Contents lists available at ScienceDirect

# **Computer Physics Communications**

journal homepage: www.elsevier.com/locate/cpc



# Implementation of the bin hierarchy method for restoring a smooth function from a sampled histogram<sup>∞</sup>



Olga Goulko a,b,\*, Alexander Gaenko c, Emanuel Gull c, Nikolay Prokof ev a,d, Boris Svistunov a,d,e

- <sup>a</sup> Department of Physics, University of Massachusetts, Amherst, MA 01003, USA
- <sup>b</sup> Department of Physics, Boise State University, Boise, ID 83725, USA
- <sup>c</sup> Department of Physics, University of Michigan, Ann Arbor, MI 48109, USA
- <sup>d</sup> National Research Center "Kurchatov Institute", 123182 Moscow, Russia
- e Wilczek Quantum Center, School of Physics and Astronomy and T. D. Lee Institute, Shanghai Jiao Tong University, Shanghai 200240, China

#### ARTICLE INFO

#### Article history: Received 18 April 2018 Received in revised form 4 September 2018 Accepted 25 September 2018 Available online 30 October 2018

#### ABSTRACT

We present BHM, a tool for restoring a smooth function from a sampled histogram using the bin hierarchy method. It is particularly useful for the analysis of data from large-scale numerical simulations of physical systems, such as diagrammatic Monte Carlo simulations of quantum many-body problems. The theoretical background of the method is presented in Goulko et al., (2018). The code automatically generates a smooth polynomial spline with the minimal acceptable number of knots from the input data. It works universally for any sufficiently regular shaped distribution and any level of data quality (provided that the data are uncorrelated or correlations have been accounted for), requiring almost no external parameter specification. This paper explains the details of the implementation and the use of the program, including a physical example of the restoration of the Fröhlich polaron Green's function from data sampled with diagrammatic Monte Carlo.

**Program summary** Program Title: BHM

Program Files doi: http://dx.doi.org/10.17632/dvj8gxsxpk.1

Licensing provisions: GPLv3 Programming language: C++

External routines/libraries: CMake, GSL

Nature of problem: Restoring a smooth function from a sampled histogram.

Solution method: To make use of all information contained in the sampled data, the BHM algorithm generates a hierarchy of overlapping bins of different sizes from the initially supplied fine histogram. The bin hierarchy is fitted to a polynomial spline with the minimal acceptable number of knots, the positions of which are determined automatically. The output is a smooth function with error band.

© 2018 Elsevier B.V. All rights reserved.

### 1. Introduction

Numerical approaches to problems in condensed matter and quantum many-body physics often involve generating data points according to an unknown probability density f(x), which needs to be restored from the sampled data. The amount of data generated in large-scale quantum Monte Carlo simulations is usually so large that it is impossible (or at least impractical) to store the complete

E-mail address: goulko@umass.edu (O. Goulko).

list of sampled data points  $x_i$ , in order to use, for instance, density estimation protocols or methods based on the cumulative distribution function [1-6] to recover f(x).

Instead, data points are typically collected into a histogram, the histogram bins representing integrals over the sampled distribution. This does not involve any significant loss of information, as long as the bins are sufficiently small to resolve the features of the distribution (which is always possible provided that f(x) is sufficiently smooth). More sophisticated sampling methods exist, which retain more information about the individual points, but these are in general less efficient and require a case-dependent implementation. We provide a universal and efficient program to restore a smooth distribution, which uses the standard histogram as input.

This paper and its associated computer program are available via the Computer Physics Communication homepage on ScienceDirect (http://www.sciencedirect. com/science/journal/00104655).

Corresponding author at: Department of Physics, Boise State University, Boise, ID 83725, USA.

BHM is an implementation of the bin hierarchy method, introduced in [7]. It is

## 1. unbiased:

- utilizes all relevant information contained in the data;
- non-parametric fit automatically adjusts to data quality;
- provides maximally featureless solution (least acceptable number of spline knots);

#### 2. efficient:

- based on regular histogram, which is efficient to sample;
- fast analysis;

#### 3. automatic:

- very little user input;
- no adjustment for different types of sampled functions;
- no adjustment with simulation time as more data is collected.

The bin hierarchy method is applicable to sufficiently regular shaped distributions, meaning that the distribution cannot contain features such as kinks or  $\delta$ -function peaks, which cannot be resolved by a smooth spline. Moreover, sufficient statistics must be collected on the entire domain of the distribution. In particular, this means that special care is required if the distribution changes by many orders of magnitude over a small interval. In such cases, one either requires a large overall statistics, or needs to apply an appropriate reweighing procedure to ensure that the small scales are sufficiently represented in the sampling.

The paper is organized as follows. The general problem setup is presented in Section 2. In Section 3, we give an overview of the algorithm. We explain how to use the program in Section 4, giving a detailed explanation of the input and output formats, as well as possible parameter specifications. Several examples are presented in Section 5, including the physical example of the reconstruction of the smooth Fröhlich polaron Green's function from data generated with diagrammatic Quantum Monte Carlo.

## 2. Problem setup

The central object in BHM is a smooth function f(x) defined on a bounded domain  $\mathcal{D}$ . Statistical sampling with a probability density p(x) is performed to generate samples for f(x) according to  $f(x_j)/p(x_j)$  with p-distributed  $x_j$ . This means that for each value  $x_j$  generated from the probability density p(x), the value  $f(x_j)/p(x_j)$  is sampled. In the simplest case, when f(x) itself is a normalized probability distribution, p(x) = f(x) can be chosen, implying  $f(x_j)/p(x_j) = 1$ . The samples are binned into a histogram with  $2^K$  bins. We are interested in restoring a smooth function  $\tilde{f}(x)$  from this histogram.

Each histogram bin i with bin boundaries  $x_{i,\min}$  and  $x_{i,\max}$  represents the stochastic integral

$$I_i = \int_{x_{i,\text{min}}}^{x_{i,\text{max}}} f(x) dx \tag{1}$$

through the following relations:

$$I_i = \bar{f}_i \frac{N_i}{N},\tag{2}$$

$$M_2(I_i) = M_2(f_i) + \bar{f}_i^2 \frac{N_i(N - N_i)}{N},$$
 (3)

$$Var(I_i) = \frac{M_2(I_i)}{N-1},\tag{4}$$

$$\delta I_i = \sqrt{\frac{\text{Var}(I_i)}{N}},\tag{5}$$

where  $N_i$  is the number of samples in bin i, N the total number of samples,  $\bar{f_i}$  is the average over all samples in the bin i,

$$\bar{f}_i = \frac{1}{N_i} \sum_{x_i \in \text{bin } i} \frac{f(x_j)}{p(x_j)},\tag{6}$$

and the "scaled variance"  $M_2(f_i) = (N_i - 1) \text{Var}(f_i)$  is the sum of squares of differences from the mean,

$$M_2(f_i) = \sum_{x_i \in \text{bin } i} \left( \frac{f(x_j)}{p(x_j)} - \bar{f}_i \right)^2.$$
 (7)

Note that in the simplest case p(x) = f(x) we have  $\bar{f}_i = 1$  and  $M_2(f_i) = 0$ , so that all quantities in Eqs. (2)–(5) are determined through  $N_i$  and N alone.

The goal is to find a function  $\tilde{f}(x)$  whose *integrals* over different parts of the domain  $\mathcal{D}$  agree with the sampled integrals. Working with integrals rather than interpolated function values allows us to include combinations of histogram bins into the fitting. Rebinning data to larger bin sizes leads to a reduction of statistical noise, while retaining small bins results in a higher resolution due to smaller discretization errors.

The resulting fit  $\tilde{f}(x)$  is a polynomial spline of order m, where m is the highest power with non-zero coefficient. The spline function and its derivatives up to order m-1 are continuous, to account for the smoothness of the original f(x).

## 3. Overview of the algorithm

In this section we give a brief overview of the algorithm. More details on the theoretical background of the method can be found in [7]. A flowchart of the algorithm is shown in Fig. 1.

- The algorithm starts from a list of  $2^K$  histogram bins supplied in an input file (for a detailed format description, see Section 4). Typical values of K are 7–15. It should be noted that the bins are not required to have the same size; however, in practice there is no need to have variable size bins. The bins must not overlap or leave gaps.
- From this input the code generates a hierarchy of histogram bins of increasing size. Combining two neighboring bins of the  $2^K$  initial bins leads to  $2^{K-1}$  larger bins with, on average. twice as many entries. For example, combining initial bin 1 and initial bin 2 results in a new bin with  $N_1 + N_2$  samples and with the sample average  $(f_1N_1 + f_2N_2)/(N_1 + N_2)$ . Successive repetitions of this rebinning result in a hierarchy of levels with  $2^{K-2}, \ldots, 2, 1$  bins on each level, respectively, where the final level consists of one bin over the entire domain containing all sampled data, as illustrated in Fig. 2. Bins that do not contain enough data for meaningful statistic, i.e. when the bin counter  $N_i$  is smaller than a user defined minimal value, are excluded from the fitting process. Likewise, levels that do not contain enough usable bins (the minimal fraction can be defined by the user) are also excluded. This implies that in general fitting starts with a level K' < K so that the original binning can be chosen to be very fine without introducing noise into the final fit. For bins that will be used for fitting, the bin integrals and their errors are computed via Eqs. (2)–(5).
- The code checks if the data is compatible with zero on the whole domain. There is an option not to proceed with the fit if this is the case. This feature is particularly useful for data suffering from a severe sign problem.

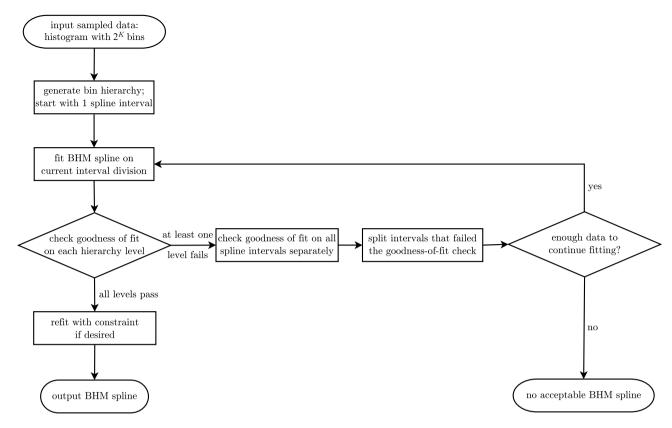


Fig. 1. Flowchart of the algorithm.

• The next step is fitting a spline of order m on the given spline interval division. The starting point is one spline interval, which means that one polynomial  $P(x) = \sum_{k=0}^{m} a_k x^k$  is fitted on the whole domain. The fit minimizes

$$\sum_{n=0}^{K} \frac{\chi_n^2}{2^n},\tag{8}$$

where  $\chi_n^2$  is defined for bins on hierarchy level n in the usual way:

$$\chi_n^2 = \sum_{\text{bins } i \text{ on level } n} \left( \frac{I_i - I_i^{(P)}}{\delta I_i} \right)^2, \tag{9}$$

where  $I_i^{(p)} = \int_{x_{i,\min}}^{x_{i,\max}} P(x) dx$  is the integral of the spline over the bin *i*.

 Afterwards the goodness of fit is evaluated on each hierarchy level individually. The criterion is

$$\frac{\chi_n^2}{\tilde{n}} \le 1 + T\sqrt{\frac{2}{\tilde{n}}},\tag{10}$$

where T is the fit acceptance threshold (input parameter) and  $\tilde{n}$  the number of bins on level n that were used for fitting. The expression  $\sqrt{2/\tilde{n}}$  corresponds to one standard deviation of the  $\chi^2$ -distribution.

- If at least one level fails the global goodness-of-fit check, the goodness-of-fit is then evaluated on each spline interval separately (again level by level). Spline intervals on which the fit was acceptable remain unchanged, while the others are split into two parts, by introducing a spline knot in the middle ("number of bins"-wise).
- If any of the resulting intervals is too small, meaning that there is not enough data to fit on that interval, the code exits without having produced an acceptable spline. Otherwise the BHM fit is repeated on the new interval division.

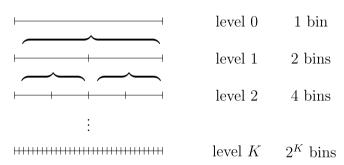


Fig. 2. Illustration of the bin hierarchy.

- Once an acceptable spline has been found, there is an option to refit the data on the same interval division with an additional constraint that aims to minimize the jump in the highest derivative.
- The resulting BHM spline is output (spline coefficients and error coefficients). In addition, the spline values can be output evaluated on a grid.

### 4. Input and output

## 4.1. Running the program

Instructions for compiling the program and executing unit tests can be found in the README file.

The program executable requires 1 argument, the name of the parameter file, e.g.:

# \$ ./bhm in.param

In particular, the parameter file determines the name of the input file with the histogram data and the name of the output file for the BHM spline (see below).

#### Parameter File in.param

DataPointsMin=100	#minimal number of data points per bin	1
SplineOrder=3	#spline order	2
MinLevel=2	#minimal number of levels per interval	3
Threshold=2.0	#minimal goodness of fit threshold	4
ThresholdMax=4.0	#maximal goodness of fit threshold	5
ThresholdSteps=4	#number of steps to take from Threshold to ThresholdMax	6
UsableBinFraction=0.25	#minimal proportion of good bins for a level to be used	7
JumpSuppression=false	${\it \#suppress\ highest\ order\ derivative\ (the\ error\ is\ unreliable)}$	8
Verbose=true	#verbose output	9
PrintFitInfo=true	#print the fit info	10
FailOnBadFit=true	#do not proceed if the fit is bad	1
FailOnZeroFit=true	#do not proceed if the data is consistent with zero	12
Data="histogram.dat"	#location of histogram input data	13
OutputName="spline.dat"	#location of the output file	14
<pre>GridOutput="spline_plot.dat"</pre>	#location of the spline-on-a-grid output	15
GridPoints=1024	#number of points for spline-on-a-grid output	16

Fig. 3. Sample parameter file.

As a special case, if the parameter file name is an empty string, the default parameters will be used which are suitable for most applications:

#### \$ ./bhm "" <histogram.dat >spline.dat

In this case, the histogram data input is expected to be provided at the standard input, and the results will be printed to the standard output. In the example above, the standard input is redirected from file histogram.dat, and the standard output is redirected to file spline.dat. Note that in this particular case the fit information is printed to standard error rather than to standard output.

Without an argument, the program prints a short help message and exits.

## 4.2. Histogram input format

The input histogram data is text-based, line-oriented, and has the following format:

A N <sub>exc</sub>	1
$x_{1,\min} N_1 \bar{f}_1 M_2(f_1)$	2
$x_{2,\min} N_2 \bar{f}_2 M_2(f_2)$	3
	4
$x_{i,\min}$ $N_i$ $\bar{f}_i$ $M_2(f_i)$	5
	6
$\chi_{max}$	7

where the first line specifies an overall normalization factor A and the number  $N_{\rm exc}$  of samples outside of the histogram bounds. The normalization step is omitted if either A=1 or A=0. Otherwise, all values  $\bar{f}_i$  and  $M_2(f_i)$  are divided by A and  $A^2$ , respectively, before constructing the BHM fit. The value  $N_{\rm exc}$  is used to calculate the total number of samples  $N=N_{\rm exc}+\sum_i N_i$ , which is needed for Eqs. (2)–(5).  $N_{\rm exc}$  can be zero.

Starting from the second line, each line, except the last one, contains 2 or 4 blank-separated values, specifying the left bin boundary, the number of samples in the bin, and, optionally, mean value and scaled variance. For example, line 5 of the listing corresponds to a bin i with the left boundary  $x_{i,\min}$ , number of samples  $N_i$ , mean value  $\bar{f}_i$  and scaled variance  $M_2(f_i)$ . The latter two quantities are defined in Eqs. (6) and (7). If the mean value and the scaled variance are both omitted, they are assumed to be  $\bar{f}_i = 1$  and  $M_2(f_i) = 0$ , which corresponds to only ever adding 1 to bin counters, or in other

words p(x) = f(x), as detailed in Section 2. The last line of the file (line 7 of the listing) must contain a single entry  $x_{\text{max}}$ , the right boundary of the last bin.

The numbers  $x_{1,\min} \ldots < x_{i,\min} \ldots < x_{\max}$  must form a strictly monotonically increasing sequence, corresponding to non-overlapping, finite-size bins with no gaps. In the current implementation, the number of bins must be a power of 2 (in the later versions we may remove this limitation).

It is important to note that all supplied variances are assumed to be *uncorrelated*. If correlations are present in the sampled data, they have to be removed prior to the BHM fit, for example through appropriate blocking analysis or by scaling the variances with the estimated correlation factor.

## 4.3. Input parameter format

The input parameter file is a text-based, line-oriented file that has a key = value format. An example input is shown in Fig. 3. The keys are case-insensitive; the string values may be enclosed in quotes; the # symbol starts a comment which is ignored until the end of the line. The meaning of each parameter is indicated in the figure in the corresponding comment. Below we provide more detailed explanations for some of the parameters.

DataPointsMin in line 1 specifies the minimal number of data points that a bin must contain in order to be used for fitting. Bins that contain fewer sampled points are ignored (but still contribute in combination with other bins at higher hierarchy levels). DataPointsMin must be at least 10, in order to ensure that meaningful statistics can be made from the data. The default value is 100. If a hierarchy level does not contain enough usable bins (the minimal number is given by the parameter UsableBinFraction in line 7, times the total number of bins on that level) then this level and all subsequent levels are completely omitted from the fitting.

Whenever a spline interval fails the goodness-of-fit test, it is divided in half by placing an additional spline knot in the middle ("number of bins"-wise). Such a division is no longer allowed when there are too few bins left that are *fully inside* the spline interval. Since the algorithm is designed in such a way that the spline knots always coincide with a boundary between two bins on some hierarchy level, this can be quantified by specifying the maximal such level allowed. This is achieved via the parameter MinLevel in line 3. For example, if there are  $2^K$  elementary bins, MinLevel=2

means that the smallest possible spline intervals coincide with the bins on hierarchy level K-2. This means that there is one bin on level K-2 which is fully inside such an interval (its boundaries are identical with the interval boundaries), two bins on level K-1 and four bins on the base level K, corresponding to a total of seven bins fully inside the smallest possible interval. MinLevel must be at least 2; moreover, MinLevel must be large enough to ensure that the fit is underdetermined for each interval, in other words that there are more bins fully inside the smallest possible interval than there are spline parameters.

The fit acceptance threshold T (lines 4–6) can be either set to a fixed value, or to a range of values between Threshold and ThresholdMax. In the latter case, a BHM fit is first attempted with the smallest value Threshold. If no acceptable fit is found, the threshold value is successively increased in ThresholdSteps equidistant steps, until either an acceptable fit is produced or ThresholdMax is reached. Setting ThresholdMax to be smaller or equal to Threshold and/or setting ThresholdSteps=0 corresponds to only using one fixed value of T. Note that threshold values that are too low can result in overfitting (too many spline pieces) or the failure to produce an acceptable fit. Values that are too high can result in underfitting (too few spline pieces and a poor fit with underestimated error bars). These issues are illustrated in Example 5.2. The value T = 2.0 is good generic choice. Specifying a range of threshold values reduces the statistical chance that there is no acceptable BHM fit with a given threshold, even though the data quality is adequate. The default range between T = 2.0 and T = 4.0 is suitable for most data sets.

Generally, the default parameter values in the example file are suitable for all types of sampled functions, and hence there is no need to change any of the parameters unless specifically desired.

# 4.4. Output format

The default verbose output is printed to standard output and contains auxiliary information such as values of the input parameters, a brief description of the input histogram, and the log of the fitting process. The fitting log is described in detail in Example 5.1. If requested by the PrintFitInfo input parameter, information about the final fit is also printed to the standard output.

The output of the program is both human and machine-readable, and has the following text-based, line-oriented, blank-separated format:

```
# Arbitrary comments
                                                                                               1
# ...
                                                                                               2
m s
                                                                                               3
x_1 x_2 \ldots x_s
# spline piece 1
a_0 a_1 a_2 ... a_m
\varepsilon_0 \varepsilon_1 \varepsilon_2 ... \varepsilon_{2m}
                                                                                               7
                                                                                               8
# spline piece i
                                                                                               9
a_0 a_1 a_2 ... a_m
                                                                                               10
                                                                                               11
\varepsilon_0 \varepsilon_1 \varepsilon_2 ... \varepsilon_{2m}
# spline piece (i+1)
                                                                                               12
```

Any lines at the beginning of the file that start with # are considered comments and are ignored. The first significant line of the file (line 3 of the listing) specifies the spline polynomial order m and the number of splines pieces s; the next line (line 4 of the listing) lists all (s+1) spline piece boundaries  $x_1, \ldots, x_{s+1}$ . The following lines form s sections describing each spline piece  $\tilde{f_i}$ , for  $i=1\ldots s$ . Each section (lines 5–7, 9–11 of the listing) consists of 3 lines:

- 1. Header (starts with #) specifying the spline piece number (i),
- 2. (m+1) numbers specifying the spline piece coefficients  $a_0 \dots a_m$   $(\tilde{f}_i(x) = \sum_{k=0}^m a_k x^k)$ ,
- 3. (2m+1) numbers  $\varepsilon_0 \dots \varepsilon_{2m}$  specifying the error bar  $E_i(x) = \sqrt{\sum_{k=0}^{2m} \varepsilon_k x^k}$ .

#### 4.5. Plotting the resulting spline

The simplest way to plot the resulting spline is to use the provided Python3 script bhm\_spline.py, as follows:

```
$ python3 bhm_spline.py spline.dat
```

On the other hand, it may be convenient to customize the plot and/or compare it with a known function, or plot it interactively (e.g., from a Jupyter notebook). For this purpose the script can be imported as a module that provides a BHM Spline class. The following listing demonstrates a possible way of using the module.

```
import numpy as np
import matplotlib.pyplot as plt
                                                   2
from bhm_spline import BHMSpline
                                                   3
                                                   4
spline=BHMSpline("spline.dat")
                                                   5
x=np.linspace(*spline.domain())
                                                   6
# reference function:
                                                   7
def fn(x): return (x**4-0.8*x*x)/0.171964
                                                   8
# plot the spline and the reference:
                                                   9
plt.plot(x,spline(x), x,fn(x))
                                                   10
# plot the errorbar:
                                                   11
plt.plot(x,spline.errorbar(x))
                                                   12
# plot the spline with errorbars:
                                                   13
spline.plot()
                                                   14
# plot the spline and a reference:
                                                   15
spline.plot(fn)
                                                   16
# plot difference between spline and reference
                                                   17
    with error bar:
spline.plot_difference(fn)
                                                   18
```

In line 3 the class BHMSpline is imported; line 5 creates the object representing the spline. In line 6 an interval of x-values is created corresponding to the domain of the spline. Line 8 defines a reference function to compare with the spline. In line 10 the spline and the reference function are plotted using the Matplotlib plotting library; in line 12 the error bar E(x) is plotted. The class also provides a convenience plotting method: when called without arguments (as on line 14), the spline is plotted along with the error bars; when a function is passed as an argument (line 16), its graph is plotted also. It is also possible to plot the difference between the spline and the reference function with error bar (line 18).

## 4.6. Grid output

If the GridOutput parameter in the parameter file is set to a non-empty filename, the program also outputs to the specified file the values and the error bars of the spline computed on a one-dimensional grid of points. A plotting program, such as gnuplot, can then be used to plot the generated function and the error bars and to compare them with a reference function; for example:

```
$ gnuplot
gnuplot> quartic(x)=(x**4-0.8*x*x)/0.171964 2
gnuplot> plot "spline_plot.dat" with errors 3
gnuplot> replot quartic(x) 4
```

In this example, line 1 of the listing starts the gnuplot program; line 2 defines a reference function (quartic polynomial); line 3 plots the grid output file generated by BHM; and line 4 plots the reference function on the same graph.

#### Parameter File generator.param

SampleSize=10000	#total number of points sampled	1
Function=exponential	#name of test function	2
PowerBins=10	#generates uniform histogram with 2^PowerBins bins	3
RandomSeed=956475	#specify random seed (O means random initialization)	4
Output="histogram.dat"	#histogram output file (missing means "standard output")	5
<pre>GridOutput="function.dat"</pre>	#output test function on a grid into this file (if provided,	) 6
GridPoints=1024	#number of grid points for test function output	7

Fig. 4. Sample parameter file to generate example input.

## 5. Examples

In this section we present three detailed examples of the features of BHM illustrated on different distributions f(x). We provide a program to generate the input data for these examples (as well as for several additional test functions). Calling the program without arguments:

#### \$ ./generator

prints a brief help message, which includes a list of the functions supported by the program.

Calling the program with a single file argument:

## \$ ./generator generator.param

generates the histogram data for a given analytical function according to the parameters listed in the generator.param file. For all examples discussed below, the parameters are the same as shown in the example generator parameter file shown in Fig. 4 (including the random number generator seed), except when stated otherwise.

Calling the program as:

# \$ ./generator -python name

(where *name* is the name of the function, possibly abbreviated) prints the Python code that corresponds to the function, which is convenient for plotting the analytical function against the approximating spline in an interactive Python environment (as has been discussed in Section 4.5).

If the GridOutput parameter in the parameter file is set to a non-empty filename, the program also outputs the values of the function computed on a one-dimensional grid to the specified file; a plotting program, such as gnuplot, can then be used to plot the generated function; for example:

In this example, line 1 of the listing starts the gnuplot program; line 2 plots the generated function; and line 3 plots the content of the spline\_plot.dat generated by BHM as discussed in Section 4.6.

#### 5.1. Example 1

This example demonstrates BHM fits for different choices of spline order m.

The original function is a quartic polynomial (Function=quartic\_polynomial):

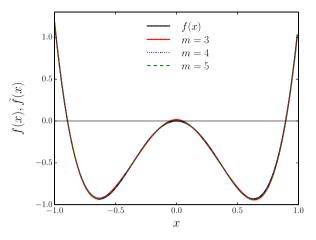
$$f(x) = \alpha(x^4 - 0.8x^2). \tag{11}$$

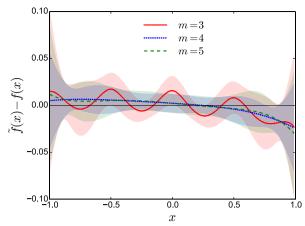
Because f(x) changes sign, sampling on the interval [-1, 1] is performed with the probability density p(x) = |f(x)| and  $\alpha = 0.171964$  is chosen to ensure normalization of p(x) on this interval.

The histogram data is fitted with BHM using the default parameters, with the exception of SplineOrder which is set to 3, 4, and 5 respectively. The fit results are shown in Fig. 5. From the output files "spline.dat" it can be seen that the cubic spline has four spline pieces; the quartic spline has one spline piece, as expected; the quintic spline also has one spline piece, its coefficients up to quartic order are similar to the ones obtained via quartic fit, and its highest spline coefficient is small.

We explain in detail the verbose output for the cubic fit m=3. At the beginning of the output, the fit parameters are listed, as well as general information about the input histogram. Then follows information about the goodness-of-fit at the different fitting stages:

				1
BHM :	fit:			2
Begi	n BHM f	itting with th	reshold T = 2	3
Chec	king se	parate chi_n^2	/n in spline fit	4
leve	l n	chi_n^2/n	max chi_n^2/n	5
0	1	9.7585	3.8284	6
1	2	2.7736	3.0000	7
2	4	1.7636	2.4142	8
3	8	832.1519	2.0000	9
4	14	539.3412	1.7559	10
5	24	210.2518	1.5774	11
6	41	118.5452	1.4417	12
7	54	67.7739	1.3849	13
Chec	king in	terval O (orde	r: 0, number: 0)	14
0	ĭ	9.7585	3.8284	15
This	interva	al fit is not	good	16
			/n in spline fit	17
leve	l n	chi_n^2/n	max chi_n^2/n	18
0	1	0.0020	3.8284	19
1	2	0.0006	3.0000	20
2	4	1.6409	2.4142	21
3	8	7.0923	2.0000	22
4	14	7.8734	1.7559	23
5	24	5.2920	1.5774	24
6	41	3.3034	1.4417	25
7	54	2.0987	1.3849	26
Chec	king in		r: 1, number: 0)	27
1	ĺ	0.0005	3.8284	28
2	2	1.7172	3.0000	29
3	4	6.3727	2.4142	30
This	interva	al fit is not		31
			r: 1, number: 1)	32
1	ĭ	0.0006	3.8284	33
2	2	1.5645	3.0000	34
3	4	7.8119	2.4142	35
This	interva	al fit is not	good	36
			/n in spline fit	37
leve		chi_n^2/n	max chi_n^2/n	38
0	1	0.0001	3.8284	39
1	2	0.0002	3.0000	40
2	4	0.0055	2.4142	41
3	8	0.0519	2.0000	42
4	14	0.3837	1.7559	43
5	24	0.8437	1.5774	44





**Fig. 5.** Quartic polynomial test function f(x) and BHM fits  $\tilde{f}(x)$  with different spline orders m. The left panel shows the function and the fits; the right panel shows the difference between them.

First a fit is attempted with one spline piece on the whole domain (lines 4–13). This fit is not acceptable because  $\chi_n^2/\tilde{n}$  (third column in the output) exceeds the maximally allowed value  $1+T\sqrt{2/\tilde{n}}$  (fourth column in the output) for most of the levels. The second column lists  $\tilde{n}$ , the number of available bins at each level. This number is in general smaller than  $2^n$ , because some bins do not contain enough data to be used for fitting. Also, hierarchy levels below n=7 were omitted because the fraction of usable bins on these levels was below the set UsableBinFraction value.

Since the first fit was unsuccessful,  $\chi^2$  is evaluated on each spline interval separately (lines 14–16). In this case, this yields no new information, since only one interval is present. As soon as a level is found where the fit is unacceptable (level 0 in this case), this check stops without proceeding to lower levels, since this is enough to identify a bad interval.

After the interval is divided, another BHM fit is attempted on two intervals (lines 17–26). This fit already has smaller  $\chi_n^2/\tilde{n}$  values than the previous one, but still fails the threshold on several levels. Both spline intervals are then again checked separately (lines 27–31 and 32–36, respectively) and both fail the goodness-of-fit check on level 3. Note that level 0 is not present in the individual interval checks, because the bin on this level is larger than each of the spline intervals.

The intervals are numbered consecutively, but additional information is provided so that their location can be recovered (see e.g. lines 27 and 32). The boundaries of an interval always coincide with the boundaries of a bin on a certain hierarchy level (denoted by "order") and "number" denotes the number of this bin.

After the intervals are again divided, the resulting BHM fit (lines 38–47) is acceptable. No separate interval checks need to be performed and the code exits with the fit result. If PrintFitInfo is requested, the goodness-of-fit information of the final result is output again at the end. This includes the  $\chi_n^2/\tilde{n}$  values on each level n, the unit standard deviation  $\sqrt{2/\tilde{n}}$  of the corresponding  $\chi^2$ -distribution, as well as the number of standard deviations by which  $\chi_n^2/\tilde{n}$  exceeds 1 on each level (last column). If  $\chi_n^2/\tilde{n} \leq 1$  the latter value is 0.

#### 5.2. Example 2

This example demonstrates BHM fits for different choices of the threshold T. The sampled distribution is a decaying exponential

(Function=exponential),

$$f(x) = \alpha \exp(-3x),\tag{12}$$

normalized on the interval [1, 3], which implies  $\alpha = 3e^9/(e^6-1)$ . The function is sampled on the interval [1, 2.8], so that there is a finite number of values  $N_{\rm exc}$  sampled outside of the histogram bounds. The total number of sampled points in this example is SampleSize=100000.

The histogram data is fitted with BHM using the default parameters, with the exception of the parameters defining the fit acceptance threshold, which is set to be fixed at T=0, 2, and 8, respectively. This can be achieved by either setting the value of ThresholdMax to be equal or less than the value of Threshold, or by setting ThresholdSteps=0. The fit results are shown in Fig. 6.

For all threshold values an acceptable fit exists, but with different interval divisions. The extremely low threshold value T=0 (which means that only fits with  $\chi_n^2/\tilde{n} \leq 1$  are accepted) yields an overfitted spline with 12 spline pieces. The value T=2 produces a suitable fit with 3 spline pieces that captures the shape of the test function well. The very high value T=8 yields an underfitted spline with only 2 pieces. This spline deviates strongly from the true function and the error on the spline is severely underestimated.

#### 5.3. Example 3

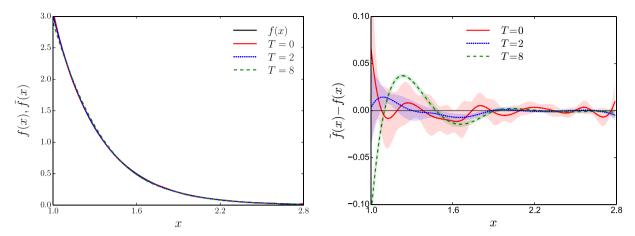
This example demonstrates that BHM works for both uniform and non-uniform input histograms. The sampled distribution,

$$f(x) = 0.2G(0, 0.2) + 0.4[G(2, 1) + G(-2, 1)],$$
(13)

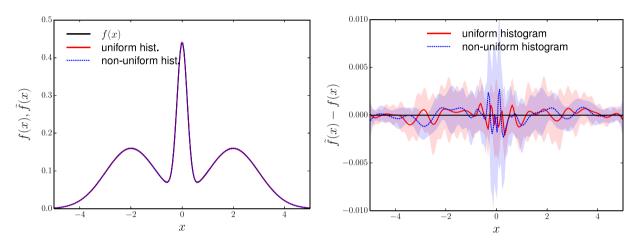
is a linear combination of three Gaussians  $G(\mu, \sigma)$  with mean  $\mu$  and standard deviation  $\sigma$  (Function=triple\_gaussian). It has several distinct features and resembles a physically relevant case.

We sample SampleSize=1000000 data points on the interval [-5,5] into a uniform and a non-uniform histogram, both with 2<sup>8</sup> bins. Note that the non-uniform histogram binning is predefined and cannot be adjusted by changing the PowerBins entry. The non-uniform histogram bins are smaller in the center of the domain (where the sampled function has a sharp feature) and increase exponentially in size towards the domain boundaries. The smallest bin size is equal to the domain length divided by 2<sup>12</sup>. The non-uniform histogram is always collected in addition to the customizable uniform histogram if Function=triple\_gaussian is chosen and is output into the file nonuniform\_histogram.dat.

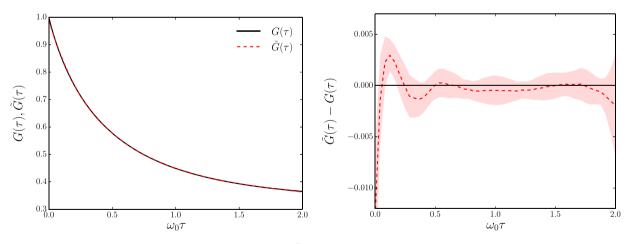
The fit results are shown in Fig. 7. Both histogram divisions produce fits of similar quality that reproduce the tested distribution



**Fig. 6.** Decaying exponential test function f(x) and BHM fits  $\tilde{f}(x)$  of the test function with different goodness-of-fit thresholds T. The left panel shows the function and the fits; the right panel shows the difference between them.



**Fig. 7.** Triple Gaussian test function f(x) and BHM fits  $\tilde{f}(x)$  of the test function based on a uniform histogram and a histogram with bins of different sizes. The left panel shows the function and the fits; the right panel shows the difference between them.



**Fig. 8.** Fröhlich polaron Green's function. Reference calculation  $G(\tau)$  and BHM fit  $\tilde{G}(\tau)$  using default parameters (left panel) and the difference between BHM fit and reference (right panel).

well. Since BHM automatically considers combinations of elementary bins, there is no need for a case-specific implementation of a non-uniform histogram grid. Note that sampling the same data in a uniform histogram with  $2^{12}$  bins produces nearly the same fit as when using  $2^8$  uniform bins in this example.

# 5.4. Example 4

We now demonstrate the effectiveness of BHM on a real physical example: the restoration of the Fröhlich polaron [8] Green's function from data sampled with diagrammatic Monte Carlo [9]. In this

example, we consider the zero momentum imaginary time Green's function of the Fröhlich Hamiltonian,

$$H = H_e + H_{ph} + H_{eph}, (14)$$

$$H_e = \sum_{\mathbf{k}} \frac{k^2}{2} a_{\mathbf{k}}^{\dagger} a_{\mathbf{k}}, \tag{15}$$

$$H_{ph} = \sum_{\mathbf{q}} \omega_0 b_{\mathbf{q}}^{\dagger} b_{\mathbf{q}}, \tag{16}$$

$$H_{eph} = \sum_{\mathbf{k},\mathbf{q}} \frac{i\sqrt{2^{3/2}\alpha\pi}}{q} (b_{\mathbf{q}}^{\dagger} - b_{-\mathbf{q}}) a_{\mathbf{k}-\mathbf{q}}^{\dagger} a_{\mathbf{k}}, \tag{17}$$

with the dimensionless coupling constant  $\alpha=2$  and chemical potential  $\mu=-2.07\omega_0$ , where  $\omega_0$  is the phonon frequency. This Hamiltonian describes an electron (with annihilation operator  $a_{\bf k}$ ) coupled to a bath of phonons (with annihilation operator  $b_{\bf g}$ ).

The Green's function is a fundamental quantity for diagrammatic Monte Carlo, since it gives the most complete information about the system, from which other observables of interest can be extracted. We sample approximately  $2 \cdot 10^6$  data points on the interval [0, 2], divided into  $2^{10}$  elementary bins.

The full Green's function cannot be computed analytically, so we use a very long Monte Carlo run as the reference. The errors on the reference run are several orders of magnitude smaller than the errors on the data used for BHM. Fig. 8 shows the reference function and the BHM fit result of the sampled data using default

parameters. It can be clearly seen the fit is in excellent agreement with the reference function.

#### Acknowledgments

This work was supported by the Simons Collaboration, United States on the Many Electron Problem and by the National Science Foundation, United States under the grant DMR-1720465 (N.P. and B.S.). O.G. also acknowledges support by the US-Israel Binational Science Foundation, Israel (Grants 2014262 and 2016087).

#### References

- [1] I. Narsky, F.C. Porter, Statistical Analysis Techniques in Particle Physics, John Wiley & Sons, 2013.
- [2] D.W. Scott, Multivariate Density Estimation, John Wiley & Sons, 2015.
- [3] B.W. Silverman, Density Estimation for Statistics and Data Analysis, Chapman and Hall, London, 1986.
- [4] B.A. Berg, R.C. Harris, Comput. Phys. Comm. 179 (6) (2008) 443–448, http://dx.doi.org/10.1016/j.cpc.2008.03.010, URL http://www.sciencedirect.com/science/article/pii/S0010465508001458.
- [5] R. van Zon, J. Schofield, J. Chem. Phys. 132 (15) (2010) 154110, http://dx.doi. org/10.1063/1.3366523.
- [6] L.K. Saul, S.T. Roweis, J. Mach. Learn. Res. 4 (2003) 119–155, http://dx.doi.org/ 10.1162/153244304322972667.
- [7] O. Goulko, N. Prokof'ev, B. Svistunov, Phys. Rev. E 97 (2018) 053305, http://dx.doi.org/10.1103/PhysRevE.97.053305, arXiv:1707.07625.
- [8] H. Fröhlich, H. Pelzer, S. Zienau, Phil. Mag. 41 (314) (1950) 221–242, http://dx.doi.org/10.1080/14786445008521794.
- [9] A.S. Mishchenko, N.V. Prokof'ev, A. Sakamoto, B.V. Svistunov, Phys. Rev. B 62 (2000) 6317–6336, http://dx.doi.org/10.1103/PhysRevB.62.6317, arXiv:cond-mat/9910025.