

# Randomized Estimation Algorithms

This version of the document is dated 2023-02-09.

Peter Occil

## 1 Introduction

Suppose we have an endless stream of numbers, each generated at random and independently from each other, and we can sample as many numbers from the stream as we want. These numbers are called *random variates*. This page presents general-purpose algorithms for estimating the mean value ("long-run average") of those variates, or estimating the mean value of a function of those numbers. The estimates are either *unbiased* (they have no systematic bias from the true mean value), or they come close to the true value with a user-specified error tolerance.

The algorithms are described to make them easy to implement by programmers.

Not yet covered are the following algorithms:

- Unbiased mean estimation algorithms that take a sequence of estimators that get better and better at estimating the desired mean (for example, estimators that average an increasing number of sample points). See, for example, Vihola (2018)<sup>1</sup>.

### 1.1 About This Document

**This is an open-source document; for an updated version, see the [source code](#) or its [rendering on GitHub](#). You can send comments on this document on the [GitHub issues page](#)\*\*.**

My audience for this article is **computer programmers with mathematics knowledge, but little or no familiarity with calculus.**

I encourage readers to implement any of the algorithms given in this page, and report their implementation experiences. In particular, **I seek comments on the following aspects:**

- Are the algorithms in this article easy to implement? Is each algorithm written so that someone could write code for that algorithm after reading the article?
- Does this article have errors that should be corrected?
- Are there ways to make this article more useful to the target audience?

Comments on other aspects of this document are welcome.

## 2 Concepts

The following concepts are used in this document.

Each algorithm takes a stream of independent random variates (numbers). These variates follow a *probability distribution* or simply *distribution*, or a rule that says which kinds of numbers have greater probability of occurring than others. A distribution has the following properties.

- The *expectation*, *expected value*, or *mean* is the "long-run average" value of the distribution. It is expressed as  $\mathbf{E}[X]$ , where  $X$  is a number taken from the stream. If  $X$  is non-negative and if we take independent random samples and then take their average, then with probability 1, the average will approach the expected value as the number of samples gets large (the *law of large numbers*).
- An  $n^{\text{th}}$  *moment* is the expected value of  $X^n$ . In other words, take random samples, raise them to the power  $n$ , then take their average. Then with probability 1, the average approaches the  $n^{\text{th}}$  moment as  $n$  gets large.
- An  $n^{\text{th}}$  *central moment (about the mean)* is the expected value of  $(X - \mu)^n$ , where  $\mu$  is the distribution's mean. The 2nd central moment is called *variance*.
- An  $n^{\text{th}}$  *central absolute moment (c.a.m.)* is the expected value of

$\text{abs}(X - \mu)^n$ , where  $\mu$  is the distribution's mean. This is the same as the central moment when  $n$  is even.

Some distributions don't have an  $n^{\text{th}}$  moment for a particular  $n$ . This usually means the  $n^{\text{th}}$  power of the stream's numbers varies so wildly that it can't be estimated accurately. If a distribution has an  $n^{\text{th}}$  moment, it also has a  $k^{\text{th}}$  moment for every  $k$  in the interval  $[1, n)$ .

The *relative error* of an estimation algorithm is  $\text{abs}(\text{est}/\text{trueval}) - 1$ , where *est* is the estimate and *trueval* is the true expected value.

### 3 A Relative-Error Algorithm for a Bernoulli Stream

The following algorithm from Huber (2017)<sup>2</sup> estimates the probability that a stream of random zeros and ones produces the number 1. The algorithm's relative error is independent of that probability, however, and the algorithm produces *unbiased* estimates. Specifically, the stream of numbers has the following properties:

- The stream produces only zeros and ones (that is, the stream follows the **Bernoulli distribution**).
- The stream of numbers can't take on the value 0 with probability 1.
- The stream's mean (expected value) is unknown.

The algorithm, also known as *Gamma Bernoulli Approximation Scheme*, has the following parameters:

- $\varepsilon, \delta$ : Both parameters must be greater than 0, and  $\varepsilon$  must be 3/4 or less, and  $\delta$  must be less than 1.

With this algorithm, the relative error will be no greater than  $\varepsilon$  with probability  $1 - \delta$  or greater. However, the estimate can be higher than 1 with probability greater than 0.

The algorithm, called **Algorithm A** in this document, follows.

1. Calculate the minimum number of samples  $k$ . There are two suggestions. The simpler one is  $k = \text{ceil}(-6 \cdot \ln(2/\delta) / (\varepsilon^2 \cdot (4 \cdot \varepsilon - 3)))$ . A more complicated one is the smallest integer  $k$  such that  $\text{gammainc}(k, (k-1)/(1+\varepsilon)) + (1 - \text{gammainc}(k, (k-1)/(1-\varepsilon))) \leq \delta$ ,

where  $\text{gammaln}$  is the regularized lower incomplete gamma function.

2. Take samples from the stream until  $k$  1's are taken this way. Let  $r$  be the total number of samples taken this way.
3. Generate  $g$ , a gamma random variate with shape parameter  $r$  and scale 1, then return  $(k-1)/g$ .

**Notes:**

1. As noted in Huber 2017, if we have a stream of random variates that take on values in the interval  $[0, 1]$ , but have unknown mean, we can transform each number by—
  1. generating a  $\text{uniform}(0, 1)$  random variate  $u$ , then
  2. changing that number to 1 if  $u$  is less than that number, or 0 otherwise,

and we can use the new stream of zeros and ones in the algorithm to get an unbiased estimate of the unknown mean.

2. As can be seen in Feng et al. (2016)<sup>3</sup>, the following is equivalent to steps 2 and 3 of *Algorithm A*: "Let  $G$  be 0. Do this  $k$  times: 'Flip a coin until it shows heads, let  $r$  be the number of flips (including the last), generate a gamma random variate with shape parameter  $r$  and scale 1, and add that variate to  $G$ .' The estimated probability of heads is then  $(k-1)/G$ .", and the following is likewise equivalent if the stream of random variates follows a (zero-truncated) "geometric" distribution with unknown mean: "Let  $G$  be 0. Do this  $k$  times: 'Take a sample from the stream, call it  $r$ , generate a gamma random variate with shape parameter  $r$  and scale 1, and add that variate to  $G$ .' The estimated mean is then  $(k-1)/G$ ." (This is with the understanding that the geometric distribution is defined differently in different academic works.) The geometric algorithm produces unbiased estimates just like *Algorithm A*.
3. The generation of a gamma random variate and the division by that variate can cause numerical errors in practice, such as rounding and cancellations, unless care is taken.
4. Huber proposes another algorithm that claims to be faster when the mean is bounded away from zero; see (Huber 2022)<sup>4</sup>.

## 4 A Relative-Error Algorithm for a Bounded Stream

The following algorithm comes from Huber and Jones (2019)<sup>5</sup>; see also Huber (2017)<sup>6</sup>. It estimates the expected value of a stream of random variates with the following properties:

- The numbers in the stream lie in the closed interval  $[0, 1]$ .
- The stream of numbers can't take on the value 0 with probability 1.
- The stream's mean (expected value) is unknown.

The algorithm has the following parameters:

- $\varepsilon, \delta$ : Both parameters must be greater than 0, and  $\varepsilon$  must be  $1/8$  or less, and  $\delta$  must be less than 1.

With this algorithm, the relative error will be no greater than  $\varepsilon$  with probability  $1 - \delta$  or greater. However, the estimate has a nonzero probability of being higher than 1.

This algorithm is not guaranteed to produce unbiased estimates.

The algorithm, called **Algorithm B** in this document, follows.

1. Set  $k$  to  $\text{ceil}(2 \cdot \ln(6/\delta)/\varepsilon^{2/3})$ .
2. Set  $b$  to 0 and  $n$  to 0.
3. (Stage 1: Modified gamma Bernoulli approximation scheme.) While  $b$  is less than  $k$ :
  1. Add 1 to  $n$ .
  2. Take a sample from the stream, call it  $s$ .
  3. With probability  $s$  (for example, if a newly generated  $\text{uniform}(0, 1)$  random variate is less than  $s$ ), add 1 to  $b$ .
4. Set  $gb$  to  $k + 2$ , then generate a gamma random variate with shape parameter  $n$  and scale 1, then divide  $gb$  by that variate.
5. (Find the sample size for the next stage.) Set  $c1$  to  $2 \cdot \ln(3/\delta)$ .
6. Generate a Poisson random variate with mean  $c1/(\varepsilon \cdot gb)$ , call it  $n$ .
7. Run the standard deviation sub-algorithm (given later)  $n$  times. Set  $A$  to the number of 1's returned by that sub-algorithm this way.
8. Set  $csquared$  to  $(A / c1 + 1 / 2 + \text{sqrt}(A / c1 + 1 / 4)) * (1 + \varepsilon^{1/3})^2 * \varepsilon / gb$ .
9. Set  $n$  to  $\text{ceil}((2 \cdot \ln(6/\delta)/\varepsilon^2)/(1 - \varepsilon^{1/3}))$ , or an integer greater than this.

10. (Stage 2: Light-tailed sample average.) Set  $e0$  to  $\varepsilon^{1/3}$ .
11. Set  $\mu0$  to  $gb/(1-e0^2)$ .
12. Set  $\alpha$  to  $\varepsilon/(csquared*\mu0)$ .
13. Set  $w$  to  $n*\mu0$ .
14. Do the following  $n$  times:
  1. Get a sample from the stream, call it  $g$ . Set  $s$  to  $\alpha*(g-\mu0)$ .
  2. If  $s \geq 0$ , add  $\ln(1+s+s*s/2)/\alpha$  to  $w$ . Otherwise, subtract  $\ln(1-s+s*s/2)/\alpha$  from  $w$ .
15. Return  $w/n$ .

The standard deviation sub-algorithm follows.

1. Generate an unbiased random bit. If that bit is 1 (which happens with probability 1/2), return 0.
2. Get two samples from the stream, call them  $x$  and  $y$ .
3. Generate a uniform(0, 1) random variate, call it  $u$ .
4. If  $u$  is less than  $(x-y)^2$ , return 1. Otherwise, return 0.

#### Notes:

1. As noted in Huber and Jones, if the stream of random variates takes on values in the interval  $[0, m]$ , where  $m$  is a known number, we can divide the stream's numbers by  $m$  before using them in *Algorithm B*, and the algorithm will still work.
2. While this algorithm is exact in theory (assuming computers can store real numbers of any precision), practical implementations of it can cause numerical errors, such as rounding and cancellations, unless care is taken.

## 5 An Absolute-Error Adaptive Algorithm

The following algorithm comes from Kunsch et al. (2019)<sup>7</sup>. It estimates the mean of a stream of random variates with the following properties:

- The distribution of numbers in the stream has a finite  $q^{\text{th}}$  c.a.m. and  $p^{\text{th}}$  c.a.m.
- The exact  $q^{\text{th}}$  c.a.m. and  $p^{\text{th}}$  c.a.m. need not be known in advance.
- The  $q^{\text{th}}$  c.a.m.'s  $q^{\text{th}}$  root divided by the  $p^{\text{th}}$  c.a.m.'s  $p^{\text{th}}$  root is no more than  $\kappa$ , where  $\kappa$  is 1 or greater. (The  $q^{\text{th}}$  c.a.m.'s  $q^{\text{th}}$  root is

also known as *standard deviation* if  $q = 2$ , and *mean absolute deviation* if  $q = 1$ ; similarly for  $p$ .)

The algorithm works by first estimating the  $p^{\text{th}}$  c.a.m. of the stream, then using the estimate to determine a sample size for the next step, which actually estimates the stream's mean.

This algorithm is not guaranteed to produce unbiased estimates.

The algorithm has the following parameters:

- $\varepsilon, \delta$ : Both parameters must be greater than 0, and  $\delta$  must be less than 1. The algorithm will return an estimate within  $\varepsilon$  of the true expected value with probability  $1 - \delta$  or greater, and the estimate will not go beyond the bounds of the stream's numbers. The algorithm is not guaranteed to maintain a finite mean squared error or expected error in its estimates.
- $p$ : The degree of the  $p^{\text{th}}$  c.a.m. that the algorithm will estimate to determine the mean.
- $q$ : The degree of the  $q^{\text{th}}$  c.a.m.  $q$  must be greater than  $p$ .
- $\kappa$ : Maximum value allowed for the following value: the  $q^{\text{th}}$  c.a.m.'s  $q^{\text{th}}$  root divided by the  $p^{\text{th}}$  c.a.m.'s  $p^{\text{th}}$  root. (If  $p = 2$  and  $q = 4$ , this is the maximum value allowed for the kurtosis's 4th root (Hickernell et al. 2012)<sup>8 9</sup>.)  $\kappa$  may not be less than 1.

Both  $p$  and  $q$  must be 1 or greater and are usually integers.

For example:

- With parameters  $p = 2, q = 4, \varepsilon = 1/10, \delta = 1/16, \kappa = 1.1$ , the algorithm assumes the stream's numbers are distributed so that the kurtosis's 4th root, that is, the 4th c.a.m.'s 4th root ( $q=4$ ) divided by the standard deviation ( $p=2$ ), is no more than 1.1 (or alternatively, the kurtosis is no more than  $1.1^4 = 1.4641$ ), and will return an estimate that's within  $1/10$  of the true mean with probability at least  $(1 - 1/16)$  or  $15/16$ .
- With parameters  $p = 1, q = 2, \varepsilon = 1/10, \delta = 1/16, \kappa = 2$ , the algorithm assumes the stream's numbers are distributed so that the standard deviation ( $q=2$ ) divided by the mean deviation ( $p=1$ ) is no more than 2, and will return an estimate that's within  $1/10$  of the true mean with probability at least  $(1 - 1/16)$  or  $15/16$ .

The algorithm, called **Algorithm C** in this document, follows.

1. If  $\kappa$  is 1:
  1. Set  $n$  to  $\text{ceil}(\ln(1/\delta)/\ln(2))+1$  (or an integer greater than this).
  2. Get  $n$  samples from the stream and return  $(mn + mx)/2$ , where  $mx$  is the highest sample and  $mn$  is the lowest.
2. Set  $k$  to  $\text{ceil}((2*\ln(1/\delta))/\ln(4/3))$ . If  $k$  is even<sup>10</sup>, add 1 to  $k$ .
3. Set  $kp$  to  $k$ .
4. Set  $\kappa$  to  $\kappa^{(p*q/(q-p))}$ .
5. If  $q$  is 2 or less:
  - Set  $m$  to  $\text{ceil}(3*\kappa*48^{1/(q-1)})$  (or an integer greater than this); set  $s$  to  $1+1/(q-1)$ ; set  $h$  to  $16^{1/(q-1)*\kappa/\epsilon^s}$ .
6. If  $q$  is greater than 2:
  - Set  $m$  to  $\text{ceil}(144*\kappa)$ ; set  $s$  to 2; set  $h$  to  $16*\kappa/\epsilon^s$ .
7. (Stage 1: Estimate  $p^{\text{th}}$  c.a.m. to determine number of samples for stage 2.) Create  $k$  many blocks. For each block:
  1. Get  $m$  samples from the stream.
  2. Add the samples and divide by  $m$  to get this block's sample mean, *mean*.
  3. Calculate the estimate of the  $p^{\text{th}}$  c.a.m. for this block, which is:  $((\text{block}[0] - \text{mean})^p + \text{block}[1] - \text{mean})^p + \dots + \text{block}[k-1] - \text{mean})^p)/m$ , where  $\text{block}[i]$  is the sample at position  $i$  of the block (positions start at 0).
8. (Find the median of the  $p^{\text{th}}$  c.a.m. estimates.) Sort the estimates calculated by step 7 in ascending order, and set *median* to the value in the middle of the sorted list (at position  $\text{floor}(k/2)$  with positions starting at 0); this works because  $k$  is odd.
9. (Calculate sample size for the next stage.) Set  $mp$  to  $\max(1, \text{ceil}(h * \text{median}^s))$ , or an integer greater than this.
10. (Stage 2: Estimate of the sample mean.) Create  $kp$  many blocks. For each block:
  1. Get  $mp$  samples from the stream.
  2. Add the samples and divide by  $mp$  to get this block's sample mean.
11. (Find the median of the sample means. This is definitely an unbiased estimate of the mean when  $kp$  is 1 or 2, but unfortunately, it isn't one for any  $kp > 2$ .) Sort the sample means from step 10 in ascending order, and return the value in the middle of the sorted list (at position  $\text{floor}(kp/2)$  with positions starting at 0); this works because  $kp$  is odd.

**Notes:**



1. The interval  $[\hat{\mu} - \epsilon, \hat{\mu} + \epsilon]$  is also known as a *confidence interval* for the mean, with *confidence level* at least  $1 - \delta$  (where  $\hat{\mu}$  is an estimate of the mean returned by *Algorithm C*).
2. If the stream of random variates meets the condition for *Algorithm C* for a given  $q$ ,  $p$ , and  $\kappa$ , then it still meets that condition when those variates are multiplied by a constant or a constant is added to them.
3. Theorem 3.4 of Kunsch et al. (2019)<sup>11</sup> shows that there is no mean estimation algorithm that—
  - produces an estimate within a user-specified error tolerance (in terms of *absolute error*, as opposed to *relative error*) with probability greater than a user-specified value, and
  - works for all streams whose distribution is known only to have finite moments (the moments are bounded but the bounds are unknown).
4. There is also a mean estimation algorithm for very high dimensions, which works if the stream of multidimensional variates has a finite variance (Lee and Valiant 2022)<sup>12</sup>, but this algorithm is impractical — it requires millions of samples at best.

#### Examples:

1. To estimate the probability of heads of a coin that produces either 1 with an unknown probability in the interval  $[\mu, 1 - \mu]$ , or 0 otherwise, we can take  $q = 4$ ,  $p = 2$ , and  $\kappa \geq (1/\min(\mu, 1 - \mu))^{1/4}$  (Kunsch et al. 2019, Lemma 3.6).
2. The kurtosis of a Poisson distribution with mean  $\mu$  is  $(3 + 1/\mu)$ . Thus, for example, to estimate the mean of a stream of Poisson variates with mean  $\nu$  or greater but otherwise unknown, we can take  $q = 4$ ,  $p = 2$ , and  $\kappa \geq (3 + 1/\nu)^{1/4}$ .
3. The kurtosis of an exponential distribution is 9 regardless of its rate. Thus, to estimate the mean of a stream of exponential variates with unknown mean, we can take  $q = 4$ ,  $p = 2$ , and  $\kappa \geq 9^{1/4} = \sqrt[4]{9}$ .

## 6 Estimating the Mode

Suppose there is an endless stream of items, each generated at random and independently from each other, and we can sample as many items from the stream as we want. Then the following algorithm estimates the most frequently occurring item, called the *mode*. (Dutta and Goswami 2010)<sup>13</sup> This assumes the following are known:

- Exactly one item must occur more frequently than the others.
- $\epsilon$  is greater than 0 and less than one half of the smallest possible difference between the mode's probability and the next most frequent item's probability.
- $\delta$  is greater than 0 and less than 1.
- $n$  is the number of distinct items that can be taken.

The following algorithm correctly estimates the mode with probability  $1 - \delta$ .

1. Calculate  $m = \text{ceil}(\frac{(4\epsilon + 3)(\ln(\frac{n}{\delta}) + \ln(2))}{6\epsilon^2})$ .
2. Take  $m$  items from the stream. If one item occurs more frequently than any other item taken this way, return the most frequent item. Otherwise, return an arbitrary but fixed item (among the items the stream can take).

## 7 Estimating a Function of the Mean

*Algorithm C* can be used to estimate a function of the mean of a stream of random variates with unknown mean. Specifically, the goal is to estimate  $f(\mathbf{E}[\mathbf{z}])$ , where:

- $\mathbf{z}$  is a number produced by the stream. Each number produced by the stream must lie in the interval  $[0, 1]$ .
- $f$  is a known continuous function that maps the closed interval  $[0, 1]$  to  $[0, 1]$ .
- The stream's numbers can take on a single value with probability 1.

The following algorithm takes the following parameters:

- $p$ ,  $q$ , and  $\kappa$  are as defined in *Algorithm C*.
- $\epsilon$ ,  $\delta$ : The algorithm will return an estimate within  $\epsilon$  of  $f(\mathbf{E}[\mathbf{z}])$  with probability  $1 - \delta$  or greater, and the estimate will be in the interval  $[0, 1]$ .

The algorithm, like *Algorithm C*, works only if the stream's distribution has the following technical property: The  $q^{\text{th}}$  c.a.m.'s  $q^{\text{th}}$  root divided by the  $p^{\text{th}}$  c.a.m.'s  $p^{\text{th}}$  root is no more than  $\kappa$ , where  $\kappa$  is 1 or greater. The algorithm, called **Algorithm D** in this document, follows.

1. Calculate  $\gamma$  as a number equal to or less than  $\psi(\varepsilon)$ , or the *inverse modulus of continuity*, which is found by taking the so-called *modulus of continuity* of  $f(x)$ , call it  $\omega(h)$ , and solving the equation  $\omega(h) = \varepsilon$  for  $h$ .
  - Loosely speaking, a modulus of continuity  $\omega(h)$  gives the maximum range of  $f$  in a window of size  $h$ .
  - For example, if the slope of  $f$  is no "steeper" than that of the function  $M \cdot h$ , then  $f$  is *Lipschitz continuous* with Lipschitz constant  $M$ , so that its modulus of continuity is  $\omega(h) = M \cdot h$ . The solution for  $\psi$  is then  $\psi(\varepsilon) = \varepsilon/M$ .
  - Because  $f$  is continuous on a closed interval, it's guaranteed to have a modulus of continuity (by the Heine–Cantor theorem; see also a [related question](#)).
2. Run *Algorithm C* with the given parameters  $p$ ,  $q$ ,  $\kappa$ , and  $\delta$ , but with  $\varepsilon = \gamma$ . Let  $\mu$  be the result.
3. Return  $f(\mu)$ .

A simpler version of *Algorithm D* was given as an answer to the linked-to question; see also Jiang and Hickernell (2014)<sup>14</sup>. As with *Algorithm D*, this algorithm will return an estimate within  $\varepsilon$  of  $f(\mathbf{E}[\mathbf{z}])$  with probability  $1 - \delta$  or greater, and the estimate will be in the interval  $[0, 1]$ . The algorithm, called **Algorithm E** in this document, follows.

1. Calculate  $\gamma$  as given in step 1 of *Algorithm D*.
2. (Calculate the sample size.) Set  $n$  to  $\text{ceil}(\ln(2/\delta)/(2\gamma^2))$ . (As the answer notes, this sample size is based on Hoeffding's inequality.)
3. (Calculate the sample mean.) Get  $n$  samples from the stream, sum them, then divide the sum by  $n$ , then call the result  $\mu$ . Return  $f(\mu)$ .

If the stream is **unbounded** (can take on any real number) and its distribution has a **known upper bound on the standard deviation  $\sigma$  (or the variance  $\sigma^2$ )**, then a similar algorithm follows from Chebyshev's inequality. This was mentioned as Equation 14 in Hickernell et al. (2012/2013)<sup>15</sup>, but is adapted to find the mean for  $f(x)$ , which must be bounded and continuous on every closed interval

of the real line. The algorithm will return an estimate within  $\varepsilon$  of  $f(\mathbf{E}[\mathbf{z}])$  with probability  $1 - \delta$  or greater, and the estimate will not go beyond the bounds of the stream's numbers. The algorithm, called **Algorithm F** in this document, follows.

1. Calculate  $\gamma$  as given in step 1 of *Algorithm D*.
2. (Calculate the sample size.) Set  $n$  to  $\text{ceil}(\sigma^2/(\delta*\gamma^2))$ .
3. (Calculate the sample mean.) Get  $n$  samples from the stream, sum them, then divide the sum by  $n$ , then call the result  $\mu$ . Return  $f(\mu)$ .

#### Notes:

1. *Algorithm D* and *Algorithm E* won't work in general when  $f(x)$  has jump discontinuities (this happens in general when  $f$  is piecewise continuous, or made up of independent continuous pieces that cover all of  $[0, 1]$ ), at least when  $\varepsilon$  is equal to or less than the maximum jump among all the jump discontinuities (see also a [related question](#)).
2. *Algorithm D* and *Algorithm E* can be adapted to apply to streams outputting numbers in a bounded interval  $[a, b]$  (where  $a$  and  $b$  are known rational numbers), but with unknown mean, and with  $f$  being a continuous function that maps  $[a, b]$  to  $[a, b]$ , as follows:
  - For each number in the stream, subtract  $a$  from it, then divide it by  $(b - a)$ .
  - Instead of  $\varepsilon$ , take  $\varepsilon/(b - a)$ .
  - If the algorithm would return  $f(\mu)$ , instead return  $g(\mu)$  where  $g(\mu) = f(a + (\mu*(b - a)))$ .
3. *Algorithm E* and *Algorithm F* are not unbiased estimators in general. However, when  $f(x) = x$ , the sample mean used by both algorithms is an unbiased estimator of the mean as long as the sample size  $n$  is unchanged.

#### Examples:

1. Take  $f(x) = \sin(\pi*x^4)/2 + 1/2$ . This is a Lipschitz continuous function with Lipschitz constant  $2*\pi$ , so for this  $f$ ,  $\psi(\varepsilon) = \varepsilon/(2*\pi)$ . Now, if we have a coin that produces heads with an unknown probability in the interval  $[\mu, 1-\mu]$ , or 0 otherwise, we can run *Algorithm D* or *Algorithm E* with  $q = 4$ ,  $p = 2$ , and  $\kappa \geq (1/\min(\mu, 1-\mu))^{1/4}$  (see the section on *Algorithm C*).

2. Take  $f(x) = x$ . This is a Lipschitz continuous function with Lipschitz constant 1, so for this  $f$ ,  $\psi(\varepsilon) = \varepsilon/1$ .
3. The variance of a Poisson distribution with mean  $\mu$  is  $\mu$ . Thus, for example, to estimate the mean of a stream of Poisson variates with mean  $\nu$  or less but otherwise unknown, we can take  $\sigma = \sqrt{\nu}$  so that the sample size  $n$  is  $\text{ceil}(\sigma^2/(\delta^2\varepsilon^2))$ , in accordance with *Algorithm F*.

## 8 Randomized Integration

Monte Carlo integration is a randomized way to estimate the integral ("area under the graph") of a function.

This time, suppose we have an endless stream of *vectors* ( $n$ -dimensional points), each generated at random and independently from each other, and we can sample as many vectors from the stream as we want.

*Algorithm C* can be used to estimate an integral of a function  $h(\mathbf{z})$ , where  $\mathbf{z}$  is a vector from the stream, with the following properties:

- $h(\mathbf{z})$  is a multidimensional function that takes an  $n$ -dimensional vector and returns a real number.  $h(\mathbf{z})$  is usually a function that's easy to evaluate but whose integral is hard to calculate.
- $\mathbf{z}$  is an  $n$ -dimensional vector chosen at random in the sampling domain.

The estimate will come within  $\varepsilon$  of the true integral with probability  $1 - \delta$  or greater, as long as the following conditions are met:

- The  $q^{\text{th}}$  c.a.m. for  $h(\mathbf{z})$  is finite. That is,  $\mathbf{E}[\text{abs}(h(\mathbf{z}) - \mathbf{E}[h(\mathbf{z})])^q]$  is finite.
- The  $q^{\text{th}}$  c.a.m.'s  $q^{\text{th}}$  root divided by the  $p^{\text{th}}$  c.a.m.'s  $p^{\text{th}}$  root is no more than  $\kappa$ , where  $\kappa$  is 1 or greater.

Unfortunately, these conditions may be hard to verify in practice, especially when the distribution  $h(\mathbf{z})$  is not known. (In fact,  $\mathbf{E}[h(\mathbf{z})]$ , as seen above, is the unknown integral to be estimated.)

For this purpose, each number in the stream of random variates is generated as follows (see also Kunsch et al.):

1. Set  $\mathbf{z}$  to an  $n$ -dimensional vector (list of  $n$  numbers) chosen at random in the sampling domain, independently of any other choice. Usually,  $\mathbf{z}$  is chosen *uniformly* at random this way (see note later in this section).
2. Calculate  $h(\mathbf{z})$ , and set the next number in the stream to that value.

Alternatively, if  $h(\mathbf{z})$  can take on only numbers in the closed interval  $[0, 1]$ , the much simpler *Algorithm E* can be used on the newly generated stream (taking  $f(x) = x$ ), rather than *Algorithm C*.

The following example (coded in Python for the SymPy computer algebra library) shows how to find parameter  $\kappa$  for estimating the integral of  $\min(Z1, Z2)$  where  $Z1$  and  $Z2$  are each uniformly chosen at random in the interval  $[0, 1]$ . It assumes  $p = 2$  and  $q = 4$ . (This is a trivial example because we can calculate the integral directly —  $1/3$  — but it shows how to proceed for more complicated cases.)

<h1>Distribution of Z1 and Z2</h1>

```
u1=Uniform('U1',0,1)
```

```
u2=Uniform('U2',0,1)
```

<h1>Function to estimate</h1>

```
func = Min(u1,u2)
```

```
emean=E(func)
```

```
p = S(2) # Degree of p-moment
```

```
q = S(4) # Degree of q-moment
```

<h1>Calculate value for kappa</h1>

```
kappa = E(Abs(func-emean)**q)**(1/q) / E(Abs(func-emean)**p)**(1/p)
```

```
pprint(Max(1,kappa))
```

**Note:** As an alternative to the usual process of choosing a point uniformly in the *whole* sampling domain, *stratified sampling* (Kunsch and Rudolf 2018)<sup>16</sup>, which divides the sampling domain in equally sized boxes and finds the mean of random points in those boxes, can be described as follows (assuming the sampling domain is the  $d$ -dimensional hypercube  $[0, 1]^d$ ):

1. For a sample size  $n$ , set  $m$  to  $\text{floor}(n^{1/d})$ , where  $d$  is the number of dimensions in the sampling domain (number of components of each point). Set  $s$  to 0.
2. For each  $i[1]$  in  $[0, m)$ , do: For each  $i[2]$  in  $[0, m)$ , do: ..., For

each  $i[d]$  in  $[0, m)$ , do:

1. For each dimension  $j$  in  $[1, d]$ , set  $p[j]$  to a number in the half-open interval  $[i[j]/m, (i[j]+1)/m)$  chosen uniformly at random.
2. Add  $f(p[1], p[2], \dots, p[j])$  to  $s$ .
3. Return  $s/m^d$ .

The paper also implied a sample size  $n$  for use in stratified sampling when  $f$  is  $\beta$ -Hölder continuous (is continuous and no "steeper" than  $\mathbf{z}^\beta$ ) and is defined on  $[0, 1]^d$ , namely  $n = \text{ceil}((\ln(2/\delta)/2 * \varepsilon^2)^{d/(2*\beta+d)})$ .

## 9 Finding Coins with Maximum Success Probabilities

Given  $m$  coins each with unknown probability of heads, the following algorithm finds the  $k$  coins with the greatest probability of showing heads, such that the algorithm correctly finds them with probability at least  $1 - \delta$ . It uses the following parameters:

- $k$  is the number of coins to return.
- $\delta$  is the confidence level; the algorithm correctly finds the coins with the greatest probability of showing heads with probability at least  $1 - \delta$ .
- $D$  is a *gap parameter* or a lesser number, but must be greater than 0. The *gap parameter* is the difference between the  $k^{\text{th}}$  most probable coin to show heads and the  $(k+1)^{\text{th}}$  most probable coin to show heads. Practically speaking,  $D$  is the smallest possible difference between one probability of heads and another.
- $r$  is the number of rounds to run the algorithm and must be an integer 1 or greater.

In this section,  $\text{ilog}(a, r)$  means either  $a$  if  $r$  is 0, or  $\max(\ln(\text{ilog}(a, r-1)), 1)$  otherwise.

Agarwal et al. (2017)<sup>17</sup> called this algorithm "aggressive elimination", and it can be described as follows.

1. Let  $t$  be  $\text{ceil}((\text{ilog}(m, r) + \ln(8*k/\delta)) * 2/(D*D))$ .
2. For each integer  $i$  in  $[1, m]$ , flip the coin labeled  $i$ ,  $t$  many times, then set  $P[i]$  to a list of two items: first is the number of times coin

- $i$  showed heads, and second is the label  $i$ .
3. Sort the  $P[i]$  in decreasing order by their values.
  4. If  $r$  is 1, return the labels to the first  $k$  items in the list  $P$ , and the algorithm is done.
  5. Set  $\mu$  to  $\text{ceil}(k + m/\text{ilog}(m, r-1))$ .
  6. Let  $C$  be the coins whose labels are given in the first  $\mu$  items in the list (these are the  $\mu$  many coins found to be the most biased by this algorithm).
  7. If  $\mu \leq 2*k$ , do a recursive run of this algorithm, using only the coins in  $C$  and with  $\delta = \delta/2$  and  $r = 1$ .
  8. If  $\mu > 2*k$ , do a recursive run of this algorithm, using only the coins in  $C$  and with  $\delta = \delta/2$  and  $r = r - 1$ .

## 10 Requests and Open Questions

Let  $X$  be an endless stream of random variates and let  $f(x)$  be a known continuous function.

1. Is there an algorithm, besides *Algorithm C* or *Algorithm F*, that can find  $\mathbf{E}[X]$  (or  $f(\mathbf{E}[X])$ ) with either a high probability of a "small" absolute error or one of a "small" relative error, when the distribution of  $X$  is unbounded, and additional assumptions on the distribution of  $X$  apply, such as—
  - being unimodal (having one peak) and symmetric (mirrored on each side of the peak), and/or
  - following a geometric distribution, and/or
  - having decreasing or nowhere increasing probabilities?

Notice that merely having finite moments is not enough (Theorem 3.4, Kunsch et al. 2019<sup>18</sup>). Here, the accuracy tolerances for small error and high probability are user-specified. A relative-error algorithm for  $\mathbf{E}[X]$  for the geometric distribution was given already in a note.

2. How can *Algorithm D* or *Algorithm E* be adapted to a known discontinuous function  $g$ , so that the algorithm finds  $g(\mathbf{E}[X])$  with either a high probability of a "small" absolute error or one of a "small" relative error at all points in  $[0, 1]$  except at a "negligible" area around  $g$ 's discontinuities? Is it enough to replace  $g$  with a continuous function  $f$  that equals  $g$  everywhere except at that "negligible" area? Here, the accuracy tolerances for small error,



high probability, and "negligible" area are user-specified. Perhaps the tolerance could be defined as the integral ("area under the graph") of absolute differences between  $f$  and  $f$  instead of "negligible area"; in that case, how should the continuous  $f$  be built?

3. Is it true that *Algorithm F* remains valid when the sample size  $n$  is  $\text{ceil}(\text{abs}(M)/(\delta \cdot \gamma^k))$ , given that the stream's distribution is known to have a maximum  $k^{\text{th}}$  central absolute moment of  $M$ ?

## 11 Notes

## 12 License

Any copyright to this page is released to the Public Domain. In case this is not possible, this page is also licensed under **Creative Commons Zero**.

- 
1. Vihola, M., 2018. Unbiased estimators and multilevel Monte Carlo. *Operations Research*, 66(2), pp.448-462.↵
  2. Huber, M., 2017. A Bernoulli mean estimate with known relative error distribution. *Random Structures & Algorithms*, 50(2), pp.173-182. (preprint in arXiv:1309.5413v2 [math.ST], 2015).↵
  3. Feng, J. et al. "Monte Carlo with User-Specified Relative Error." (2016).↵
  4. Huber, M., "**Tight relative estimation in the mean of Bernoulli random variables**", arXiv:2210.12861 [cs.LG], 2022.↵
  5. Huber, Mark, and Bo Jones. "Faster estimates of the mean of bounded random variables." *Mathematics and Computers in Simulation* 161 (2019): 93-101.↵
  6. Huber, Mark, "**An optimal( $\epsilon$ ,  $\delta$ )-approximation scheme for the mean of random variables with bounded relative variance**", arXiv:1706.01478, 2017.↵

7. Kunsch, Robert J., Erich Novak, and Daniel Rudolf. "Solvable integration problems and optimal sample size selection." *Journal of Complexity* 53 (2019): 40-67. Also in <https://arxiv.org/pdf/1805.08637.pdf> .↵
8. Hickernell, F.J., Jiang, L., et al., "**Guaranteed Conservative Fixed Width Intervals via Monte Carlo Sampling**", arXiv:1208.4318v3 [math.ST], 2012/2013.↵
9. As used here, kurtosis is the 4th c.a.m. divided by the square of the 2nd c.a.m.↵
10. "k is even" means that k is divisible by 2. This is true if  $k - 2 \cdot \text{floor}(k/2)$  equals 0, or if the least significant bit of  $\text{abs}(x)$  is 0.↵
11. Kunsch, Robert J., Erich Novak, and Daniel Rudolf. "Solvable integration problems and optimal sample size selection." *Journal of Complexity* 53 (2019): 40-67. Also in <https://arxiv.org/pdf/1805.08637.pdf> .↵
12. Lee, J.C. and Valiant, P., 2022. **Optimal Sub-Gaussian Mean Estimation in Very High Dimensions**. In 13th Innovations in Theoretical Computer Science Conference (ITCS 2022). Schloss Dagstuhl-Leibniz-Zentrum für Informatik.↵
13. Dutta, Santanu, and Alok Goswami. "Mode estimation for discrete distributions." *Mathematical Methods of Statistics* 19, no. 4 (2010): 374-384.↵
14. Jiang, L., Hickernell, F.J., "**Guaranteed Monte Carlo Methods for Bernoulli Random Variables**", arXiv:1411.1151 [math.NA], 2014.↵
15. Hickernell, F.J., Jiang, L., et al., "**Guaranteed Conservative Fixed Width Intervals via Monte Carlo Sampling**", arXiv:1208.4318v3 [math.ST], 2012/2013.↵
16. Kunsch, R.J., Rudolf, D., "**Optimal confidence for Monte Carlo integration of smooth functions**", arXiv:1809.09890, 2018.↵
17. Agarwal, A., Agarwal, S., et al., "Learning with Limited Rounds of Adaptivity: Coin Tossing, Multi-Armed Bandits, and Ranking from Pairwise Comparisons", *Proceedings of Machine Learning Research* 65 (2017).↵

18. Kunsch, Robert J., Erich Novak, and Daniel Rudolf. "Solvable integration problems and optimal sample size selection." Journal of Complexity 53 (2019): 40-67. Also in [\*\*https://arxiv.org/pdf/1805.08637.pdf\*\*](https://arxiv.org/pdf/1805.08637.pdf) .↵