

lintsampler: Easy random sampling via linear interpolation

Aneesh P. Naik¹ and Michael S. Petersen¹

¹ Institute for Astronomy, University of Edinburgh, UK ¶ Corresponding author

DOI: [10.xxxxxx/draft](https://doi.org/10.xxxxxx/draft)

Software

- [Review](#)
- [Repository](#)
- [Archive](#)

Editor: [Open Journals](#)

Reviewers:

- [@openjournals](#)

Submitted: 01 January 1970

Published: unpublished

License

Authors of papers retain copyright and release the work under a Creative Commons Attribution 4.0 International License ([CC BY 4.0](#)).

Summary

lintsampler provides a Python implementation of a technique we term ‘linear interpolant sampling’: an algorithm to efficiently draw pseudo-random samples from an arbitrary probability density function (PDF). First, the PDF is evaluated on a grid-like structure. Then, it is assumed that the PDF can be approximated between grid vertices by the (multidimensional) linear interpolant. With this assumption, random samples can be efficiently drawn via inverse transform sampling (Devroye, 1986).

lintsampler is primarily written with numpy (Harris et al., 2020), drawing some additional functionality from scipy (Virtanen et al., 2020). Under the most basic usage of lintsampler, the user provides a Python function defining the target PDF and some parameters describing a grid-like structure to the LintSampler class, and is then able to draw samples via the sample method. Additionally, there is functionality for the user to set the random seed, employ quasi-Monte Carlo sampling, or sample within a premade grid (DensityGrid) or tree (DensityTree) structure.

Statement of need

Below is a (non-exhaustive) list of ‘use cases’, i.e., situations where a user might find lintsampler (and/or the the linear interpolant sampling algorithm underpinning it) to be preferable over random sampling techniques such as importance sampling, rejection sampling or Markov chain Monte Carlo (MCMC). MCMC in particular is a powerful class of methods with many excellent Python implementations (Coullon & Nemeth, 2022; Foreman-Mackey et al., 2019; Marignier, 2023; Patil et al., 2010). In certain use cases as described below, lintsampler can offer more convenient and/or more efficient sampling compared with these various techniques.

We’ll assume that the target PDF the user wishes to sample from does not have its own exact sampling algorithm (such as the Box-Muller transform for a Gaussian PDF). The power of lintsampler lies in its applicability to arbitrary PDFs for which tailor-made sampling algorithms are not available.

Use Cases

1. Expensive PDF

If the PDF being sampled from is expensive to evaluate and a large number of samples is desired, then lintsampler might be the most cost-effective option. This is because lintsampler does not evaluate the PDF for each sample (as would be the case for other random sampling techniques), but on the nodes of the user-chosen grid. Particularly in a low-dimensional setting

38 where the grid does not have too many nodes, this can mean far fewer PDF evaluations. This
39 point is demonstrated in the [first example notebook](#) in the lintsampler docs.¹

40 2. Multimodal PDF

41 If the target PDF has a highly complex structure with multiple, well-separated modes, then
42 lintsampler might be the *easiest* option (in terms of user configuration). In such scenarios,
43 MCMC might struggle unless the walkers are carefully preconfigured to start near the modes.
44 Similarly, rejection sampling or importance sampling would be highly suboptimal unless the
45 proposal distribution is carefully chosen to match the structure of the target PDF. With
46 lintsampler, the user need only ensure that the resolution of their chosen grid is sufficient to
47 resolve the PDF structure, and so the setup remains straightforward. This is demonstrated in
48 the [second example notebook](#) in the lintsampler docs.²

49 3. PDF with large dynamic range

50 If the target PDF has a very large dynamic range, then the DensityTree object provided by
51 lintsampler might be an effective solution. Here, the PDF is evaluated not on a fixed grid,
52 but on the leaves of a tree. The tree is refined such that regions of concentrated probability
53 are more finely resolved, based on accuracy criteria. Such an example is shown in the [third
54 example notebook](#) in the lintsampler docs.

55 4. Noise needs to be minimised

56 In Quasi-Monte Carlo (QMC) sampling, one purposefully generates more 'balanced' (and
57 thus less random) draws from a target PDF, so that sampling noise decreases faster than
58 $N^{-1/2}$. lintsampler allows easy QMC sampling with arbitrary PDFs. We are not aware of
59 such capabilities with any other package. We give an example of using lintsampler for QMC
60 in the [fourth example notebook](#) in the lintsampler docs.

61 'Real World' Example

62 Any one of the four use cases above would serve by itself as a sufficient case for choosing
63 lintsampler, but here we give an example of a real-world scenario that combines all of the
64 use cases. It is drawn from our own primary research interests in computational astrophysics.

65 Much of computational astrophysics consists of large-scale high performance computational
66 simulations of gravitating systems. For example, simulations of planets evolving and interacting
67 within a solar system, simulations of stars interacting within a galaxy, or vast cosmological
68 simulations in which a whole universe is grown *in silico*.

69 There exists a myriad of codes used to run these simulations, each using different algorithms
70 to solve the governing equations. One class of simulation code that has gained much attention
71 in recent years is the class of code employing basis function expansions ([Petersen et al., 2022](#);
72 [Vasiliev, 2019](#)). In these codes, the matter density at any point in space is represented by a
73 sum over basis functions (not unlike a Fourier series), with the coefficients in the sum changing
74 over space and time. As such, the matter comprising the system is represented everywhere
75 as a smooth, continuous fluid, but for many applications and/or downstream analyses of the
76 simulated system, one needs to instead represent the system as a set of discrete particles.
77 These particles can be obtained by drawing samples from the continuous density field.

¹Similarly, there might be situations where the user is not so concerned about strict statistical representative-
ness but wants to generate a huge number of samples from a target PDF with the least possible computational cost
(such as e.g., generating realistic point cloud distributions in video game graphics). They can use lintsampler
with a coarse grid (so minimal PDF evaluations), then `sample()` to their heart's content.

²It is worth noting that in these kinds of complex, multimodal problems, a single fixed grid might not be
the most cost-effective sampling domain. For this reason, lintsampler also provides simple functionality for
sampling over multiple disconnected grids.

78 This is a scenario that satisfies all four of the use cases list above. To explain further:

- 79 ▪ The PDF we are sampling from (i.e., the basis expansion representation of the matter
80 density field) can be expensive to evaluate if a large number of terms are included in the
81 sum.
- 82 ▪ The PDF can be highly multimodal when the system we are simulating comprises many
83 distinct gravitating substructures, such as stellar clusters.
- 84 ▪ The PDF can have a large dynamic range. Astrophysical structures such as galaxies and
85 dark matter 'haloes' often have power-law density profiles, such as the Navarro-Frenk-
86 White profile (Navarro et al., 1997). Further complicating the issue is that a typical dark
87 matter halo will host several 'subhaloes', which in turn might host 'subsubhaloes', and
88 so on. In short, a range of spatial scales needs to be resolved.
- 89 ▪ If the particle set being sampled is to be used for further simulation, it can be helpful to
90 draw as 'noiseless' a sample as possible for reasons of numerical stability.

91 For these reasons, this kind of astrophysical simulation code provides an excellent example
92 of a 'real world' application for `lintsampler`. Here, one would cover the simulation domain
93 with a `DensityTree` instance (or several instances, one for each primary structure), call the
94 `refine` method to better resolve the high-density regions, then feed the tree to a `LintSampler`
95 instance and call `sample` to generate particles. The `qmc` flag can be passed to the sampler in
96 order to employ Quasi-Monte Carlo sampling.

97 Caveats

98 In all use cases listed above, it is assumed that the dimension of the problem is not too high.
99 `lintsampler` works by evaluating a given PDF on the nodes of a grid (or grid-like structure,
100 such as a tree), so the number of evaluations (and memory occupancy) grows exponentially
101 with the number of dimensions. As a consequence, many of the efficiency arguments given for
102 `lintsampler` below don't apply to higher dimensional problems. We probably wouldn't use
103 `lintsampler` in more than 6 dimensions, but there is no hard limit here: the question of how
104 many dimensions is too many will depend on the problem at hand.

105 Usage

106 `lintsampler` is designed with a interface that makes sampling from an input PDF straightfor-
107 ward. For example, if you have PDF with multiple separated peaks:

```
import numpy as np
from scipy.stats import norm

def gmm_pdf(x):
    mu = np.array([-3.0, 0.5, 2.5])
    sig = np.array([1.0, 0.25, 0.75])
    w = np.array([0.4, 0.25, 0.35])
    return np.sum([w[i] * norm.pdf(x, mu[i], sig[i]) for i in range(3)], axis=0)
```

108 `lintsampler` can efficiently draw samples from it on some defined interval:

```
from lintsampler import LintSampler

grid = np.linspace(-7,7,100)
samples = LintSampler(grid,pdf=gmm_pdf).sample(N=10000)
```

109 `samples` is then an array of 10000 samples drawn from the PDF. Apart from defining the PDF,
110 `lintsampler` enables creating discrete samples from a continuous PDF in a small handful of
111 lines.

Features

Although `lintsampler` is written in pure Python, making the code highly readable, the methods make extensive use of numpy functionality to provide rapid sampling. After the structure spanning the domain has been constructed, sampling proceeds with computational effort scaling linearly with number of sample points.

We provide two methods to define the domain, both optimised with numpy functionality for efficient construction. The `DensityGrid` class takes highly flexible inputs for defining a grid. In particular, the grid need not be evenly spaced (or even continuous) in any dimension; the user can preferentially place grid elements near high-density regions. The `DensityTree` class takes error tolerance parameters and constructs an adaptive structure to achieve the specified tolerance. We also provide a base class (`DensityStructure`) such that the user could extend the methods for spanning the domain.

Documentation for `lintsampler`, including example notebooks demonstrating a range of problems, is available via a [readthedocs page](#). The documentation also has an extensive explanation of the interfaces, including optimisation parameters for increasing the efficiency in sampling.

Acknowledgements

We would like to thank Sergey Koposov for useful discussions. APN acknowledges funding support from an Early Career Fellowship from the Leverhulme Trust. MSP acknowledges funding support from a UKRI Stephen Hawking Fellowship.

References

- Coullon, J., & Nemeth, C. (2022). SGMCMCJax: A lightweight JAX library for stochastic gradient Markov chain Monte Carlo algorithms. *Journal of Open Source Software*, 7(72), 4113. <https://doi.org/10.21105/joss.04113>
- Devroye, L. (1986). *Non-uniform random variate generation*. Springer-Verlag. <https://doi.org/10.1007/978-1-4613-8643-8>
- Foreman-Mackey, D., Farr, W., Sinha, M., Archibald, A., Hogg, D., Sanders, J., Zuntz, J., Williams, P., Nelson, A., de Val-Borro, M., Erhardt, T., Pashchenko, I., & Pla, O. (2019). emcee v3: A Python ensemble sampling toolkit for affine-invariant MCMC. *The Journal of Open Source Software*, 4(43), 1864. <https://doi.org/10.21105/joss.01864>
- Harris, C. R., Millman, K. J., Walt, S. J. van der, Gommers, R., Virtanen, P., Cournapeau, D., Wieser, E., Taylor, J., Berg, S., Smith, N. J., Kern, R., Picus, M., Hoyer, S., Kerkwijk, M. H. van, Brett, M., Haldane, A., Río, J. F. del, Wiebe, M., Peterson, P., ... Oliphant, T. E. (2020). Array programming with NumPy. *Nature*, 585(7825), 357–362. <https://doi.org/10.1038/s41586-020-2649-2>
- Marignier, A. (2023). PxMCMC: A Python package for proximal Markov chain Monte Carlo. *The Journal of Open Source Software*, 8(87), 5582. <https://doi.org/10.21105/joss.05582>
- Navarro, J. F., Frenk, C. S., & White, S. D. M. (1997). A Universal Density Profile from Hierarchical Clustering. *490*(2), 493–508. <https://doi.org/10.1086/304888>
- Patil, A., Huard, D., & Fonnesbeck, C. J. (2010). PyMC: Bayesian stochastic modelling in python. *Journal of Statistical Software*, 35(4), 1–81. <https://doi.org/10.18637/jss.v035.i04>
- Petersen, M. S., Weinberg, M. D., & Katz, N. (2022). EXP: N-body integration using basis function expansions. *510*(4), 6201–6217. <https://doi.org/10.1093/mnras/stab3639>

- 155 Vasiliev, E. (2019). AGAMA: action-based galaxy modelling architecture. *482*(2), 1525–1544.
156 <https://doi.org/10.1093/mnras/sty2672>
- 157 Virtanen, P., Gommers, R., Oliphant, T. E., Haberland, M., Reddy, T., Cournapeau, D.,
158 Burovski, E., Peterson, P., Weckesser, W., Bright, J., van der Walt, S. J., Brett, M., Wilson,
159 J., Millman, K. J., Mayorov, N., Nelson, A. R. J., Jones, E., Kern, R., Larson, E., ... SciPy
160 1.0 Contributors. (2020). SciPy 1.0: Fundamental Algorithms for Scientific Computing in
161 Python. *Nature Methods*, *17*, 261–272. <https://doi.org/10.1038/s41592-019-0686-2>

DRAFT