLE in a loss function framework, convert it into a minimization problem and take the log to ease the
1787 computational burden:

$$L_{MLE}(\boldsymbol{\beta}|\mathbf{x}) = (-1) \cdot \sum_{i=1}^n \log g(x_i|\boldsymbol{\beta}) \quad (27)$$

1788 21.6 Maximize Resistance to an Attack

1789 Define the potency of an attack on the untraceability of Monero transactions for a specified $f_D(x)$ as $\mathcal{P}(f_D(x))$,
1790 the unconditional probability of correctly guessing the real spend. Then the corresponding loss function can be
1791 defined as

$$L_{Attack}(f_D(x)) = \mathcal{P}(f_D(x)) \quad (28)$$

1792

1793

1794 Now that the set \mathcal{L} of loss functions has been defined, the set \mathcal{D} of allowable estimators will be defined. Continuing
1795 with the θ_i notation, the allowable estimators will be selected from the set of parametric probability density functions
1796 (PDF) $f(x|\boldsymbol{\beta})$ such that each θ_i shall be estimated by the image of i under $f(x|\boldsymbol{\beta})$. In other words, $\hat{\theta}_i = f(i|\boldsymbol{\beta})$
1797 for all i .

1798 Strictly speaking, the estimator here is the minimizer of a specified loss function in the set \mathcal{L} where $a_i = f(i)$
1799 for all i for some specified parametric PDF $f(x|\boldsymbol{\beta})$. The set of all PDFs (and probability mass functions (PMFs),
1800 for that matter) can be all PDFs whose support is the set $[0, \infty)$. The Wikipedia page on probability distributions

lists over 40 distributions that have such a support.²⁷ So let us say that the target number of elements of the set of PDFs used in estimation is 40. Note that many distributions are special cases of more general distributions and therefore the actual number of distributions to fit may be lower than 40 in practice. Define the set \mathcal{F} of these PDFs. Let \mathcal{B} be the set of parameters of each $f(x|\beta)$ under consideration.

21.7 Formalized Optimization Criteria

First, assume we have a good estimate of $f_S(x)$. Then OSPEAD is the set of procedures that perform the following minimizations:

$$\begin{aligned}
 \forall L \in \mathcal{L}, & \quad \text{Loss functions} \\
 \forall f(x|\beta) \in \mathcal{F}, & \quad \text{Parametric distributions} \\
 \forall \lambda \in \boldsymbol{\lambda}, & \quad \text{Weight parameters for } w(\theta_i, \lambda) \\
 \forall \alpha \in \boldsymbol{\alpha}, & \quad \text{Exponents for FGT poverty indicator} \\
 \forall \eta \in \boldsymbol{\eta}, & \quad \text{Parameters for CRRA utility} \\
 \min_{\beta \in \mathcal{B}} L(f(x|\beta)) & \quad \text{Numerical optimization problems}
 \end{aligned} \tag{29}$$

where

$$\begin{aligned}
 \mathcal{L} &= \{L_{FGT_\alpha}, L_{Welfare_\eta}, L_{Gini}, L_{Worst\ case}, L_{MLE}, L_{Attack}\} && \text{Equations (20), (22), (23), (26), (27), (28)} \\
 \mathcal{F} &= \{\text{Roughly 40 PDF/PMFs that have support } [0, \infty)\} \\
 \boldsymbol{\lambda} &= \{0, 0.5, 0.9, 0.95, 0.99, 0.999, 0.9999, 1\} \\
 \boldsymbol{\alpha} &= \{1, 2, 3\} \\
 \boldsymbol{\eta} &= \{0.5, 1, 5, 10, 50\} \\
 \mathcal{B} &= \{\text{Sets of parameters associated with each element of } \mathcal{F}\}
 \end{aligned}$$

For each loss function there will likely be a unique minimizer $f^*(x|\beta)$. However, it is unlikely that every loss function will have the same minimizer. Therefore, determining which candidate $f(x|\beta)$ is “best” will be the result of a judgment call, taking into account the totality of evidence and theory.

Why stop with single PDFs when we could also examine mixture distributions? It would be useful to test some mixture distributions composed of two or three different parametric PDFs. The main difficulty is the combinatorics: The number of two-PDF mixture distributions would be $\binom{40}{2} = 780$. Given the loss functions and their variations, that would amount to about 58,500 numerical minimization procedures. There are already practical difficulties with carrying out the minimizations with no mixture distributions. Starting values have to be defined and robust minimization algorithms must be chosen and checked for problems. Even with no mixture distributions, the number of minimizations that I set out to perform, as stated in (29), is about 3,000.²⁸

Therefore, it could make sense to choose a small set of the most promising or complementary PDFs to do a second round of OSPEAD with mixture distributions. Another possible enhancement could be incorporating some periodic component in a mixture distribution to account for users’ sleep-wake cycle.

²⁷[https://en.wikipedia.org/wiki/List_of_probability_distributions#Supported_on_semi-infinite_intervals,_usually_\[0,%E2%88%9E](https://en.wikipedia.org/wiki/List_of_probability_distributions#Supported_on_semi-infinite_intervals,_usually_[0,%E2%88%9E)

²⁸Let $|\mathbf{S}|$ denote the cardinality of set \mathbf{S} . Note that $|\boldsymbol{\lambda}| = 8$, $|\boldsymbol{\alpha}| = 3$, $|\boldsymbol{\eta}| = 5$, and $|\mathcal{F}| = 40$. Taking into account the different flavors of the loss functions, $|L_{FGT_\alpha}(\theta, \alpha, \lambda)| = |\boldsymbol{\lambda}| \cdot |\boldsymbol{\alpha}| = 24$, $|L_{Welfare_\eta}(\theta, \alpha, \lambda)| = |\boldsymbol{\lambda}| \cdot |\boldsymbol{\eta}| = 40$, $|L_{Gini}(\theta, \alpha, \lambda)| = |\boldsymbol{\lambda}| = 8$, $|L_{Worst\ case}(\theta, \alpha)| = 1$, $|L_{MLE}(\theta|\mathbf{x})| = 1$, and $|L_{RRRA}(f_D(x))| = 1$. Therefore, $|\mathcal{L}| = 75$ and hence $|\mathcal{L}| \cdot |\mathcal{F}| = 3,000$.

1821 22 Dry Run with Old Moser et al. (2018) Data

1822 To see how the minimizers in (29) would work in practice, we can perform a dry run on partially de-anonymized
 1823 Monero spend age data provided by [Möser et al., 2018] from before February 2017. The final version of OSPEAD
 1824 will use data from about September 2021 to October 2022, but the dry run is useful as a demonstration.

1825 In the dry run I use a subset of the loss function and distribution functions described in (29). For the loss
 1826 functions I select L_{FGT_α} with $\alpha = \{1, 2\}$, $L_{Welfare_\eta}$ with $\eta = \{0.5, 1\}$, and L_{MLE} . For the distribution function
 1827 I choose the Log-gamma, Noncentral F, Right-Pareto Log-normal, Generalized Extreme Value, and Generalized
 1828 Hyperbolic distributions. I also include mixtures of Log-gamma with F, Generalized Extreme Value, and a Laplace
 1829 Periodic distribution.

1830 **Log-gamma:** A two-parameter distribution used to fit the spend age data in [Möser et al., 2018]. The current decoy
 1831 selection algorithm is based on a particular form of the Log-gamma distribution. The Gamma distribution
 1832 includes the exponential, Erlang, and chi-square distributions as special cases.

1833 **Noncentral F:** A three-parameter distribution used often in hypothesis testing.

1834 **Right-Pareto Log-normal:** A three-parameter distribution that is a special case of the Double Pareto-lognormal
 1835 distribution, which “arises as that of the state of a geometric Brownian motion (GBM), with lognormally
 1836 distributed initial state, after an exponentially distributed length of time” [Reed & Jorgensen, 2004].

1837 **Generalized Extreme Value:** A three-parameter distribution that includes the Gumbel, Fréchet and Weibull
 1838 distributions as special cases.

1839 **Generalized Hyperbolic** A six-parameter distribution that includes the Normal Inverse Gaussian, Variance
 1840 Gamma, and Generalized Hyperbolic Student-t distributions as special cases.

Table 4: Performance of Dry Run with Old Moser et al. (2018) Data

Loss function	L_FGT	L_FGT	L_Welfare	L_Welfare	L_MLE
Loss function parameter	1	2	0.5	1	
Log-gamma	0.1138	0.0574	-1.8542	0.1955	7.65e+07
F	0.1095	0.0462	-1.8681	0.1635	7.67e+07
Right-Pareto Log-normal	0.1073	0.0449	-1.8703	0.1604	7.66e+07
Generalized Extreme Value	0.1169	0.0503	-1.8608	0.1766	7.76e+07
Generalized Hyperbolic	0.1100	0.0455	-1.8684	0.1632	7.62e+07
Log-gamma + F mix	0.1081	0.0402	-1.8721	0.1622	
Log-gamma + GEV mix	0.1095	0.0460	-1.8704	0.1547	
Log-gamma + Laplace Periodic	0.112	0.0572	-1.86	0.195	

Note: Values should be compared down columns. Lower values (darker green) indicate better performance. MLE value is Akaike Information Criterion (AIC).

1841 Table 4 contains the values of the minimums for the various loss functions and parametric distributions. These
 1842 numbers are useful for demonstration purposes, but the final numbers will look very different. We can see that

1843 the Log-gamma distribution tends to compare poorly to the other distributions. The Right-Pareto Log-normal
1844 distribution tends to perform better than other pure distributions. The Log-gamma has poor performance and the
1845 F has mediocre performance, but when they are combined in a mixture distribution their performance is consistently
1846 good.

1847 Table 5 contains the optimized parameter values for the various loss functions and parametric distributions. An
1848 important fact to recognize is that a particular set of parameter values that gives good performance according to
1849 one loss function may give poor performance for a different loss function. In these exercises we are not comparing a
1850 fixed set of parameter values across various loss functions. Rather, we are allowing each distribution's parameters
1851 to be bent to minimize each loss function value.

1852 Following Table 5 is a series of plots showing the fitted distributions. The vertical black lines represent the
1853 empirical probability mass function. The vertical green lines extending below the horizontal axis represents the
1854 empirical cumulative distribution function. 76 percent of the mass of the M oser et al. (2018) data is less than
1855 10,000 blocks old.

1856 The first five plots show the fitted PDFs themselves. To compare the fitted distribution more closely with
1857 the empirical data, the second set of five plots shows the ratio of the fitted distributions $f_D(x)$ to the empirical
1858 probability mass function $f_S(x)$.

1859 The code for producing these tables and plots is at <https://github.com/Rucknium/OSPEAD>

Table 5: Optimized Parameter Values of Dry Run with Old Moser et al. (2018) Data

Distribution	Loss fn	Loss fn param	param_1	param_2	param_3
F	L_FGT	1	0.0258	0.471	7.26
F	L_FGT	2	0.0338	0.294	10.3
F	L_Welfare	0.5	0.0289	0.429	8.1
F	L_Welfare	1	0.0333	0.383	9.27
F	L_MLE	0	0.0473	0.685	11.5
Generalized Extreme Value	L_FGT	1	106	421	2.63
Generalized Extreme Value	L_FGT	2	-72.7	3.14e-29	5.74
Generalized Extreme Value	L_Welfare	0.5	-95.9	1.92e-18	3.88
Generalized Extreme Value	L_Welfare	1	-86.6	6.28e-15	4.34
Generalized Extreme Value	L_MLE	0	99.7	246	2.48
Log-gamma	L_FGT	1	6.48	0.894	
Log-gamma	L_FGT	2	5.15	0.639	
Log-gamma	L_Welfare	0.5	6.25	0.852	
Log-gamma	L_Welfare	1	5.79	0.76	
Log-gamma	L_MLE	0	6.62	0.912	
Right-Pareto Log-normal	L_FGT	1	0.235	4.12	1.38
Right-Pareto Log-normal	L_FGT	2	0.133	3.79	1.21
Right-Pareto Log-normal	L_Welfare	0.5	0.209	4.03	1.32
Right-Pareto Log-normal	L_Welfare	1	0.18	3.92	1.25
Right-Pareto Log-normal	L_MLE	0	0.444	4.99	1.83

F Distribution: param_1 is first degree of freedom parameter; param_2 is second degree of freedom parameter; param_3 is non-centrality parameter.

Generalized Extreme Value Distribution: param_1 is location parameter; param_2 is scale parameter; param_3 is shape parameter.

Log-gamma Distribution: param_1 is shape parameter; param_2 is rate parameter.

Right-Pareto Log-normal Distribution: param_1 is shape parameter; param_2 is mean parameter; param_3 is variance parameter.

Mixture distributions and Generalized Hyperbolic Distribution are omitted from this table.

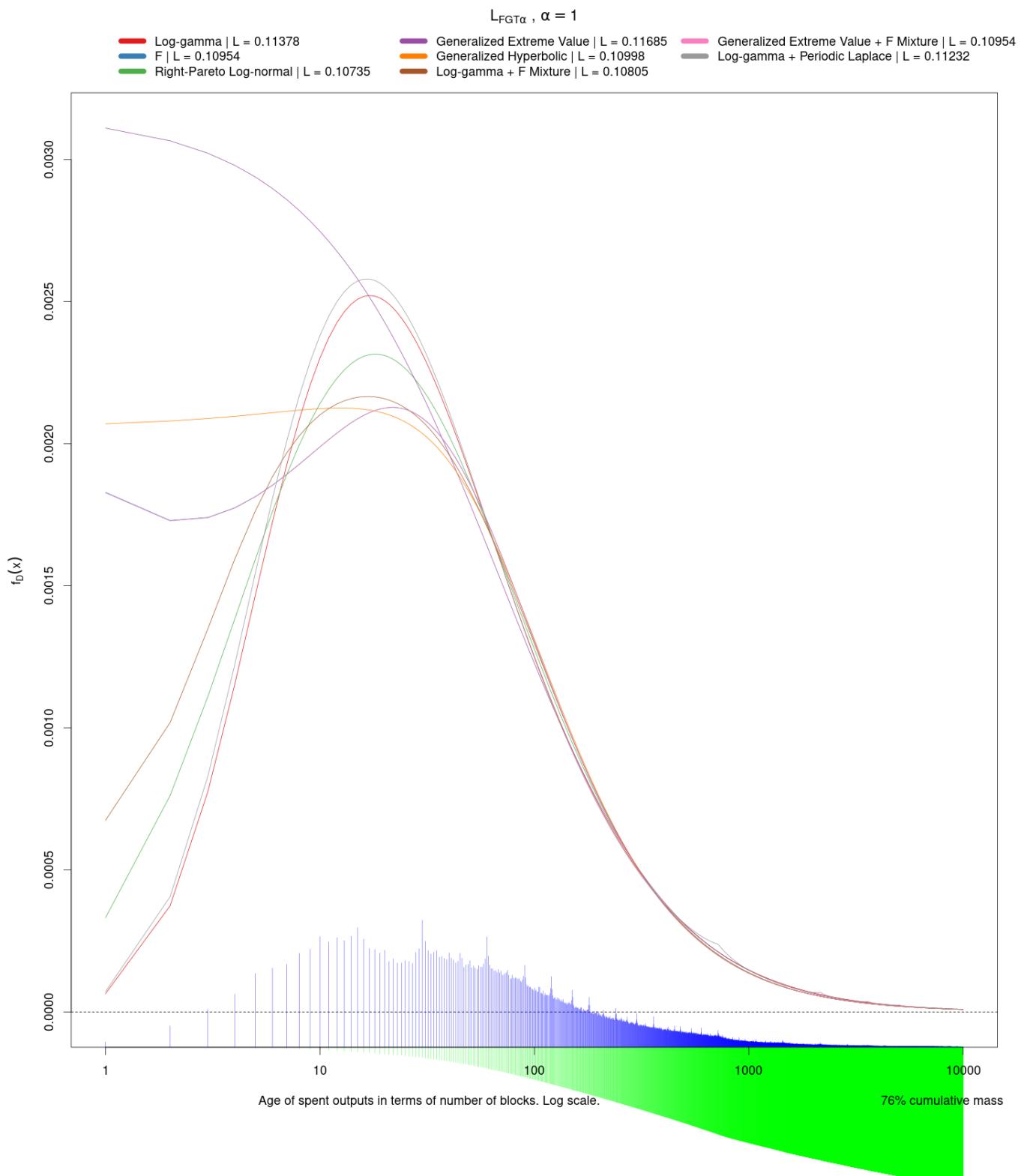
Figure 4: Optimized $f_D(x)$ for loss function L_{FGT_α} , $\alpha = 1$ 

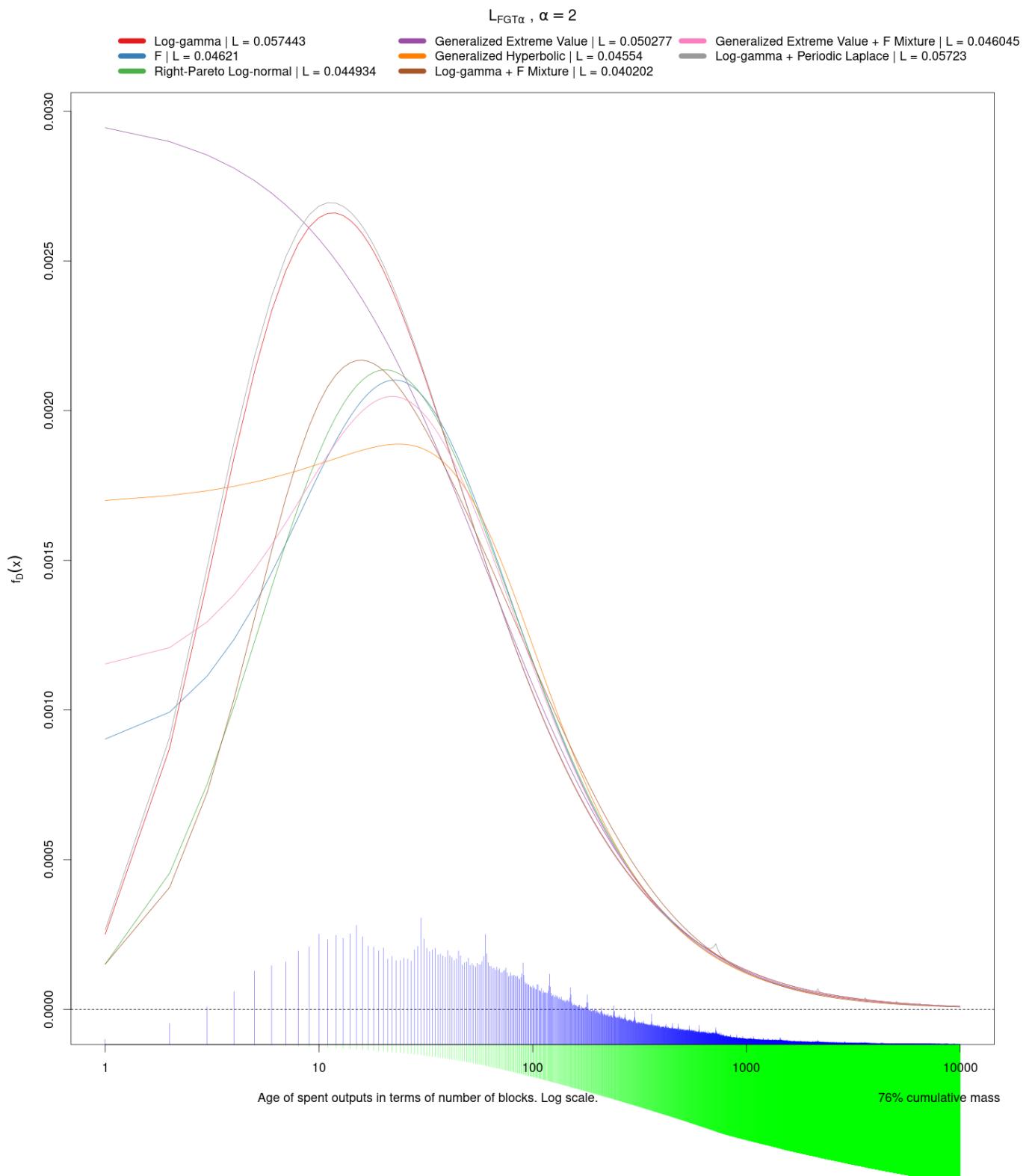
Figure 5: Optimized $f_D(x)$ for loss function L_{FGT_α} , $\alpha = 2$ 

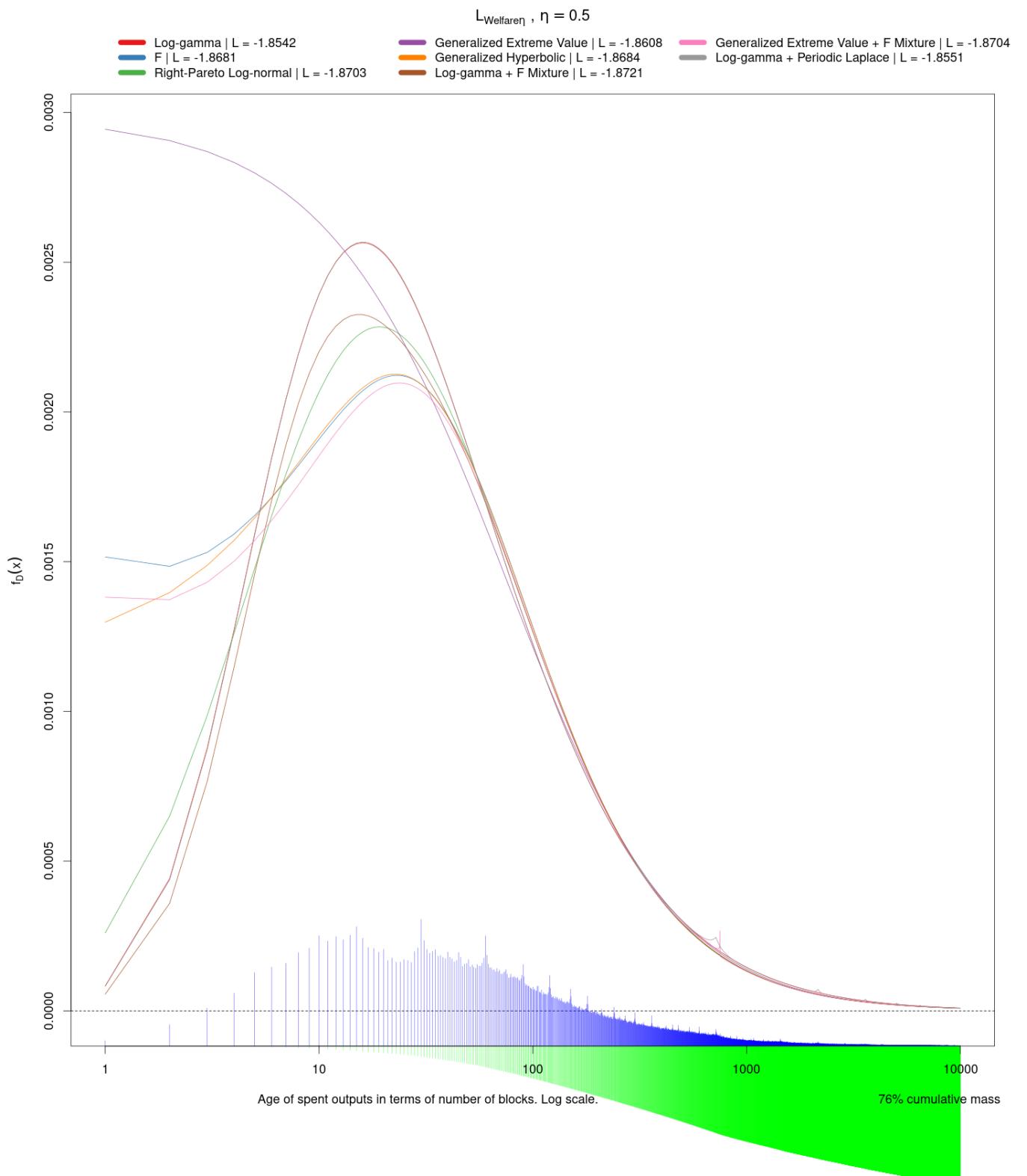
Figure 6: Optimized $f_D(x)$ for loss function $L_{Welfare_\eta}$, $\eta = 0.5$ 

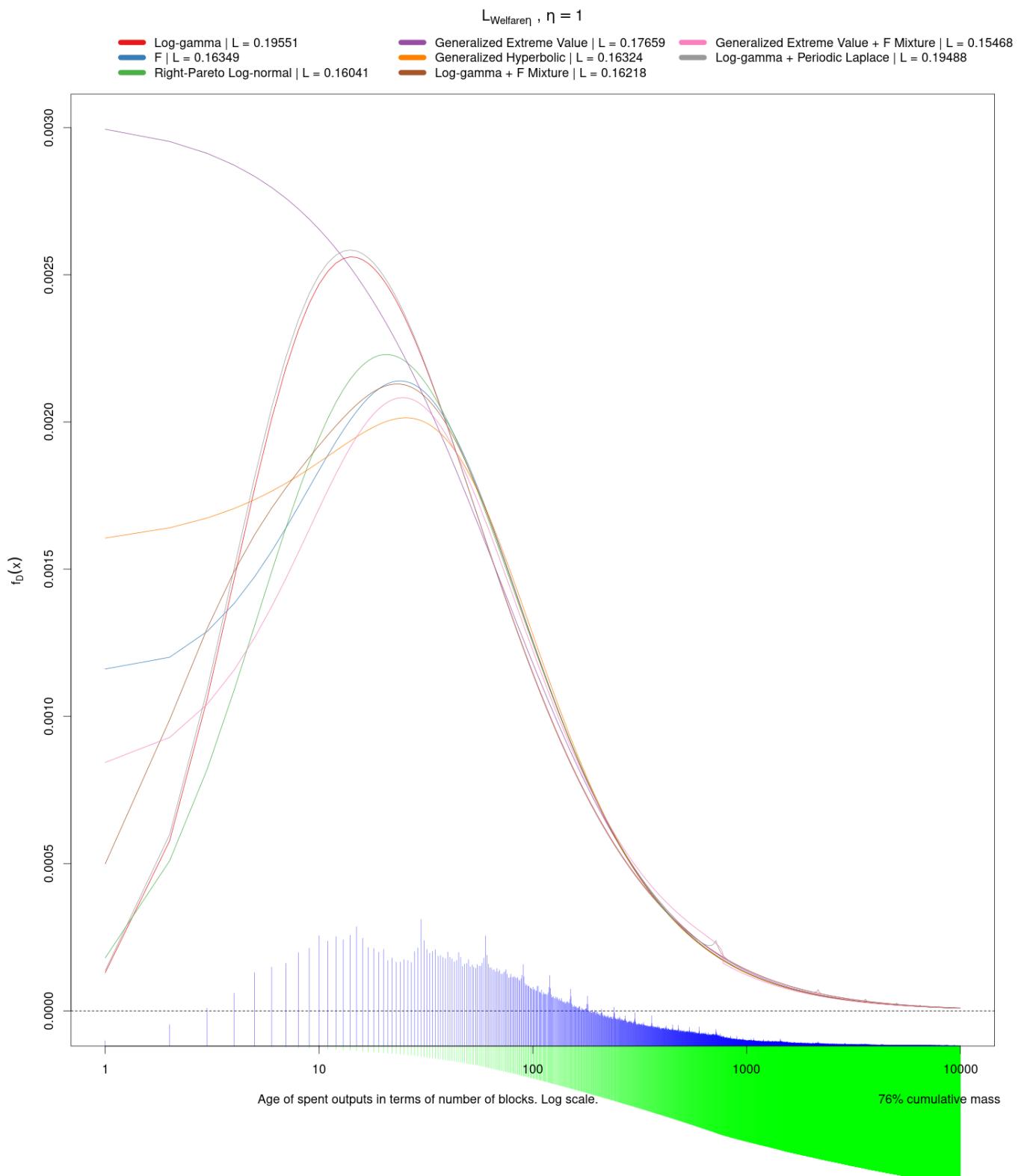
Figure 7: Optimized $f_D(x)$ for loss function $L_{Welfare_\eta}$, $\eta = 1$ 

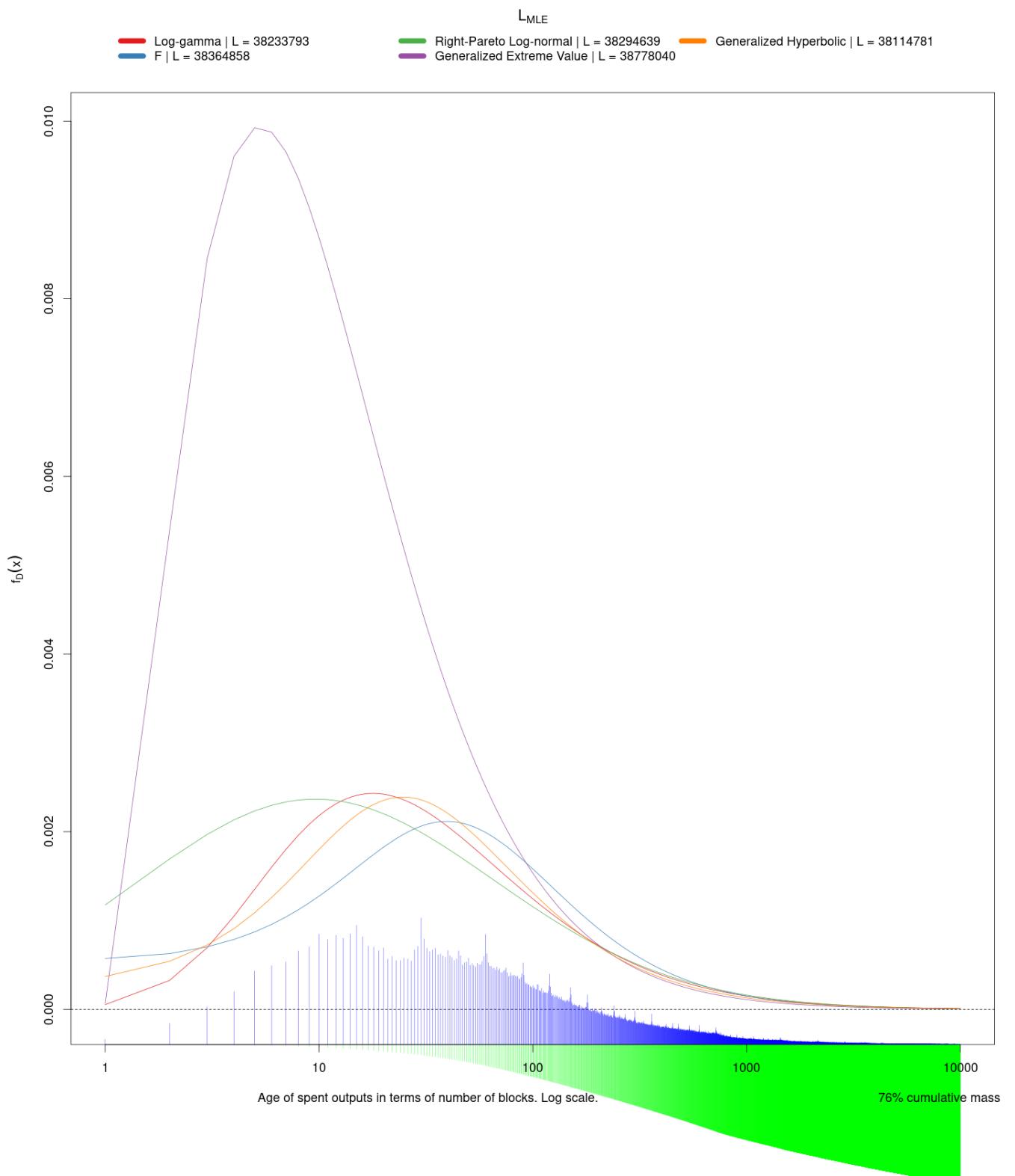
Figure 8: Optimized $f_D(x)$ for loss function L_{MLE} 

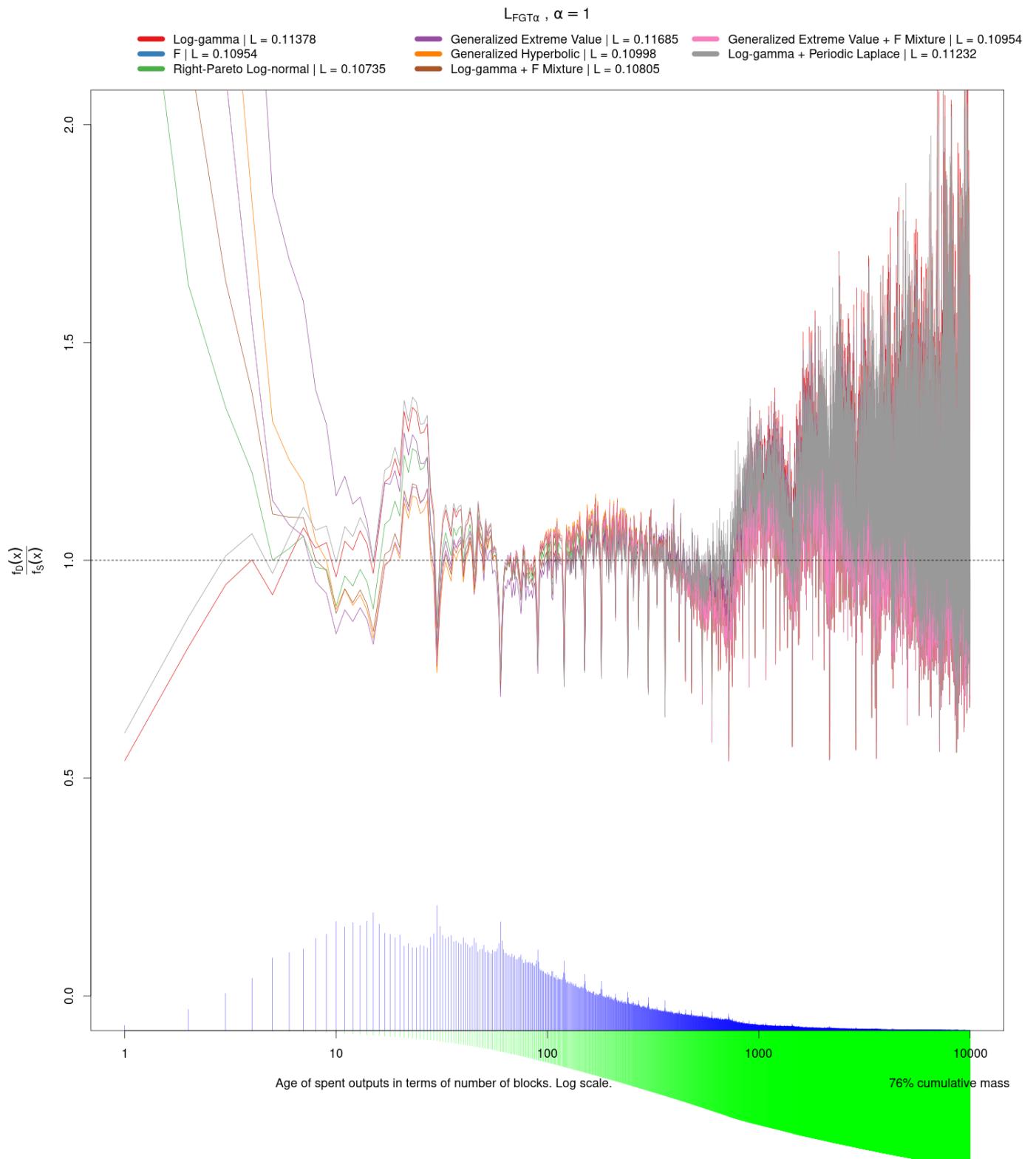
Figure 9: Optimized $f_D(x)/f_S(x)$ for loss function L_{FGT_α} , $\alpha = 1$ 

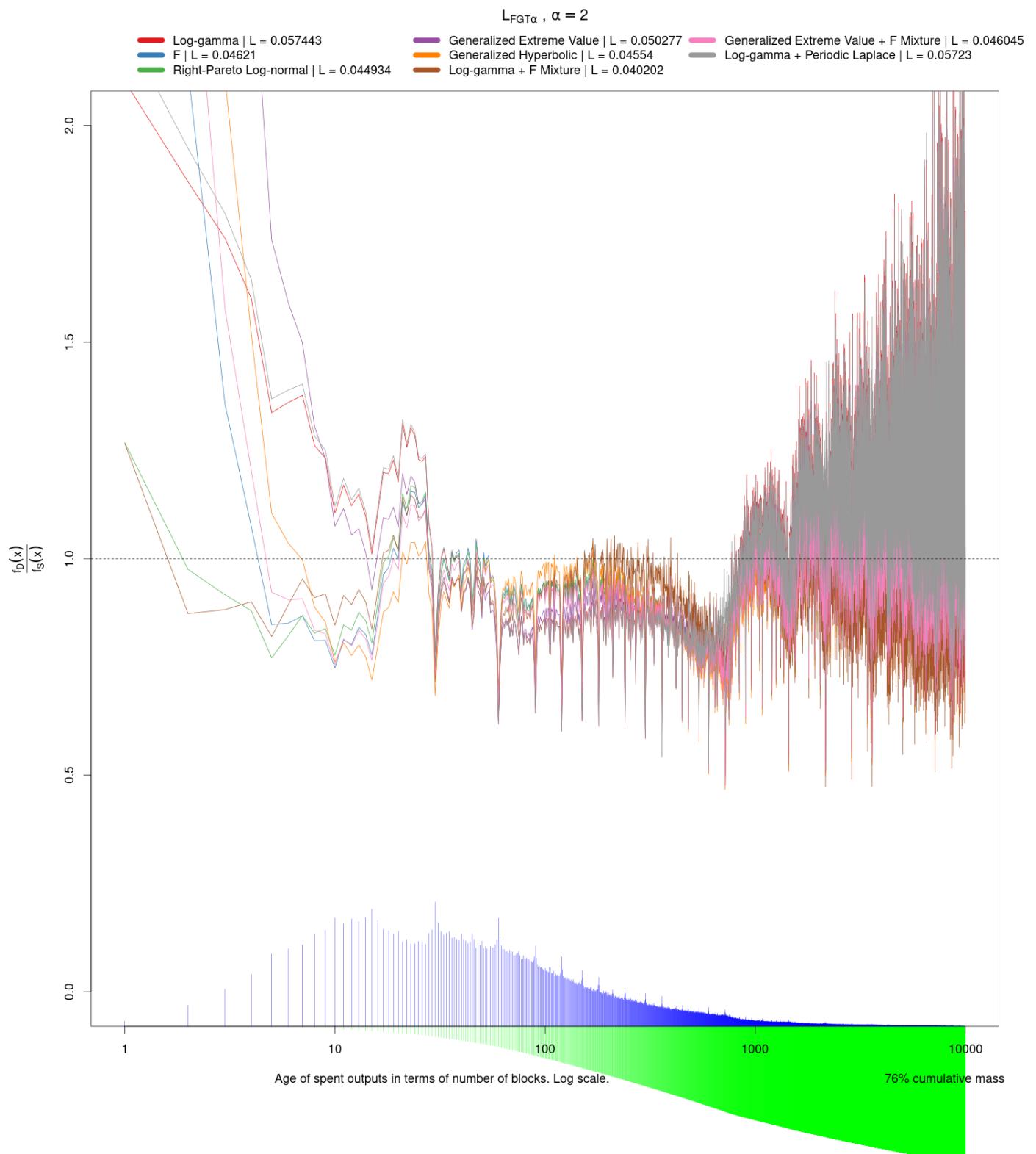
Figure 10: Optimized $f_D(x)/f_S(x)$ for loss function L_{FGT_α} , $\alpha = 2$ 

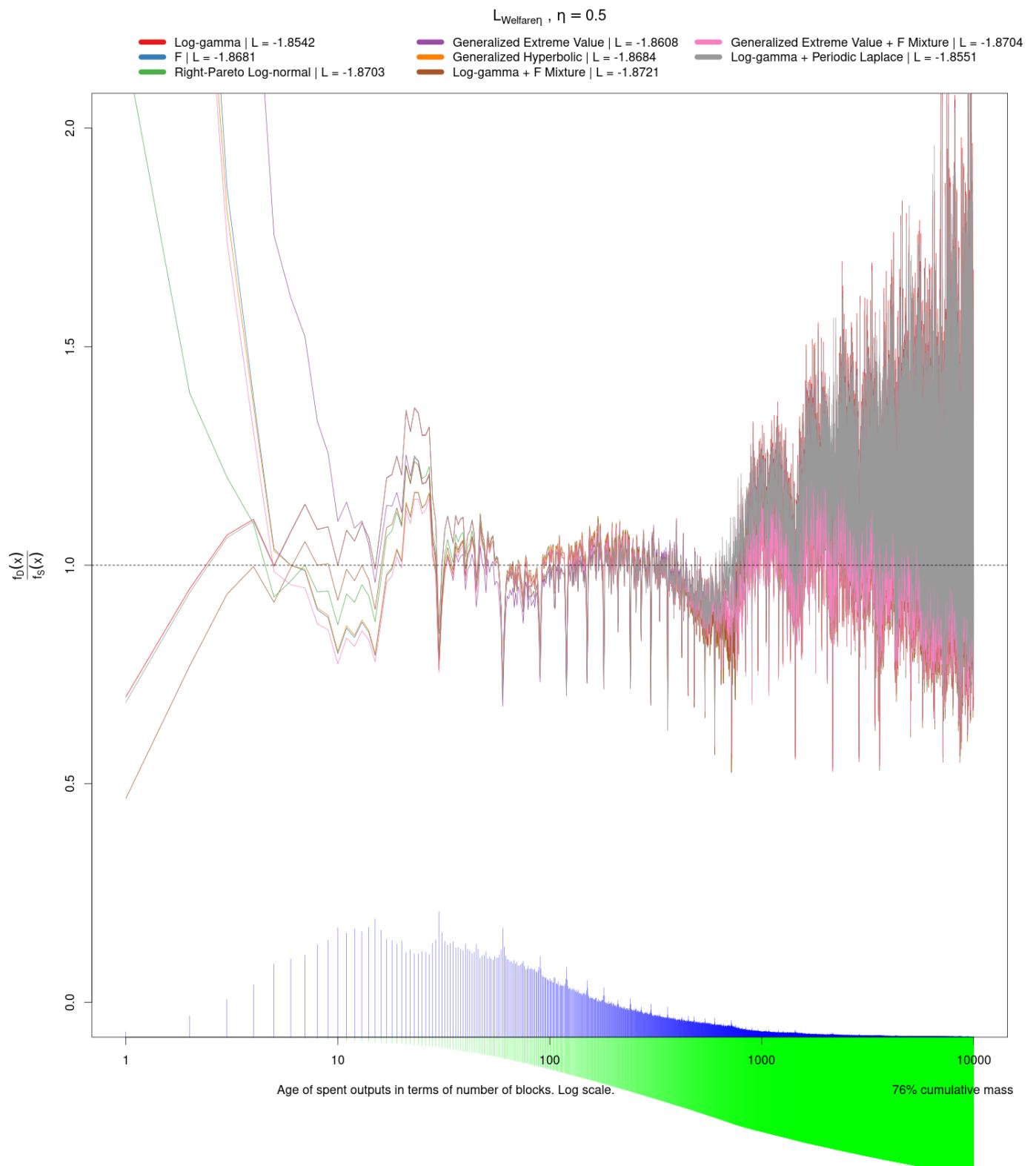
Figure 11: Optimized $f_D(x)/f_S(x)$ for loss function $L_{Welfare_\eta}$, $\eta = 0.5$ 

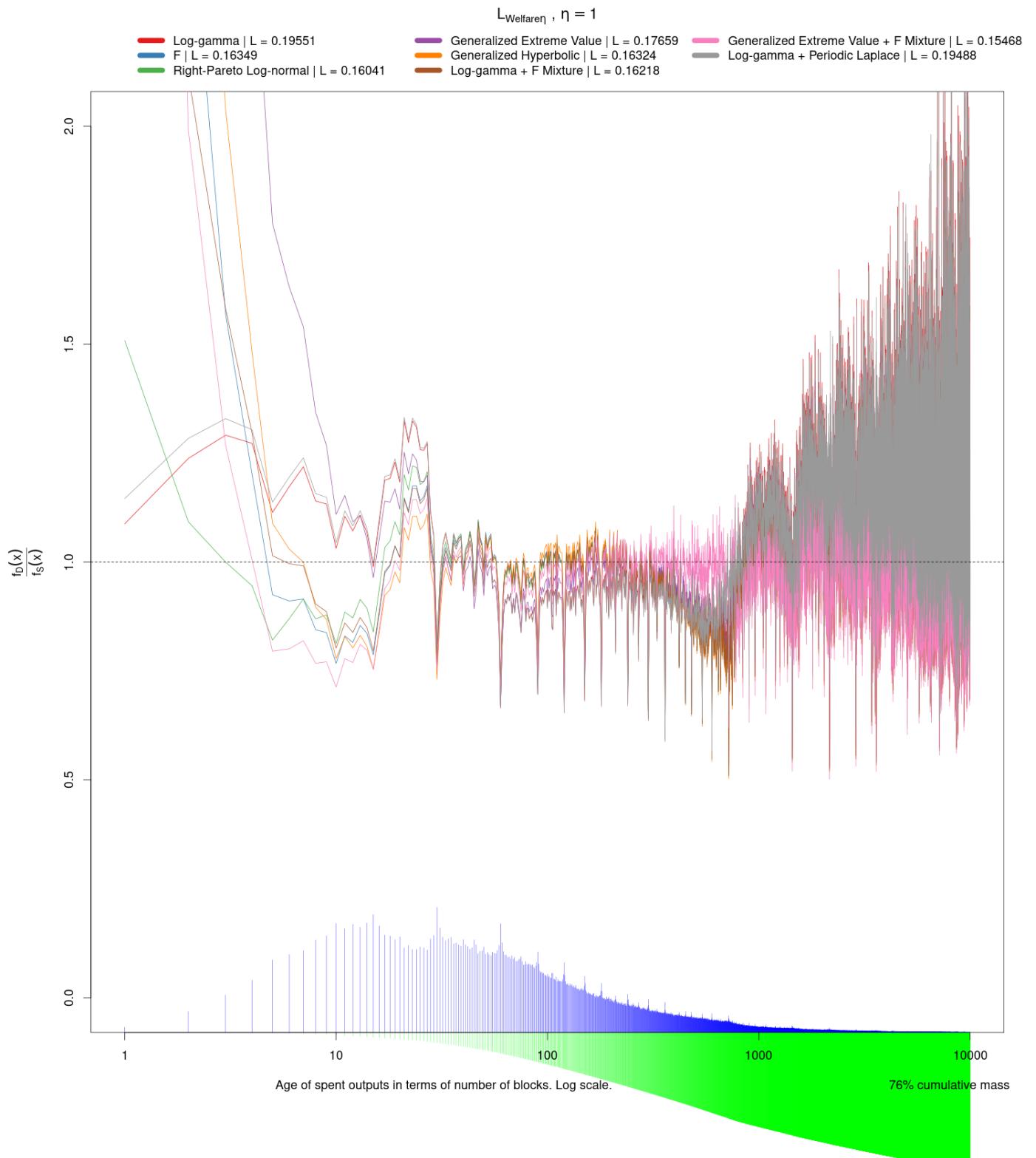
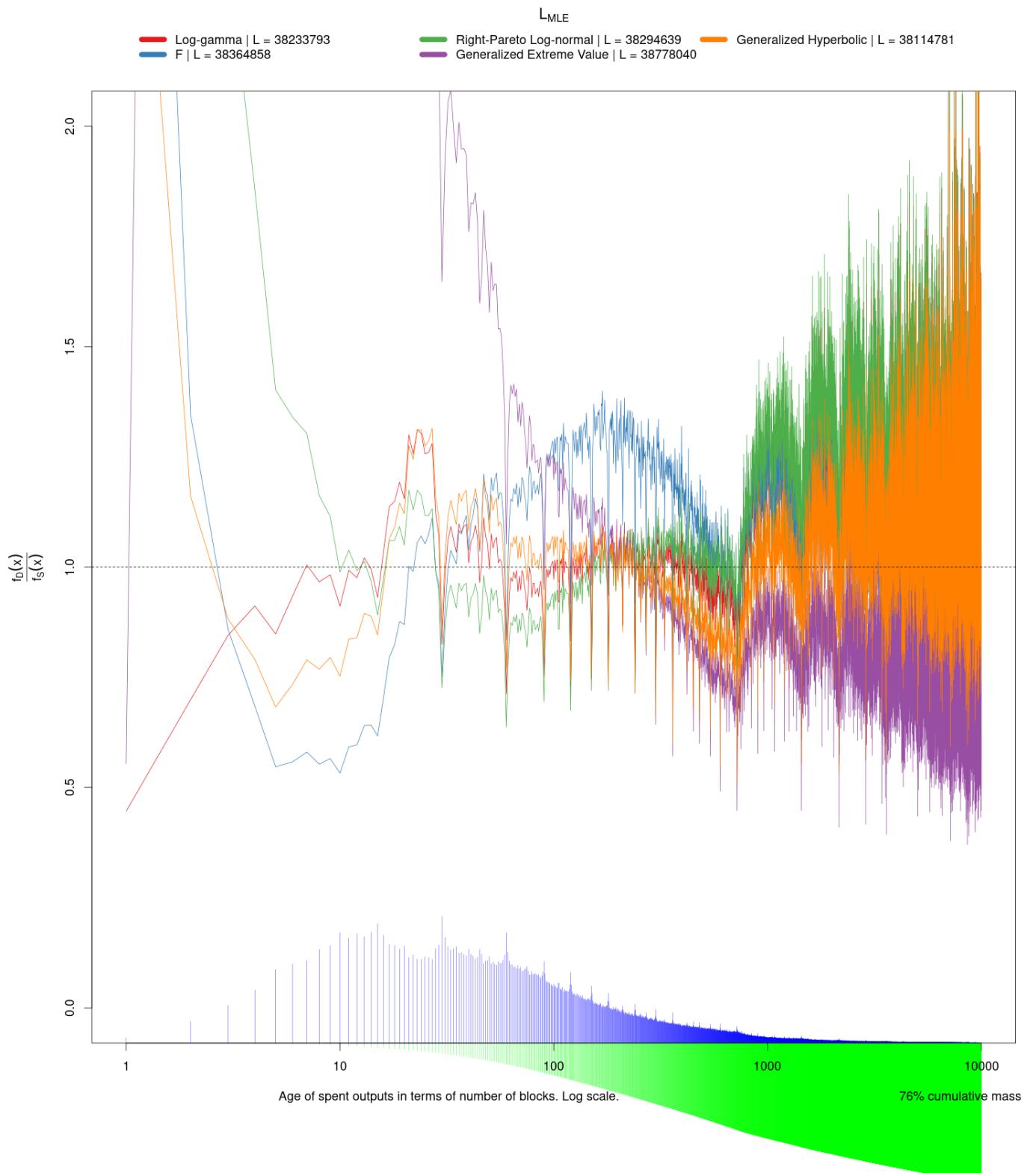
Figure 12: Optimized $f_D(x)/f_S(x)$ for loss function $L_{Welfare_\eta}$, $\eta = 1$ 

Figure 13: Optimized $f_D(x)/f_S(x)$ for loss function L_{MLE} 

23 Options for Loss Functions \mathcal{L} and Parametric Distributions \mathcal{F}

1860 In **Section 21 Criteria for Best Fit** I developed several loss functions that could be used to decide which
1861 parametric distribution (and particular parameter values) should be used to form the OSPEAD decoy selection
1862 algorithm. At this point in time, it is not necessary to decide which loss function to use. We can perform the
1863 minimizations, examine the results, and make a decision when we have all the information.

- 1865 1. Privacy impoverishment. This loss function overall appears to be one of the better choices. It is intuitive
1866 and provides options in its α parameter. To minimize risk of very low coverage of certain intervals on the
1867 probability distribution, setting $\alpha \geq 2$ would probably be best.
- 1868 2. Economic welfare. This criteria is less intuitive than privacy impoverishment, but it offers another way to
1869 adjust the risk sensitivity with its η parameter. Setting $\eta \geq 1$ will force more risk avoidance.
- 1870 3. Inequality minimization. This criteria would be most important if we seek to minimize inequality among
1871 users. One of its main drawbacks is that is can be computationally expensive.
- 1872 4. Worst-case-scenario minimization. This criteria has some appeal because worst case scenarios seem to be a
1873 focus in the field of cryptography. However, focusing on the worst case out of hundreds of thousands may be
1874 questionable.
- 1875 5. Maximum Likelihood Estimation. [Möser et al., 2018] used this criteria. For reasons stated above, I do not
1876 recommend this one.
- 1877 6. Maximize resistance to a specific attack. This criteria has some intuitive appeal: defend against an attack
1878 directly. However, it is difficult to say whether maximizing resistance to one attack may leave users vulnerable
1879 to a different attack.

1880 There is no particular statistical theoretical reason to choose one parametric distribution over another. Issues with
1881 software implementation may push us to use a simpler and more common distribution. If wallet developers do
1882 not use the `wallet2` implementation, they may have difficulty implementing a distribution defined by a mixture of
1883 unusual distributions.

24 Dynamic Risk and Forecasting

1884 The Monero Research Lab research bulletin that I quoted in the introduction continues with the following observa-
 1885 tion [Mackenzie et al., 2015]:

1887 This would suggest that we, the developers of Monero, must estimate the probability distribution gov-
 1888 erning the age of transaction outputs.

1889 This, too, is problematic. When an exchange rate is experiencing a strong long-term decline (inflation),
 1890 rational users are more likely to spend their transaction outputs, for tomorrow their outputs will be
 1891 worth less in terms of goods and services than today and hoarding is not economically rational. When
 1892 an exchange rate is experiencing a strong long-term increase (deflation), rational users are more likely to
 1893 hoard their transaction outputs for the opposite reason. Hence, the distribution of transaction output
 1894 ages will at least vary over time, and, presuming any proportion of users are rational will certainly depend
 1895 sensitively on the economic performance of the currency. It is unwise to design security recommendations
 1896 around the economic performance of our protocol.

1897 I would argue that, as long as Monero is committed to using a mimicking decoy selection algorithm, “design[ing]
 1898 security recommendations around the economic performance of our protocol” is unavoidable, if unfortunate. The
 1899 best that we can do is to minimize the risk of guessing the future distribution incorrectly. This requires an assessment
 1900 of dynamic risk with forecasting. According to empirical evidence from transparent chains, the spent output
 1901 age distribution is affected more by volatility than long-term inflation and deflation. [Makarov & Schoar, 2021]
 1902 concluded that “the vast majority of Bitcoin transactions between real entities are for trading and speculative
 1903 purposes. Starting from 2015, 75% of real bitcoin volume has been linked to exchanges or exchange-like entities
 1904 such as on-line wallets, OTC desks, and large institutional traders.” I discussed my own findings for DOGE in the
 1905 June 1, 2022 Monero Research Lab meeting.²⁹ The transaction patterns of Monero may be somewhat different due
 1906 to its absence from many centralized exchanges, but probably not completely different. Assume for a moment that
 1907 changes in the real spend age distribution are entirely caused by exchange rate volatility. In that case, forecasting
 1908 the real spend age distribution would be roughly equivalent to forecasting exchange rate volatility and therefore
 1909 would be highly challenging. The Efficient Market Hypothesis would imply that a naive forecast may be the best
 1910 forecast.

1911 We return to Equation (18) from **Section 21 Criteria for Best Fit**: the risk function $R(\theta, \delta) = E_\theta L(\theta, \delta(\mathbf{X}))$.
 1912 This is the expected value of our choice of the loss function when a particular set of parameters, e.g. the shape,
 1913 scale, and location of a parametric distribution, are set to particular values. The theoretical bounds on $R(\theta, \delta)$
 1914 can be very wide in our setting and do not grant us much guidance. Therefore, it is best to compute some sort of
 1915 empirical risk function based on data. All we have is past data, yet our goal is to compute future risk. We must
 1916 forecast future data and then determine performance of our various options of decoy selection algorithms through
 1917 forecast validation.

1918 How to do forecasting and forecast validation? A naive, simple method to forecast future values is to assume
 1919 that the values will be the same as in the current period. More sophisticated methods attempt to anticipate changes
 1920 in the forecasted values by analyzing cycles and statistical dependence of the values. In our setting with multi-
 1921 variate forecasting, some appropriate forecasting methods include Exponential Smoothing, Vector Autoregression
 1922 (VAR), Vector Autoregressive Moving Average (VARMA), Generalized Autoregressive Conditional Heteroskedas-
 1923 ticity (GARCH), and Kalman Filter.

1924 Cross-validation is a popular general method to evaluate model accuracy. There are special considerations
 1925 when the data is time series because time series data is not necessarily independent nor identically distributed.
 1926 The general consensus in the academic literature is that when the characteristics of the time series are unknown
 1927 then “out-of-sample” (OOS) cross-validation should be used. OOS cross-validation fits a forecasting model on

²⁹<https://libera.monerologs.net/monero-research-lab/20220601#c103366>

1928 all the data except the last few periods and then measures how well the model forecasts the data of the last few
 1929 periods. OOS cross-validation is also known as “last block validation”, “forward validation”, and “holdout validation”
 1930 ([Bergmeir & Benítez, 2012], [Cerqueira et al., 2020]).

1931 To be specific, what I mean by “unknown” characteristics of a time series is that it is unknown whether the time
 1932 series is stationary. Stationarity is a technical condition in time series that requires that its distribution does not
 1933 depend on time. An independent and identically distributed series is stationary. A series that is neither independent
 1934 nor identically distributed will not be stationary. Many financial time series are non-stationary. A common way to
 1935 manage non-stationarity in data is to compute the first difference (or second difference if necessary, or third...), but
 1936 it is not clear how this could be done when the object of observation is an entire distribution, as it is in our setting.

1937 [Bergmeir et al., 2018] mathematically prove that standard K -fold cross-validation is a valid technique when
 1938 the data series is stationary (plus a few other assumptions). They warn that K -fold cross-validation is not valid if
 1939 the data series is non-stationary and show through Monte Carlo simulations that OOS validation performs better
 1940 on non-stationary data. [Bergmeir & Benítez, 2012] arrive at a similar conclusion. [Cerqueira et al., 2020] separate
 1941 dozens of real data sets into stationary and non-stationary and finds that “when the time series are non-stationary,
 1942 the most accurate estimates are produced by out-of-sample methods, particularly the holdout approach repeated
 1943 in multiple testing periods....When the observations in a data set are not i.i.d. [independently and identically
 1944 distributed], the standard cross-validation approach is not directly applicable.” The top performing models in the
 1945 M5 forecasting competition generally used the last 28 day period for cross-validation ([Makridakis et al., 2022]).

1946 There is no particular reason to think that the spent output age distribution is stationary. Certainly, research
 1947 into spent output age is in its infancy. Non-stationarity is the weaker assumption. To be “safe” it is better to impose
 1948 fewer assumptions on statistical models if we cannot justify them. Furthermore, with only about 52 weeks of Monero
 1949 data it would be difficult to perform a formal hypothesis test of stationarity with high statistical power. Therefore,
 1950 it is best to use OOS cross-validation to evaluate forecast accuracy and risk. A form of OOS cross-validation is
 1951 performed in the following section.

1952 The `tsqsim` software created by Monero developer mj-xmr is capable of OOS cross-validation.³⁰ Some modifi-
 1953 cations will be necessary to support multivariate time series. Technically, spent output age is univariate, but since
 1954 we are measuring the age *distribution* at each time period, it makes sense to model it in a multivariate way.

1955 25 Inter-Temporal Stability of Spent Output Age Distribution for BTC, 1956 BCH, LTC, and DOGE

1957 Much like in **Section 22 Dry Run with Old Moser et al. (2018) Data**, it is useful to run an analysis of
 1958 similar spent output age data from other blockchains to inform statistical modeling strategies.

1959 The table below lists the share of payments made with several different cryptocurrencies with an unspent
 1960 transaction output (UTXO) model. I ignored cryptocurrencies with an account model like Ethereum because the
 1961 spending mechanism is fundamentally different. The main question that we want to answer is whether the chosen
 1962 cryptocurrencies are used in a roughly similar manner to Monero: as peer-to-peer electronic cash. This table may
 1963 suggest that BTC, BCH, LTC, and DOGE have somewhat similar usage as Monero. One major difference between
 1964 these cryptocurrencies and Monero is that there is a 10 block lock enforced at the protocol level for Monero. BTC,
 1965 BCH, LTC, and DOGE allow users to spend received funds immediately.

³⁰<https://github.com/mj-xmr/tsqsim>

UTXO-based cryptocurrency usage as payment by percent (payment processors and merchants)							
	Service	Bitpay ³¹	CoinCards ³²	Bitrefill ³³	Cake Pay ³⁴	Travala ³⁵	Shodan ³⁶
1967	TXs/month	67,000	NA	186,000	NA	NA	220 total
	Time period	June 2022	July 2022	NA	Aug 2022	NA	NA
	Metric	Transactions	USD value	Unique users	USD value	Hotel nights	Subscriptions
	BTC	53.3	41.0	40	16	9.2	47.3
	BCH	5.0	NA	NA	NA	0.5	4.1
	LTC	21.2	7.0	7	6	0.8	12.7
	DOGE	6.2	3.0	NA	NA	0.3	13.2
	XMR	NA	19.0	NA	78	1.9	NA
	DASH	NA	0.5	NA	NA	0.5	NA
	ADA	NA	NA	NA	NA	0.9	NA

1968 In the next several section I perform statistical analysis of these four UTXO cryptocurrencies. The code is
 1969 available at <https://github.com/Rucknium/OSPEAD/tree/main/General-Blockchain-Age-of-Spent-Outputs>. The
 1970 code was executed on the Monero Research Computing Server, whose hardware was improved by a CCS proposal.³⁷

1971 25.1 Summary Characterization of Distributions Over Time

1972 Figures 15 to 18 show several statistics about the age of spent outputs of BTC, BCH, LTC, and DOGE since 2015.
 1973 The age units are in terms of blocks. For BTC and BCH the interval between blocks is 10 minutes. For LTC it
 1974 is 2.5 minutes and for DOGE it is 1 minute. The unit of observation is the ISO week, a natural unit of economic
 1975 time.³⁸

1976 The first line graphs show the mean, median, standard deviation, skewness, and kurtosis. The skewness and
 1977 kurtosis statistics may be unfamiliar. They are the third and fourth standardized moments of a distribution,
 1978 respectively. The moments of a distribution is defined by a power of the expectation $E[X]$, i.e. theoretical mean, of
 1979 the distribution. The k th moment is $E[X^k]$. Moments extend the concept of moving from expectation to variance
 1980 (which is the square of the standard deviation). The mean (expectation) of a random variable is simply the first
 1981 moment:

$$\mu = E[X^1]$$

1982 The variance is the second central moment:

$$\sigma^2 = E \left[(X - E[X])^2 \right]$$

1983 The skewness is the third standardized moment (standardized by the standard deviation σ):

$$\tilde{\mu}_3 = E \left[\left(\frac{X - E[X]}{\sigma} \right)^3 \right]$$

1984 The kurtosis is the fourth standardized moment:

³¹<https://bitpay.com/stats/>

³²<https://twitter.com/CoinCards/status/1555286172385116164>

³³<https://youtu.be/bkjEcSmZKfc?t=549>

³⁴<https://twitter.com/cakewallet/status/1565370179906838528>

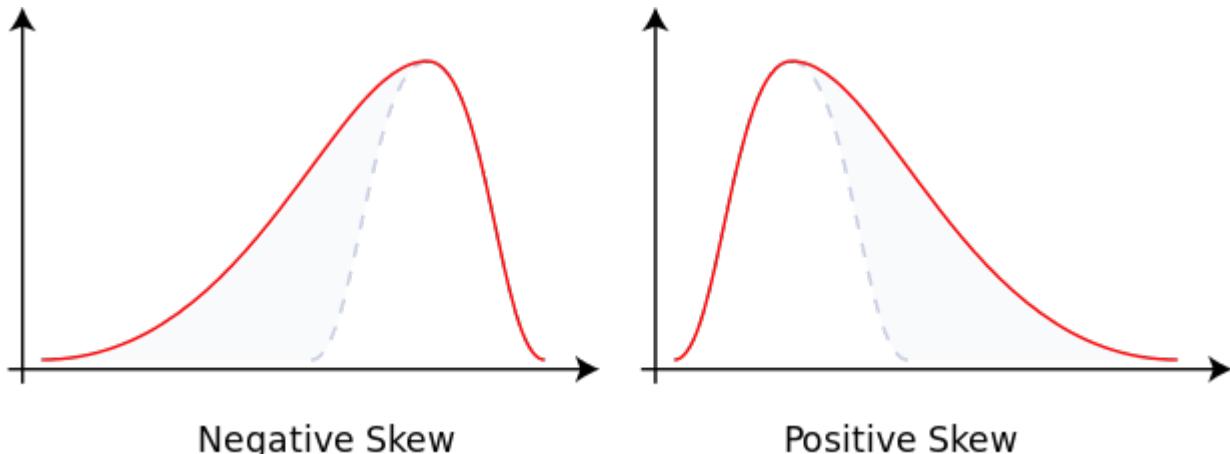
³⁵<https://travala-dashboard.com/>

³⁶<https://blog.shodan.io/accepting-crypto-a-vendor-perspective/>

³⁷https://ccs.getmonero.org/proposals/gingeropolous_z zenith _storage.html

³⁸https://en.wikipedia.org/wiki/ISO_week_date

Figure 14: Skewness



Source: [https://en.wikipedia.org/wiki/File:Negative_and_positive_skew_diagrams_\(English\).svg](https://en.wikipedia.org/wiki/File:Negative_and_positive_skew_diagrams_(English).svg)

$$\tilde{\mu}_4 = E \left[\left(\frac{X - E[X]}{\sigma} \right)^4 \right]$$

1985 The mean is a measure of the central tendency of a distribution. The standard deviation is a measure of its
 1986 dispersion (spread). The skewness and kurtosis of a distribution involve its characteristics in its tail or tails. A
 1987 positive skew means that the distribution's tail is on the right side of the distribution. All the age distributions
 1988 analyzed here tend to have a positive skew. High kurtosis means that large outliers are more likely. Kurtosis
 1989 of greater than 3 suggests that a distribution has a fatter tail than the normal distribution. The skewness and
 1990 kurtosis become relevant when very old outputs “wake up” during periods of exchange rate volatility to participate
 1991 in speculative activity, i.e. buying and selling on exchanges.

1992 The fitting function is essentially a minimum divergence estimator. This means that the parameters of the
 1993 parametric probability density function (PDF) are chosen to minimize the distance between the parametric PDF
 1994 and the PDF formed by the data (the empirical PDF), for some specified metric of “distance”.

1995 There are several measures of distance that could be used. For the purpose of this exploration of output age
 1996 distribution forecasting, the distance metric to be minimized will be the total linear sum of the mass of the estimated
 1997 parametric PDF that falls below the empirical PDF. With the “loss function” specified this way, the optimization
 1998 algorithm attempts to minimize the probability that the real spends are much more likely to come from a block of
 1999 a particular age compared to a potential decoy. Ideally, the decoy distribution would be identical to the real spend
 2000 age distribution, but parametric PDFs are not flexible enough to perfectly match an empirical PDF.

2001 Define $f_S(x)$ as the empirical spent output age distribution at block age x and $f_D(x; \beta)$ as a potential “decoy”
 2002 distribution with some parameter vector β (with the transparent blockchains presented here, there is no actual
 2003 decoy mechanism of course). Then for each week of spent output age data the parameter vector β can be chosen
 2004 to minimize this quantity:

$$L(\beta) = \sum_{i \in \{x_i : f_D(x_i; \beta) < f_S(x_i)\}}^N f_S(x_i) - f_D(x_i; \beta)$$

2005 That is, minimize the sum of the difference between the real spend age distribution for blocks x_i where the decoy
 2006 distribution is less than the real spend age distribution. The minimization is performed by computer numerical
 2007 minimization methods similar to gradient descent.

2008 The two candidate “decoy” parametric PDFs under consideration in this exercise are the Log-gamma (*lgamma*)

2009 distribution and the Right-Pareto Log-normal (*rpln*) distribution. The lgamma distribution, with two parameters,
 2010 was used in [Möser et al., 2018] to suggest a decoy distribution that was later incorporated into Monero's reference
 2011 wallet software. The rpln distribution is a more flexible distribution, with three parameters. The PDFs of these
 2012 two distributions are:

$$f_{lgamma}(x) = \frac{ratelog^{shapelog}}{\Gamma(shapelog)} \times \frac{(\ln x)^{shapelog-1}}{x^{ratelog+1}}$$

2013 with $ratelog > 0$ and $shapelog > 0$ and where Γ is the gamma function and

$$f_{rpln}(x) = shape2 \times x^{-shape2-1} e^{shape2 \times meanlog + \frac{shape2^2 \times sdlog^2}{2}} \Phi\left(\frac{\ln x - meanlog - shape2 \times sdlog^2}{sdlog}\right)$$

2014 with $shape2 > 0$ and $sdlog > 0$ and where Φ is the cumulative distribution function of the standard Normal
 2015 distribution.

Figure 15: BTC Summary Characterization of Spend Age Distribution Over Time

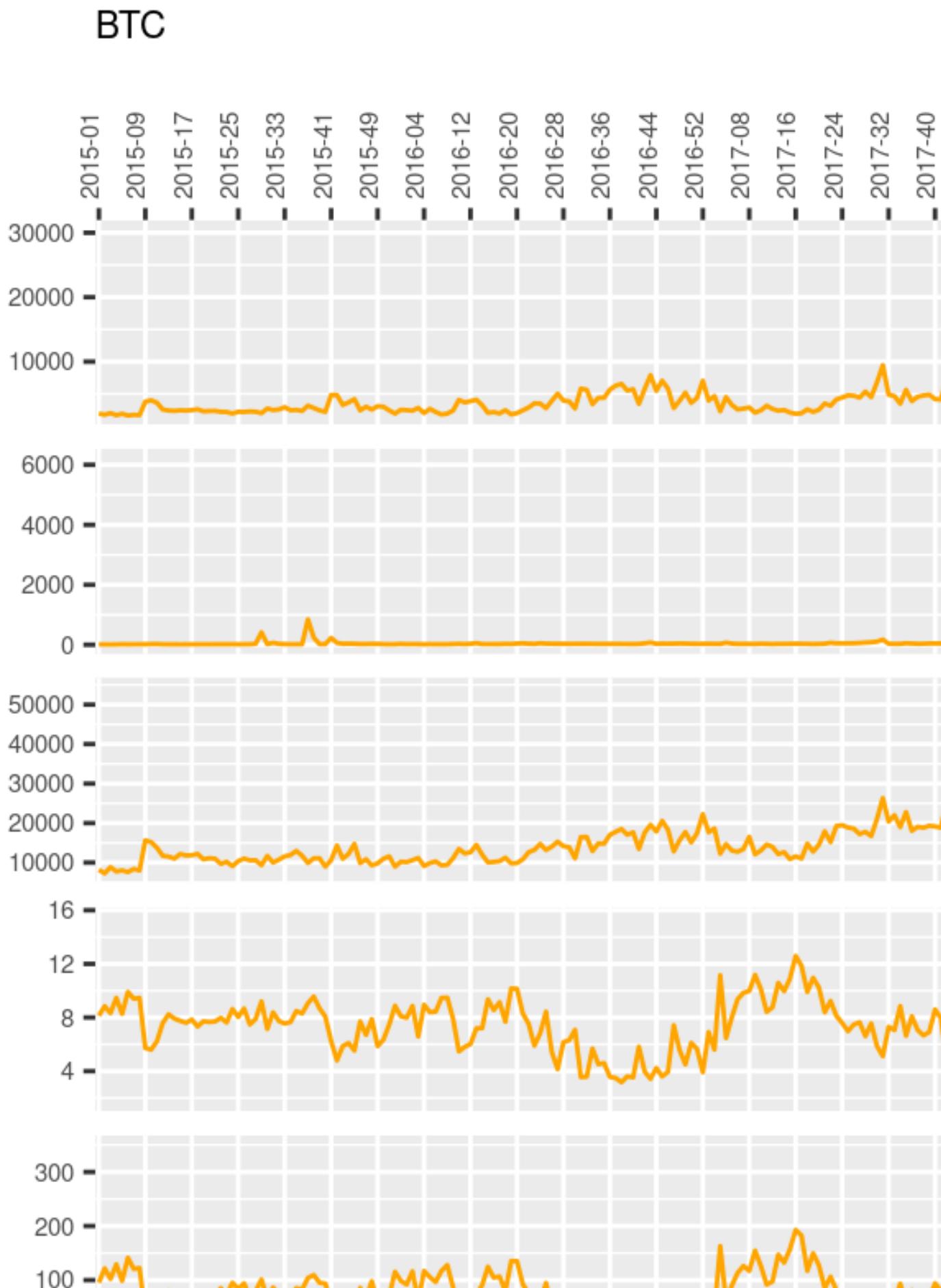


Figure 16: BCH Summary Characterization of Spend Age Distribution Over Time

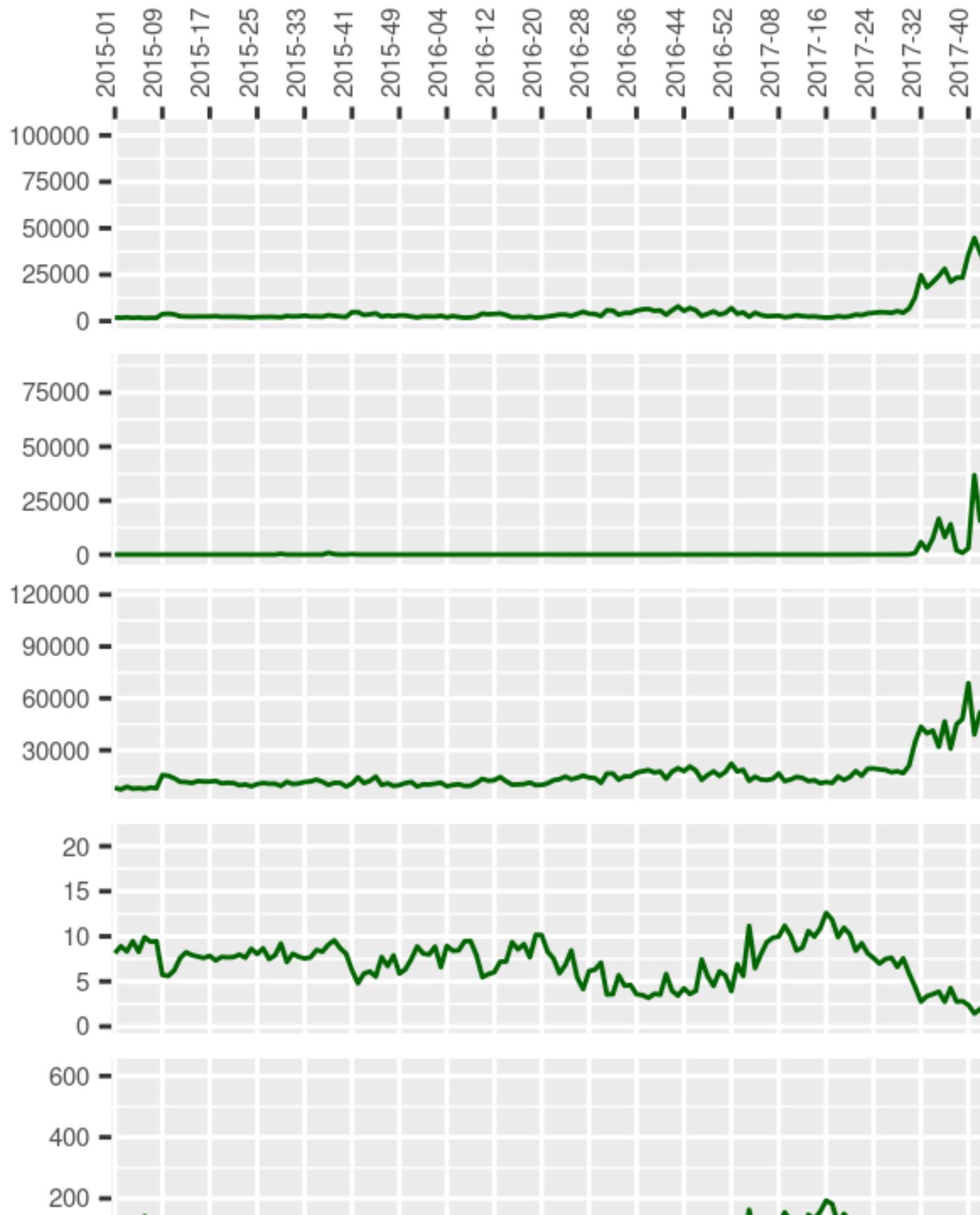
BCH

Figure 17: LTC Summary Characterization of Spend Age Distribution Over Time

LTC

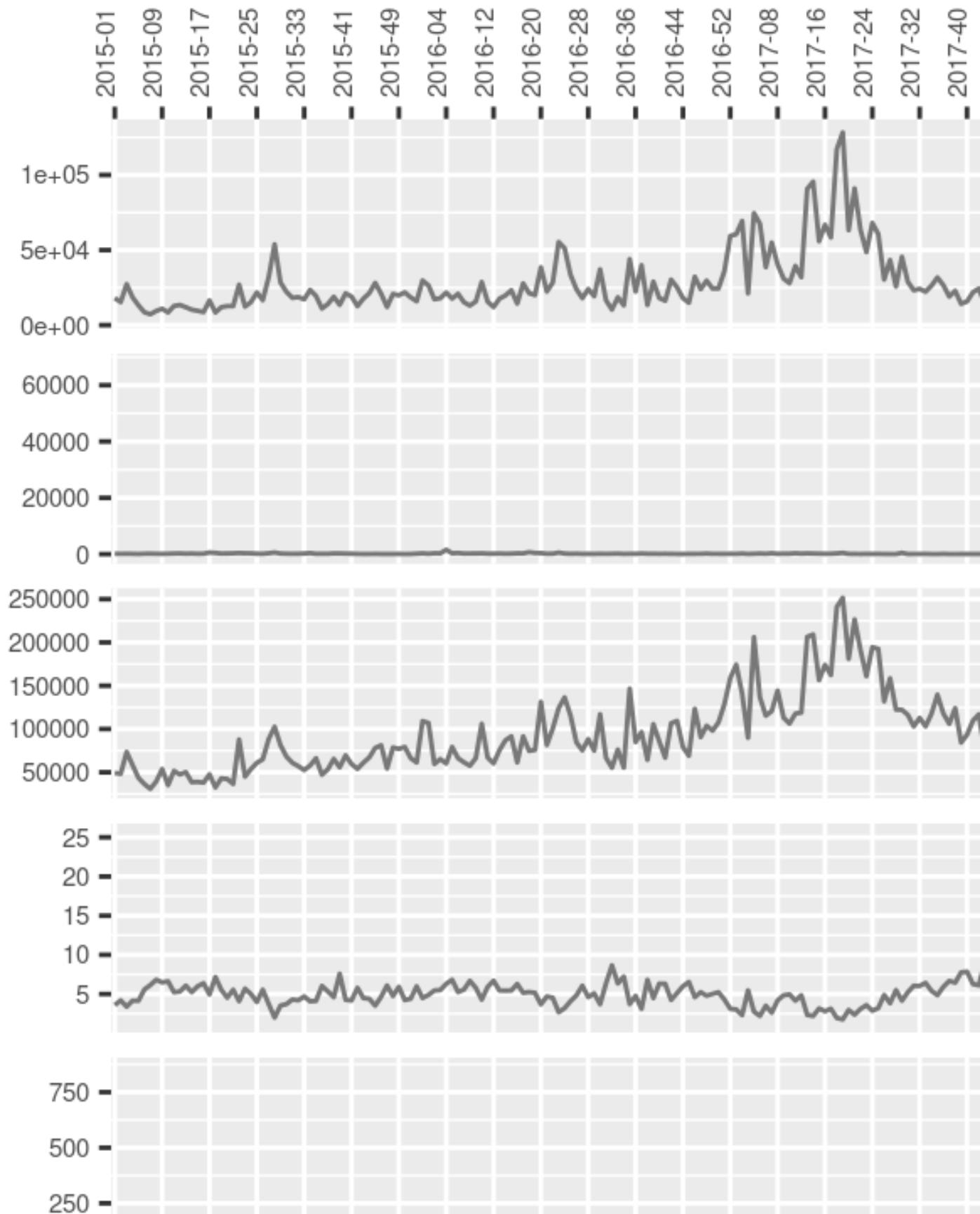


Figure 18: DOGE Summary Characterization of Spend Age Distribution Over Time

DOGE



2016 **25.2 Cross-Blockchain Correlations of Summary Statistics Across Time**

2017 The following table contains the correlation between several statistics over time for each pair of blockchains. The
 2018 “BTC&BCH” correlation included data only for January 2018 onward to avoid artificially raising the correlation by
 2019 including weeks when the BTC and BCH contained the same transaction before the August 2017 hard fork.

2020 To be precise, define vector $[x_{s,1}, x_{s,2}, \dots, x_{s,Z}] = \mathbf{X}_s$ for statistic s , e.g. the median, of blockchain x at each
 2021 week, with Z total weeks in the sample. Define \mathbf{Y}_s for blockchain y similarly. Then the quantities displayed in
 2022 Table 6 are $\text{corr}(\mathbf{X}_s, \mathbf{Y}_s)$.

2023 The two statistics that tend to have consistently high correlation for each pair of blockchains are skewness and
 2024 kurtosis.

Table 6: Cross-Blockchain Correlations of Summary Statistics Across Time

	BTC&BCH	BTC<C	BTC&DOGE	BCH<C	BCH&DOGE	LTC&DOGE
mean	-0.04	-0.04	-0.04	-0.01	-0.08	0.21
median	-0.03	-0.01	-0.01	-0.03	-0.02	-0.01
sd	-0.07	0.02	0.34	-0.02	0.02	0.35
skewness	0.09	0.24	0.19	0.11	-0.32	0.28
kurtosis	0.02	0.29	0.31	0.22	-0.26	0.23

2025 **25.3 Evaluation of Forecast Accuracy**

2026 OSPEAD needs to approximate the real spend age distribution in the future, not the past. Ideally, the decoy
 2027 selection algorithm should mimic the real spend age distribution at the time that each transaction is made. Since
 2028 the real spend age distribution is likely changing over time, some type of forecasting is needed. Here I perform some
 2029 simple evaluation of forecasting methods on transparent blockchains.

2030 Let W be the number of weeks in a forecast horizon. It is the approximate interval of time that a proposed
 2031 decoy selection algorithm should aim to be accurate. Let M be the number of weeks in the data sample. Then a
 2032 common way to evaluate a forecasting method is to use the first $M - W$ weeks to fit a forecasting model. Then
 2033 determine the performance of the method by forecasting for W weeks into the “future” and compare with the actual
 2034 sample data of the final W weeks. This is known as out-of-sample validation. For this exercise, I set $W = 8$.

2035 An entire distribution needs to be forecast rather than a single value or a set of values. Forecasting an entire
 2036 (empirical, nonparametric) distribution is extremely complex, so here I reduce the magnitude of the forecasting
 2037 problem to just the parameters of the fitted parametric distributions (lgamma and rpln).

2038 As an initial test, the sophisticated forecasting method used here was a multivariate auto-regressive(1) exoge-
 2039 nous inputs state-space model with a Koopman-Durbin Kalman filter implemented in the MARSS R package. The
 2040 forecasting method allows a forecast into each period in the future. Since the “S” in OSPEAD stands for “Static”,
 2041 only one of the forecast periods can be chosen. I chose the 4th forecast period since it is roughly in the middle of
 2042 the 8-week forecast horizon. This forecast is referred to as simply `forecast.accuracy` below.

2043 I also evaluated two naive forecasting methods. The first is to use the parameters of the $M - W$ th week, i.e. the
 2044 last week of the “training set”, as the forecast. This is called `forecast.accuracy.naive.final.week` below. The
 2045 second naive method (`forecast.accuracy.naive.horizon.period`) is to use the final W weeks of the training set
 2046 to fit the specified lgamma and rpln distributions by minimizing the $L(\beta)$ loss function for a set of weeks of data
 2047 pooled together rather than a single week.

2048 Given that the loss function $L(\beta)$ is not globally convex in the choice parameters β , the numerical optimization
 2049 algorithm can go “off the rails” and settle at a local rather than global minimum. Such a failure mode is especially
 2050 likely if the empirical distribution is unusual or if the starting values given to the optimization algorithm are far from
 2051 the global minimum. Generally, the solution to this problem is to try different optimization algorithms, determined

2052 better starting values, or use a computationally expensive grid search method. I chose to defer dealing with the
2053 issue to a later stage in the process. For the purposes of forecasting, I removed any weeks where the estimated
2054 parameters of the parametric distributions were greater than (or less than) five times the 95th (5th) percentile from
2055 the median. For example, this exclusion step caused the forecast evaluation for BCH to use the 15th through 24th
2056 weeks of 2022 rather than the 20th through 27th week as for other blockchains.

2057 The results of the forecast evaluation are in Table 7. Forecasts were evaluated based on the value of their loss
2058 function $L(\beta)$ for each out-of-sample week. Lower values indicate better performance. The minimum possible value
2059 is 0 and the maximum possible value is 1. The color code scales were calculated separately for each blockchain.
2060 Each color claims a value bin of equal size. Higher values are red, average values are yellow, and low values are
2061 green.

2062 In this preliminary exercise, `rpln.forecast.accuracy.naive.horizon.period` appears to have the most con-
2063 sistently good performance across blockchains.

Table 7: Evaluation of forecast accuracy

		BTC											
	rownames	mean	sd	2022-20	2022-21	2022-22	2022-23	2022-24	2022-25	2022-26	2022-27		
rphn_forecast_accuracy	0.1225	0.0294	0.0952	0.1086	0.1058	0.1227	0.0926	0.1814	0.1319	0.1421			
rphn_forecast_accuracy_naive_final_week	0.1213	0.0295	0.0951	0.1068	0.1037	0.1222	0.0910	0.1807	0.1313	0.1397			
rphn_forecast_accuracy_naive_horizon_period	0.0682	0.0244	0.0795	0.0504	0.0541	0.0601	0.0501	0.1237	0.0645	0.0632			
lgamma_forecast_accuracy	0.1256	0.0215	0.1216	0.1103	0.1151	0.1149	0.1059	0.1733	0.1299	0.1341			
lgamma_forecast_accuracy_naive_final_week	0.1495	0.0277	0.1179	0.1373	0.1373	0.1519	0.1188	0.1996	0.1613	0.1717			
lgamma_forecast_accuracy_naive_horizon_period	0.1338	0.0177	0.1302	0.1183	0.1274	0.1212	0.1205	0.1716	0.1358	0.1453			
BCH													
rownames	mean	sd	2022-15	2022-16	2022-18	2022-19	2022-20	2022-21	2022-22	2022-23	2022-24	2022-25	
rphn_forecast_accuracy	0.2416	0.1923	0.1811	0.1420	0.1147	0.1598	0.1057	0.4079	0.1665	0.6547			
rphn_forecast_accuracy_naive_final_week	0.2452	0.1808	0.1882	0.1530	0.1302	0.1668	0.1208	0.3957	0.1692	0.6375			
rphn_forecast_accuracy_naive_horizon_period	0.3938	0.0463	0.3837	0.3886	0.3754	0.3830	0.3809	0.3812	0.3526	0.5051			
lgamma_forecast_accuracy	0.3654	0.1374	0.2606	0.2757	0.2574	0.3479	0.2506	0.3830	0.5578	0.5906			
lgamma_forecast_accuracy_naive_final_week	0.3573	0.1491	0.2401	0.2553	0.2367	0.3476	0.2313	0.3867	0.5643	0.5967			
lgamma_forecast_accuracy_naive_horizon_period	0.4179	0.0769	0.3675	0.3757	0.3720	0.3858	0.3473	0.4209	0.5200	0.5539			
LTC													
rownames	mean	sd	2022-20	2022-21	2022-22	2022-23	2022-24	2022-25	2022-26	2022-27			
rphn_forecast_accuracy	0.0950	0.0182	0.0736	0.0849	0.0718	0.1054	0.1139	0.1184	0.0859	0.1060			
rphn_forecast_accuracy_naive_final_week	0.0947	0.0183	0.0748	0.0864	0.0691	0.1036	0.1123	0.1197	0.0859	0.1059			
rphn_forecast_accuracy_naive_horizon_period	0.0883	0.0150	0.0762	0.0827	0.0682	0.0951	0.1015	0.1108	0.0745	0.0974			
lgamma_forecast_accuracy	0.1112	0.0181	0.0873	0.1014	0.0899	0.1231	0.1328	0.1317	0.1035	0.1201			
lgamma_forecast_accuracy_naive_final_week	0.1120	0.0181	0.0880	0.1025	0.0906	0.1239	0.1336	0.1324	0.1044	0.1208			
lgamma_forecast_accuracy_naive_horizon_period	0.1056	0.0162	0.0853	0.0949	0.0893	0.1168	0.1247	0.1248	0.0949	0.1144			
DOGE													
rownames	mean	sd	2022-20	2022-21	2022-22	2022-23	2022-24	2022-25	2022-26	2022-27			
rphn_forecast_accuracy	0.2410	0.0202	0.2446	0.2479	0.2635	0.2692	0.2431	0.2137	0.2300	0.2161			
rphn_forecast_accuracy_naive_final_week	0.1928	0.0203	0.2044	0.2070	0.2104	0.2179	0.1903	0.1626	0.1799	0.1697			
rphn_forecast_accuracy_naive_horizon_period	0.1888	0.0199	0.2089	0.2049	0.1995	0.2100	0.1786	0.1548	0.1773	0.1761			
lgamma_forecast_accuracy	0.2093	0.0177	0.2329	0.2231	0.2186	0.2154	0.1935	0.1770	0.2037	0.2106			
lgamma_forecast_accuracy_naive_final_week	0.2050	0.0174	0.2139	0.2153	0.2205	0.2262	0.2012	0.1728	0.1967	0.1934			
lgamma_forecast_accuracy_naive_horizon_period	0.1953	0.0177	0.2127	0.2091	0.2040	0.2125	0.1859	0.1619	0.1904	0.1856			

2064 25.4 Animated Evolution of Distributions

2065 I created gif animations of the empirical probability densities of each blockchain's spent output age distribution over
 2066 time. The green line is the current week's empirical density, with faint echoes of prior weeks in other colors. The
 2067 fitted lgamma distribution line is white and rpln is red. The horizontal axis of the first charts for each blockchain
 2068 has a log scale. Both the horizontal and vertical charts are log scale for the second chart of each blockchain. The
 2069 "1" on the horizontal axis should be interpreted as "0". In other words, a "1" means that the age of the spent output
 2070 is approximately zero. This would occur if (1) a child transaction spent an output from a parent transaction that
 2071 was confirmed in the same block as the child transaction; (2) the timestamps of blocks were out of order, leading
 2072 to a "negative" age; (3) the block was confirmed at a time interval less than half of the target block time, e.g. 5
 2073 minutes for BTC/BCH, leading to rounding down to zero.

2074 There is a clear daily cycle visible in the log-log scale charts.

2075 To view the animated gifs, visit:

2076 <https://rucknium.me/html/spent-output-age-btc-bch-ltc-doge.html>

2077 26 Options for Forecasting

2078 My suggestion is to reduce the dimensionality of the forecasting problem by forecasting the parameters of parametric
 2079 distributions, as in **Section 25.3 Evaluation of Forecast Accuracy**. In this case we would need forecasting
 2080 methods that can handle multivariate forecasting objectives. Beyond that requirement, there are few theoretical
 2081 reasons to favor particular forecasting methods in this setting. Out-of-sample forecasting validation should guide
 2082 the final choice of forecast method. The one exception is forecasting methods that technically require the time
 2083 series to be stationary. We may want to somewhat discount such methods and choose one only if it offers much
 2084 superior performance to a method that allows non-stationarity.

2085 Forecasting methods under consideration include:

- 2086 1. Exponential Smoothing. The `legion` R package can perform multivariate Exponential Smoothing.
- 2087 2. Vector Autoregression (VAR). The `MTS` R package can perform VAR forecasting through its `VARpred` function.
 2088 Note that VAR generally requires stationary time series unless explicit cointegration relationships are specified
 2089 as in a Vector Error Correction model.
- 2090 3. Vector Autoreregressive Moving Average (VARMA). The `MTS` R package can perform VAR forecasting through
 2091 its `VARMAPred` function. Generally, stationarity is required.
- 2092 4. Generalized Autoregressive Conditional Heteroskedasticity (GARCH). The `rmgarch` R package offers multi-
 2093 variate GARCH forecasting. Generally, stationarity is required.
- 2094 5. Kalman Filter. The `MARSS` and `KFAS` R packages provide multivariate Kalman Filter forecasting methods.
- 2095 6. A naive forecast where the forecast is just the value of the parameters of the last period in the training set.
- 2096 7. A naive forecast where the parameters are fit on multiple weeks of pooled data on the last part of the training
 2097 set.

2098 A simple way to evaluate forecast accuracy is to compute the loss function on some final W number of weeks as
 2099 I did in **Section 25.3 Evaluation of Forecast Accuracy**. Another option is to perform the evaluation on a
 2100 rolling window, called "evaluation on a rolling forecasting origin" in Section 5.10 "Time series cross-validation" of
 2101 [Hyndman & Athanasopoulos, 2021]. For now for notational simplicity we will assume the former method.

Let \mathcal{Q} be the set of forecasting methods listed above and $q(\beta)$ be a particular element of the \mathcal{Q} set. $q(\beta)$ is a function that takes as arguments fitted parameter values $\hat{\beta}$ of parametric PDFs $f(x|\beta) \in \mathcal{F}$ of past weekly values and outputs forecasts $\tilde{\beta}$ for future β . Considering the fitting procedure from Equation (29), the final decoy selection algorithm parametric distribution $f(x|\beta)$ and parameters β should be chosen by:

$$\arg \min_{f(x|\beta) \in \mathcal{F}, q(\hat{\beta}) \in \mathcal{Q}} \frac{1}{W} \sum_{i=1}^W L\left(f\left(x_i \mid q\left(\hat{\beta}\right)\right)\right), \forall q \in \mathcal{Q} \quad (30)$$

Expression (30) will find the minimizers $f(x|\tilde{\beta})$ of each loss function $L \in \mathcal{L}$ among all parametric distributions $f(x|\beta) \in \mathcal{F}$ and all forecasting methods $q(\beta) \in \mathcal{Q}$ of the mean for the out-of-sample data x_i for weeks $i = 1, \dots, W$. Once these $f(x|\tilde{\beta})$ are found, the only decision left to be made is the judgment call on the proper choice of loss function and any associated loss function parameters λ, α, η . For example, in Table 7, Expression (30) would choose `rpln.forecast.accuracy.naive.horizon.period` for BTC, LTC, and DOGE, and `rpln.forecast.accuracy` for BCH. Note that Expression (30) forms a plug-in estimator for the risk function in Equation (18). Therefore, it is a type of empirical risk minimizer.³⁹

27 Acknowledgments

This research was funded by Monero's Community Crowdfunding System:

<https://ccs.getmonero.org/proposals/Rucknium-OSPEAD-Fortifying-Monero-Against-Statistical-Attack.html>

The following people gave feedback on the research process and/or contributed in other ways: ACK-J, ArcticMine, bob, coinstudent2048, garth, gingeropolous, isthmus, jberman, kayabaNerve, koe, mj-xmr, monerobull, moneromooo, neptune, plowsof, SamsungGalaxyPlayer, SerHack, SethForPrivacy, and Syksy.

³⁹https://en.wikipedia.org/wiki/Empirical_risk_minimization#Empirical_risk_minimization

2119 References

- 2120 [Aeeneh et al., 2021] Aeeneh, S., Chervinski, J. O., Yu, J., & Zlatanov, N. (2021). New attacks on the untraceability of transactions in cryptonote-style blockchains. *2021 IEEE International Conference on Blockchain and Cryptocurrency (ICBC)*, 1–5. <https://doi.org/10.1109/ICBC51069.2021.9461130>
- 2121
2122
2123 [Allman et al., 2009] Allman, E. S., Matias, C., & Rhodes, J. A. (2009). Identifiability of parameters in latent structure models with many observed variables. *The Annals of Statistics*, 37(6A), 3099 – 3132. <https://doi.org/10.1214/09-AOS689>
- 2124
2125
2126 [Argiento et al., 2020] Argiento, R., Cremaschi, A., & Vannucci, M. (2020). Hierarchical Normalized Completely Random Measures to Cluster Grouped Data. *Journal of the American Statistical Association*, 115(529), 318–333. <https://doi.org/10.1080/01621459.2019.159>
- 2127
2128
2129 [Arias-Castro & Jiang, 2021] Arias-Castro, E. & Jiang, H. (2021). *Extending the patra-sen approach to estimating the background component in a two-component mixture model*. <https://doi.org/10.48550/ARXIV.2106.13925>
- 2130
2131 [Bagui et al., 2006] Bagui, S., Bagui, S., Chatterjee, A., & Mehra, K. (2006). Classification with multiple independent measurements under a separate sampling scheme. *Statistical Methodology*, 3, 234–251. <https://doi.org/10.1016/j.stamet.2005.10.001>
- 2132
2133
2134 [Bassett & Deride, 2019] Bassett, R. & Deride, J. (2019). Maximum a posteriori estimators as a limit of bayes estimators. *Mathematical Programming*, 174(1), 129–144. <https://doi.org/10.1007/s10107-018-1241-0>
- 2135
2136 [Benaglia et al., 2009a] Benaglia, T., Chauveau, D., & Hunter, D. R. (2009a). An em-like algorithm for semi- and nonparametric estimation in multivariate mixtures. *Journal of Computational and Graphical Statistics*, 18(2), 505–526. <https://doi.org/10.1198/jcgs.2009.07175>
- 2137
2138
2139 [Benaglia et al., 2009b] Benaglia, T., Chauveau, D., Hunter, D. R., & Young, D. S. (2009b). mixtools: An r package for analyzing mixture models. *Journal of Statistical Software*, 32(6), 1–29. <https://doi.org/10.18637/jss.v032.i06>
- 2140
2141
2142 [Bergmeir & Benítez, 2012] Bergmeir, C. & Benítez, J. M. (2012). On the use of cross-validation for time series predictor evaluation. *Information Sciences*, 191, 192–213. <https://doi.org/https://doi.org/10.1016/j.ins.2011.12.028>. Data Mining for Software Trustworthiness
- 2143
2144
2145 [Bergmeir et al., 2018] Bergmeir, C., Hyndman, R. J., & Koo, B. (2018). A note on the validity of cross-validation for evaluating autoregressive time series prediction. *Computational Statistics & Data Analysis*, 120, 70–83. <https://doi.org/https://doi.org/10.1016/j.csda.2017.11.003>
- 2146
2147
2148 [Bonhomme et al., 2016] Bonhomme, S., Jochmans, K., & Robin, J.-M. (2016). Non-parametric estimation of finite mixtures from repeated measurements. *Journal of the Royal Statistical Society. Series B (Statistical Methodology)*, 78(1), 211–229. <http://www.jstor.org/stable/24775334>
- 2149
2150
2151 [Casella & Berger, 2002] Casella, G. & Berger, R. L. (2002). *Statistical Inference* (second ed.). Duxbury Pacific Grove, CA.
- 2152
2153 [Cerqueira et al., 2020] Cerqueira, V., Torgo, L., & Mozetič, I. (2020). Evaluating time series forecasting models: an empirical study on performance estimation methods. *Machine Learning*, 109(11), 1997–2028. <https://doi.org/10.1007/s10994-020-05910-7>
- 2154
2155
2156 [Copas, 1974] Copas, J. B. (1974). On Symmetric Compound Decision Rules for Dichotomies. *The Annals of Statistics*, 2(1), 199 – 204. <https://doi.org/10.1214/aos/1176342626>
- 2157
2158 [Creedy, 2015] Creedy, J. (2015). A Note on Computing the Gini Inequality Measure with Weighted Data. Working Paper Series 4235, Victoria University of Wellington, Chair in Public Finance. <https://ideas.repec.org/p/vuw/vuwcpf/4235.html>
- 2159
2160
2161 [Crowder & Hand, 1990] Crowder, M. & Hand, D. (1990). *Analysis of Repeated Measures*. Chapman & Hall/CRC Monographs on Statistics & Applied Probability. Taylor & Francis.
- 2162

- 2163 [Deuber et al., 2022] Deuber, D., Ronge, V., & Rueckert, C. (2022). Sok: Assumptions underlying cryptocurrency
2164 deanonymizations. *Proceedings on Privacy Enhancing Technologies*, 2022(3). <https://petsymposium.org/2022/files/papers/issue3/poops-2022-0091.pdf>
- 2166 [Egger et al., 2022] Egger, C., Lai, R. W. F., Ronge, V., Woo, I. K. Y., & Yin, H. H. F. (2022). On defeating
2167 graph analysis of anonymous transactions. *Proceedings on Privacy Enhancing Technologies*, 2022(3). <https://petsymposium.org/2022/files/papers/issue3/poops-2022-0085.pdf>
- 2169 [Fix & Hodges, 1989] Fix, E. & Hodges, J. L. (1989). Discriminatory analysis. nonparametric discrimination:
2170 Consistency properties. *International Statistical Review / Revue Internationale de Statistique*, 57(3), 238–247.
2171 <http://www.jstor.org/stable/1403797>
- 2172 [Fraley & Raftery, 2002] Fraley, C. & Raftery, A. E. (2002). Model-based clustering, discriminant analysis, and
2173 density estimation. *Journal of the American Statistical Association*, 97(458), 611–631. <http://www.jstor.org/stable/3085676>
- 2175 [Frühwirth-Schnatter et al., 2021] Frühwirth-Schnatter, S., Malsiner-Walli, G., & Grün, B. (2021). Generalized
2176 Mixtures of Finite Mixtures and Telescoping Sampling. *Bayesian Analysis*, 16(4), 1279 – 1307. <https://doi.org/10.1214/21-BA1294>
- 2178 [Hall, 1981] Hall, P. (1981). On the non-parametric estimation of mixture proportions. *Journal of the Royal
2179 Statistical Society: Series B (Methodological)*, 43(2), 147–156. <https://doi.org/https://doi.org/10.1111/j.2517-6161.1981.tb01164.x>
- 2181 [Hsu et al., 2014] Hsu, J., Gaboardi, M., Haeberlen, A., Khanna, S., Narayan, A., Pierce, B. C., & Roth, A. (2014).
2182 Differential privacy: An economic method for choosing epsilon. *2014 IEEE 27th Computer Security Foundations
2183 Symposium*. <https://doi.org/10.1109/csf.2014.35>
- 2184 [Hyndman & Athanasopoulos, 2021] Hyndman, R. J. & Athanasopoulos, G. (2021). *Forecasting: Principles and
2185 Practice* (third ed.). OTexts: Melbourne, Australia. <https://otexts.com/fpp3/>
- 2186 [Kasahara & Shimotsu, 2014] Kasahara, H. & Shimotsu, K. (2014). Non-parametric identification and estimation of
2187 the number of components in multivariate mixtures. *Journal of the Royal Statistical Society. Series B (Statistical
2188 Methodology)*, 76(1), 97–111. <https://www.jstor.org/stable/24772747>
- 2189 [Krawiec-Thayer et al., 2021] Krawiec-Thayer, M. P., Neptune, Rucknium, Jberman, &
2190 Carrington (2021). *Fingerprinting a flood: forensic statistical analysis of the
2191 mid-2021 monero transaction volume anomaly*. <https://mitchellpk.medium.com/fingerprinting-a-flood-forensic-statistical-analysis-of-the-mid-2021-monero-transaction-volume-a19cbf41ce60>
Available at <https://mitchellpk.medium.com/fingerprinting-a-flood-forensic-statistical-analysis-of-the-mid-2021-monero-transaction-volume-a19cbf41ce60>
- 2195 [Kumar et al., 2017] Kumar, A., Fischer, C., Tople, S., & Saxena, P. (2017). A traceability analysis of monero's
2196 blockchain. *European Symposium on Research in Computer Security (ESORICS)*. https://doi.org/10.1007/978-3-319-66399-9_9
- 2198 [Kwon & Mbakop, 2021] Kwon, C. & Mbakop, E. (2021). Estimation of the number of components of nonparametric
2199 multivariate finite mixture models. *The Annals of Statistics*, 49(4), 2178 – 2205. <https://doi.org/10.1214/20-AOS2032>
- 2201 [Levine et al., 2011] Levine, M., Hunter, D. R., & Chauveau, D. (2011). Maximum smoothed likelihood for multi-
2202 variate mixtures. *Biometrika*, 98(2), 403–416. <http://www.jstor.org/stable/23076159>
- 2203 [Lilienfeld & Landfield, 2008] Lilienfeld, S. O. & Landfield, K. (2008). Science and pseudoscience in law enforcement:
2204 A user-friendly primer. *Criminal Justice and Behavior*, 35(10), 1215–1230. <https://doi.org/10.1177/0093854808321526>
- 2206 [Liu et al., 2019] Liu, C., He, X., Chanyaswad, T., Wang, S., & Mittal, P. (2019). Investigating statistical privacy
2207 frameworks from the perspective of hypothesis testing. *Proceedings on Privacy Enhancing Technologies*, 2019(3),
2208 233–254. <https://doi.org/10.2478/poops-2019-0045>

- 2209 [Mackenzie et al., 2015] Mackenzie, A., Noether, S., & Team, M. C. (2015). *Improving obfuscation in the cryptonote*
2210 *protocol*. Research Bulletin. <https://www.getmonero.org/resources/research-lab/pubs/MRL-0004.pdf>
- 2211 [Maji et al., 2019] Maji, A., Ghosh, A., Basu, A., & Pardo, L. (2019). Robust statistical inference based on the
2212 c-divergence family. *Annals of the Institute of Statistical Mathematics*, 71(5), 1289–1322. <https://doi.org/10.1007/s10463-018-0678-5>
- 2214 [Makarov & Schoar, 2021] Makarov, I. & Schoar, A. (2021). Blockchain analysis of the bitcoin market. Working
2215 Paper 29396, National Bureau of Economic Research. <https://doi.org/10.3386/w29396>
- 2216 [Makridakis et al., 2022] Makridakis, S., Spiliotis, E., & Assimakopoulos, V. (2022). M5 accuracy competition:
2217 Results, findings, and conclusions. *International Journal of Forecasting*. <https://doi.org/https://doi.org/10.1016/j.ijforecast.2021.11.013>
- 2219 [Malsiner-Walli et al., 2017] Malsiner-Walli, G., Frühwirth-Schnatter, S., & Grün, B. (2017). Identifying mixtures
2220 of mixtures using bayesian estimation. *Journal of Computational and Graphical Statistics*, 26(2), 285–295. <https://www.tandfonline.com/doi/full/10.1080/10618600.2016.1200472>
- 2222 [Mbakop, 2017] Mbakop, E. D. (2017). *Identification of auctions with incomplete bid data in the presence of unob-*
2223 *served heterogeneity*. Working Papers, Univ. Calgary. <https://cpb-us-e1.wpmucdn.com/sites.northwestern.edu/dist/5/1999/files/2017/10/JMP-Brouillon-1bom42o.pdf>
- 2225 [McLachlan, 1992] McLachlan, G. (1992). *Discriminant Analysis and Statistical Pattern Recognition*. Wiley Series
2226 in Probability and Statistics. Wiley.
- 2227 [McLachlan & Peel, 2000] McLachlan, G. & Peel, D. (2000). *Finite Mixture Models*. Wiley Series in Probability
2228 and Statistics. Wiley.
- 2229 [McLachlan et al., 2019] McLachlan, G. J., Lee, S. X., & Rathnayake, S. I. (2019). Finite mixture
2230 models. *Annual Review of Statistics and Its Application*, 6(1), 355–378. <https://doi.org/10.1146/annurev-statistics-031017-100325>
- 2232 [Milhaud et al., 2021] Milhaud, X., Pommeret, D., Salhi, Y., & Vandekerkhove, P. (2021). Shape constraint free
2233 two-sample contamination model testing. working paper or preprint. <https://hal.archives-ouvertes.fr/hal-03201760>. working paper or preprint
- 2235 [Möser et al., 2018] Möser, M., Soska, K., Heilman, E., Lee, K., Heffan, H., Srivastava, S., Hogan, K., Hen-
2236 nnessey, J., Miller, A., Narayanan, A., & Christin, N. (2018). An empirical analysis of traceability in the mon-
2237 ero blockchain. *Proceedings on Privacy Enhancing Technologies*, 2018(3), 143–163. <https://doi.org/doi:10.1515/popets-2018-0025>
- 2239 [Murphy & Topel, 1985] Murphy, K. M. & Topel, R. H. (1985). Estimation and inference in two-step econometric
2240 models. *Journal of Business & Economic Statistics*, 3(4), 370–379. <https://doi.org/10.1080/07350015.1985.10509471>
- 2242 [Newton, 1996] Newton, M. A. (1996). Bootstrapping phylogenies: Large deviations and dispersion effects.
2243 *Biometrika*, 83(2), 315–328. <https://www.jstor.org/stable/2337603>
- 2244 [Ni et al., 2021] Ni, W., Cheng, P., Chen, L., & Lin, X. (2021). When the recursive diversity anonymity meets the
2245 ring signature. *Proceedings of the 2021 International Conference on Management of Data*, 1359?1371. Association
2246 for Computing Machinery. <https://doi.org/10.1145/3448016.3452825>
- 2247 [O'Hagan et al., 2012] O'Hagan, A., Murphy, T. B., & Gormley, I. C. (2012). Computational aspects of fitting
2248 mixture models via the expectation–maximization algorithm. *Computational Statistics & Data Analysis*, 56(12),
2249 3843–3864. <https://www.sciencedirect.com/science/article/pii/S0167947312002101>
- 2250 [Patra & Sen, 2016] Patra, R. K. & Sen, B. (2016). Estimation of a two-component mixture model with applications
2251 to multiple testing. *Journal of the Royal Statistical Society. Series B (Statistical Methodology)*, 78(4), 869–893.
2252 <http://www.jstor.org/stable/24775367>
- 2253 [Reed & Jorgensen, 2004] Reed, W. J. & Jorgensen, M. (2004). The double pareto-lognormal distribution—a new
2254 parametric model for size distributions. *Communications in Statistics - Theory and Methods*, 33(8), 1733–1753.
2255 <https://doi.org/10.1081/STA-120037438>

- 2256 [Ritchie et al., 2020] Ritchie, A., Vandermeulen, R. A., & Scott, C. (2020). Consistent estimation of identifiable non-
2257 parametric mixture models from grouped observations. *Advances in Neural Information Processing Systems*, 33,
2258 11676–11686. <https://papers.nips.cc/paper/2020/hash/866d90e0921ac7b024b47d672445a086-Abstract.html>
- 2259
2260 [Ronge et al., 2021] Ronge, V., Egger, C., Lai, R. W. F., Schröder, D., & Yin, H. H. F. (2021). Foundations
2261 of ring sampling. *Proceedings on Privacy Enhancing Technologies*, 2021(3), 265–288. <https://doi.org/doi:10.2478/popets-2021-0047>
- 2262
2263 [Vandermeulen & Scott, 2019] Vandermeulen, R. A. & Scott, C. D. (2019). An operator theoretic approach to
2264 nonparametric mixture models. *The Annals of Statistics*, 47(5), 2704 – 2733. <https://doi.org/10.1214/18-AOS1762>
- 2265
2266 [Vankadara et al., 2021] Vankadara, L. C., Bordt, S., von Luxburg, U., & Ghoshdastidar, D. (2021). Recovery guarantees
2267 for kernel-based clustering under non-parametric mixture models. *Proceedings of The 24th International Conference on Artificial Intelligence and Statistics*, volume 130 of *Proceedings of Machine Learning Research*,
2268 3817–3825. <https://proceedings.mlr.press/v130/vankadara21a.html>
- 2269
2270 [Vijayakumaran, 2021] Vijayakumaran, S. (2021). *Analysis of cryptonote transaction graphs using the dulmage-mendelsohn decomposition*. Cryptology ePrint Archive, Report 2021/760. <https://eprint.iacr.org/2021/760>
- 2271
2272 [Wei & Nguyen, 2022] Wei, Y. & Nguyen, X. (2022). Convergence of de Finetti's mixing measure in latent structure
2273 models for observed exchangeable sequences. *The Annals of Statistics*, 50(4), 1859 – 1889. <https://doi.org/10.1214/21-AOS2120>
- 2274
2275 [Ye et al., 2020] Ye, C., Ojukwu, C., Hsu, A., & Hu, R. (2020). *Alt-coin traceability*. Cryptology ePrint Archive,
2276 Report 2020/593. <https://eprint.iacr.org/2020/593>
- 2277
2278 [Yu et al., 2019] Yu, J., Au, M. H. A., & Esteves-Verissimo, P. (2019). Re-thinking untraceability in the cryptonote-style
2279 blockchain. *2019 IEEE 32nd Computer Security Foundations Symposium (CSF)*, 94–9413. <https://doi.org/10.1109/CSF.2019.00014>