



Universität
Zürich ^{UZH}

Performance of Univariate Kernel Density Estimation methods in TensorFlow

Bachelor Thesis

Author: Marc Steiner

Supervisors: Jonas Eschle, Nicola Serra

University of Zurich

Abstract

Multiple implementations of a one-dimensional Kernel Density Estimation in TensorFlow/Zfit are proposed. Starting from the basic algorithm, several optimizations from recent papers are introduced and combined to ameliorate the efficiency of the algorithm. By comparing its accuracy and efficiency to implementations in pure Python it is shown as competitive and useful in real world applications.

Contents

Abstract	2
1 Introduction	3
1.1 Kernel Density Estimation	3
1.2 zfit and TensorFlow	4
2 Current state of the art	5
3 Implementation	5
3.1 Exact Kernel Density Estimation	5
3.2 Simple and linear Binning	6
3.3 Using convolution and the Fast Fourier Transform	7
3.4 Improved Sheather Jones Algorithm	8
4 Comparison	9
4.1 Benchmark setup	10
4.1.1 Runtime	12
4.1.2 Accuracy	13
4.2 New implementation against KDEpy	13
4.2.1 Runtime	13
4.2.2 Accuracy	13
4.3 New implementation run with GPU support	13
4.3.1 Runtime	13
4.3.2 Accuracy	13
5 Summary	13
References	13

1 Introduction

1.1 Kernel Density Estimation

In many fields of science and in physics especially, scientists need to estimate the probability density function (PDF) from which a set of data is drawn without having a approximate model of the underlying mechanisms, since they are too complex to be fully understood analytically. So called parametric methods fail, because the knowledge of the system is too poor to design a model, which parameters

can then be fitted by some goodness-of-fit criterion like log-likelihood or χ^2 . In the particle accelerator at CERN for instance, they record a whopping 25 gigabytes of data per second¹, resulting from the process of many physical interactions that occur almost simultaneously, making it impossible to anticipate features of the distribution one observes.

To combat this so called non-parametric methods like histograms are used. By summing the data up in discrete bins we can approximate the underlying parent distribution, without needing any knowledge of the physical interactions. However there are also many more sophisticated non-parametric methods, one in particular is Kernel Density Estimation (KDE), which can be looked at as a sort of generalized histogram.²

Histograms tend to produce PDFs that are highly dependent on bin width and bin positioning, meaning the interpretation of the data changes by a lot by two arbitrary parameters. KDE circumvents this problem by replacing each data point with a so called kernel that specifies how much it influences its neighbouring regions as well. The kernels themselves are centered at the data point directly, eliminating the need for arbitrary bin positioning³. Since KDE still depends on kernel width (instead of bin width), one might argue that this is not a major improvement. However, upon closer inspection, one finds that the underlying PDF does depend less strongly on the kernel width than histograms on bin width and it is much easier to specify rules for an approximately optimal kernel width than it is to do so for bin width⁴. In addition, by specifying a smooth kernel, one gets a smooth distribution as well, which is often desirable or even expected from theory.

Due to this increased robustness, KDE is particular useful in High-Energy Physics (HEP) where it has been used for confidence level calculations for the Higgs Searches at the Large Electron Positron Collider (LEP)⁵. However there is still room for improvement and certain more sophisticated approaches to Kernel Density Estimation have been proposed in dependence on specific areas of application⁵.

1.2 zfit and TensorFlow

Currently the basic principle of KDE has been implemented in various programming languages and statistical modeling tools. The standard framework used in HEP, that includes KDE is the ROOT/RooFit toolkit written in C++. However, Python plays an increasingly large role in the natural sciences due to support by corporations involved in Big Data and its superior accessibility. To elevate research in HEP, zfit, a new alternative to RooFit, was proposed. It is implemented on top of TensorFlow, one of the leading Python frameworks to handle large data and high parallelization, allowing a transparent usage of CPUs and GPUs⁶.

So far there exists no direct implementation of Kernel Density Estimation in zfit nor TensorFlow, but various implementations in Python. This implementations will be discussed in the next chapter (2) before I propose an implementation of Kernel Density Estimation in TensorFlow (3). Starting with a rather simple implementation, multiple improvements from recent papers are integrated to amelio-

rate its efficiency. Efficiency and accuracy of the implementation are then tested by comparing it to other implementations in pure Python and simple smoothed histograms (4).

In the last chapter (5) I compare its accuracy and efficiency to smoothed histograms and implementations in pure Python.

2 Current state of the art

...

For humongous data streams (like at CERN 1), Kernel Density Estimation itself needs to be approximated.

...

3 Implementation

3.1 Exact Kernel Density Estimation

The implementation of a simple Kernel Density Estimation in TensorFlow is straightforward. As described in the original Tensorflow Probability Paper⁷, a KDE can be constructed by using its Mixture-SameFamily Distribution, given sampled data as follows

```
from tensorflow_probability import distributions as tfd

f = lambda x: tfd.Independent(tfd.Normal(loc=x, scale=1.))
n = data.shape[0].value

kde = tfd.MixtureSameFamily(
    mixture_distribution=tfd.Categorical(
        probs=[1 / n] * n),
    components_distribution=f(data))
```

Interestingly, due to the smartly capsulated structure of TensorFlow Probability we can use any distribution of the loc-scale family type as a Kernel, if there exists an implementation for it in TensorFlow Probability. If the used Kernel has bounded support, the implementation proposed in this paper allows to specify the support upon instantiation of the class. If the Kernel has infinite support (like a Gaussian kernel i. e.) a practical support estimate is calculated by searching for approximative roots with the Brent's method⁸ implemented in TensorFlow the python package `tf_quant_finance` by Google. This allows us to speed up the calculation

However calculating an exact Kernel Density Estimation is not always feasible as this can take a long time with a huge collection of events, especially in high energy physics. By implementing it in TensorFlow we already get a significant speed up compared to implementations in native Python, since most of TensorFlow is actually implemented in C++ and the code is optimized before running. The computational complexity however, remains the same nonetheless.

$$\hat{f}_h(x) = \frac{1}{nh} \sum_{k=1}^n K\left(\frac{x - x_k}{h}\right) \quad (1)$$

The computational complexity of the basic exact KDE above is $\mathcal{O}(nm)$ where n is the number of sample points to estimate from and m is the number of evaluation points (the points where you want to calculate the estimate).

To combat this complexity several methods exist.

3.2 Simple and linear Binning

The most straightforward way to decrease runtime is by limiting the number of sample points. This can be done by a binning routine, where the values at a smaller number of regular grid points are estimated from the original large number of sample points. Given a set of sample points $X = \{x_0, x_1, \dots, x_k, \dots, x_{n-1}, x_n\}$ with weights w_k and a set of equally spaced grid points $G = \{g_0, g_1, \dots, g_l, \dots, g_{n-1}, g_M\}$ where $N < n$ we can assign an estimate (or a count) c_l to each grid point g_l and use the newly found g_l 's to calculate the kernel density estimation instead. This brings the computational complexity down to $\mathcal{O}(Nm)$. Depending on the number of grid points N the estimate is either more accurate and slower or less accurate and faster. However as we will see in the comparison chapter later as well, even a grid of size 1024 is enough to capture the true density with high accuracy on a million data points.⁹

As described in the excellent overview by Artur Gramacki¹⁰ simple binning or linear binning can be used, although the last is often preferred since it is more accurate and the difference in computational complexity is negligible.

Simple binning is just the standard process of taking a weighted histogram that is divided by the sum of the sample points weights (normalization). In one dimension simple binning is binary in that it assigns a data points weight (1 for an unweighted histogram) either to the bin left or right of itself. Linear binning on the other hand assigns a fraction of the whole weight to both bins on either side, proportional to the closeness of bin and data point in relation to the bin width.

Mathematically linear binning in one dimension can be calculated like this:

$$c_l = c(g_l) = \sum_{\substack{x_k \in X \\ g_l < x_k < g_{l+1}}} \frac{g_{k+1} - x_k}{g_{l+1} - g_l} \cdot w_k + \sum_{\substack{x_k \in X \\ g_{l-1} < x_k < g_l}} \frac{x_k - g_{l-1}}{g_{l+1} - g_l} \cdot w_k \quad (2)$$

Implementing linear binning efficiently with TensorFlow is a bit tricky since loops should be avoided. However with some inspiration from the excellent KDEpy package¹¹ which implements a kernel density estimation in native python/cython this can be done without using loops at all. By transforming the data such that every data point x_k can be described by an integral part (corresponding to its nearest left grid point number l) plus some fractional part (corresponding to the distance between grid point g_l and data point x_k) and applying `tf.math.bincount` twice to the transformed integral part weighting it with the fractional part times the initial weight.

The kernel density estimation can then be calculated as a mixture distribution of kernels located at the grid points, weighted with their associated grid count.

$$\hat{f}_h(x) = \frac{1}{nh} \sum_{l=1}^N c_l \cdot K\left(\frac{x - g_l}{h}\right) \quad (3)$$

3.3 Using convolution and the Fast Fourier Transform

With binning implemented, another technique to speed up the computation is rewriting the Kernel Density Estimation as convolution operation between the kernel and the grid counts calculated by the binning routine. By using the fact that a convolution is just a multiplication in Fourier space one can reduce the computational complexity down to $\mathcal{O}(\log N \cdot m)$.¹⁰

Using the equation (3) from above but also only evaluating it at grid points gives us

$$\hat{f}_h(g_j) = \frac{1}{nh} \sum_{l=1}^N c_l \cdot K\left(\frac{g_j - g_l}{h}\right) = \frac{1}{nh} \sum_{l=1}^N k_{j-l} \cdot c_l \quad (4)$$

where $k_{j-l} = K\left(\frac{g_j - g_l}{h}\right)$. If we set $c_l = 0$ for all l not in the set $\{1, \dots, N\}$ and notice that $K(-x) = K(x)$ we can extend equation (4) to a discrete convolution as follows

$$\hat{f}_h(g_j) = \frac{1}{nh} \sum_{l=-N}^N k_{j-l} \cdot c_l = \vec{c} * \vec{k} \quad (5)$$

where the two vectors look like this

```
knitr::include_graphics('figures/c_conv_k.png')
```

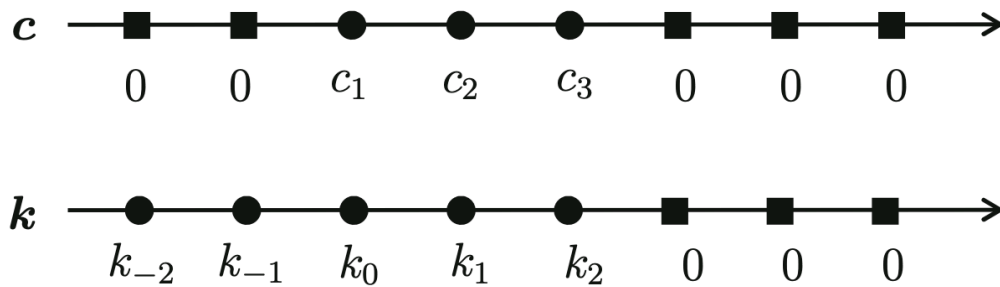


Figure 1: Vectors \vec{c} and \vec{k}

By using the well known convolution theorem we can fourier transform \vec{c} and \vec{k} multiply them and inverse fourier transform them back into real space.

In TensorFlow convolutions are efficiently implemented already in this way if we use `tf.nn.conv1d`. In benchmarking using this method proved significantly faster than using `tf.signal.rfft` and `tf.signal.irfft` to transform, multiply and inverse transform the vectors, which is implemented as an alternative as well.

This algorithm is implemented as its own class since it does not represent a complete mixture distribution anymore but calculates just the density distribution values at the specified grid points. To still infer values for other points in the range of x `tfp.math.interp_regular_1d_grid` is used which computes a linear interpolation of values between the grid.

3.4 Improved Sheather Jones Algorithm

A different take on Kernel Density Estimators is described in the paper “Kernel density estimation by diffusion” by Botev et al.¹². The authors present a new adaptive kernel density estimator based on linear diffusion processes which also includes an estimation for the optimal bandwidth.

The algorithm is quite difficult to understand, a detailed explanation is given in the “Handbook of Monte Carlo Methods”¹³ by the original paper authors. However the general idea is briefly sketched below.

A critical insight is that the Gaussian kernel density estimator $\hat{f}_{h,norm}$ is the solution of the partial differential equation

$$\frac{\partial}{\partial t} \hat{f}_{h,norm}(x, t) = \frac{1}{2} \frac{\partial^2}{\partial x^2} \hat{f}_{h,norm}(x, t), t > 0 \quad (6)$$

with $x \in \mathbb{R}$, $\lim_{x \rightarrow \pm\infty} \hat{f}_{h,norm}(x, t) = 0$ and initial condition $\hat{f}_{h,norm}(x, 0) = \Delta(x)$, where $\Delta(x) = \frac{1}{N} \sum_{k=0}^N \delta_{x_k}(x)$ is the empirical density of the given sample points $X = \{x_0, x_1, \dots, x_k, \dots, x_{n-1}, x_n\}$

and $\delta_{x_k}(x)$ is the Dirac measure at x_k . This means the kernel density estimator can be obtained by evolving the solution of the partial differential equation (6) up to time t . The key observation is that (6) can be solved on a finite domain efficiently using the fast cosine transform - an FFT-related transform.¹³

The optimal bandwidth is often defined as the one that minimizes the mean integrated square error (*MISE*)

$$MISE(t) = \mathbb{E}_f \int [\hat{f}_{h,norm}(x, t) - f(x)]^2 dx \quad (7)$$

An asymptotically optimal value t^* which minimizes a first-order asymptotic approximation of the *MISE* is then given by¹³

$$t^* = \left(\frac{1}{2N\sqrt{\pi}\|f''\|^2} \right)^{\frac{2}{5}} \quad (8)$$

Using the fact that $\|f^{(j)}\|^2 = (-1)^j \mathbb{E}_f[f^{(2j)}(X)]$, $j \geq 1$ and an initial estimation for $\|\hat{f}_{h,norm}^{(l+2)}\|^2$ for some $l \geq 3$ one can then iteratively get an estimation for $\|\hat{f}_{h,norm}^{(2)}\|^2$ which can then be used to estimate t^* instead of $\|f''\|^2$. According to their handbook $l = 7$ is a suitable value to yield good practical results.

The improvement compared to the standard Sheather-Jones plug-on method¹⁴ consists in the fact that to compute the initial estimation of $\|\hat{f}_{h,norm}^{(l+2)}\|^2$ is calculated by solving the partial differential equation using the fast cosine transform as described above, eliminating the need to assume normally distributed data for the initial estimate and leading improved performance, especially for density distributions that are far from normal as seen in the next chapter.

The implementation of the algorithm in TensorFlow proposed in this paper was also inspired a lot by the python package KDEpy¹¹ and uses Brent's method⁸ to find roots implemented in TensorFlow the python package `tf_quant_finance` as well.

One shortcoming of the improved Sheather-Jones algorithm (ISJ) is that with few data points to estimate from, it is not guaranteed to converge. If that happens, one has to use the exact, binned or FFT kernel density estimators as described above.

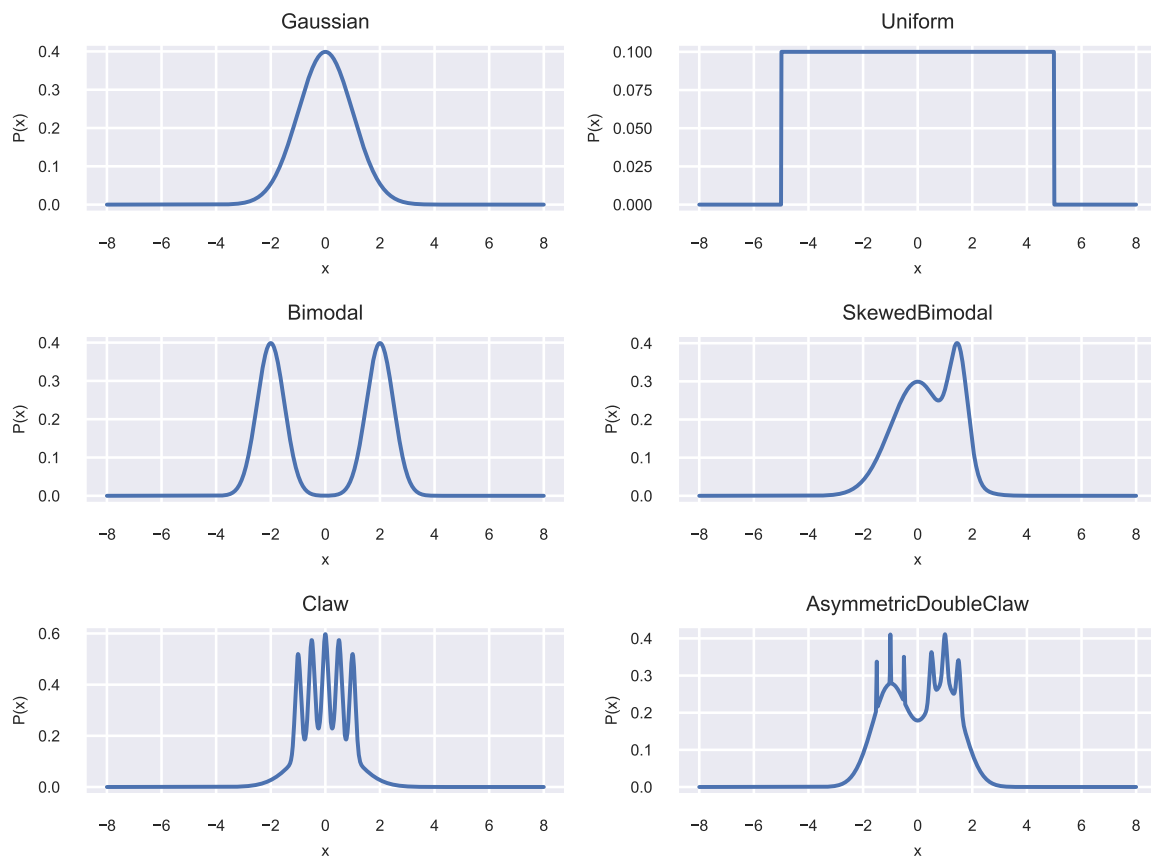
4 Comparison

To show the efficiency and performance of the Kernel Density Estimation methods implemented with TensorFlow a benchmarking suite was developed. It consists of three parts, a collection of distributions, a collection of methods to compare and a runner module that implements helper meth-

ods to execute the methods to test against the different distributions and plot the generated dataset nicely.

4.1 Benchmark setup

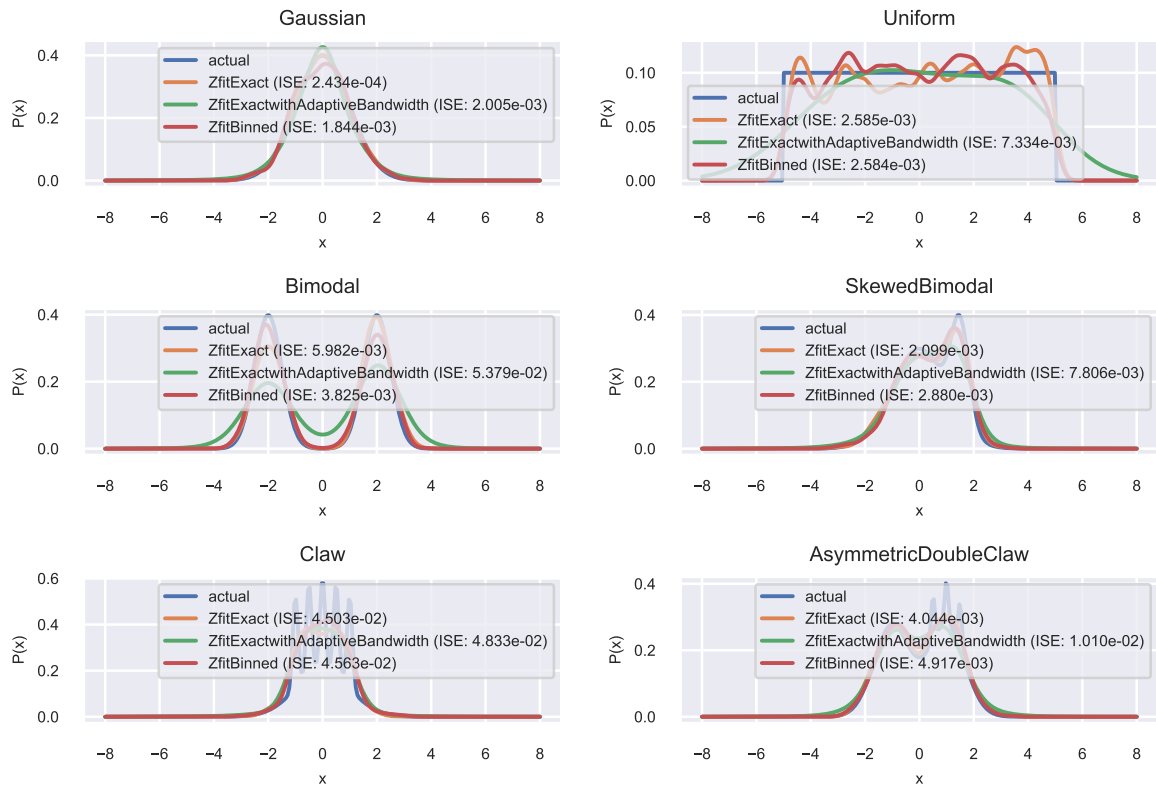
To compare the different implementations multiple popular test distributions mentioned in Wand et al.¹⁵ were used. A simple normal distribution, a simple uniform distribution, a bimodal distribution comprised of two normals, a skewed bimodal distribution, a claw distribution that has spikes and one called asymmetric double claw that has different sized spikes left and right. All comparisons were made using a standard Gaussian Kernel. Although all loc-scale family distributions of TensorFlow Probability may be used for the new implementation proposed in this paper may be used, the Gaussian kernel is the most used one and provides best reference to compare different implementations against eachother.



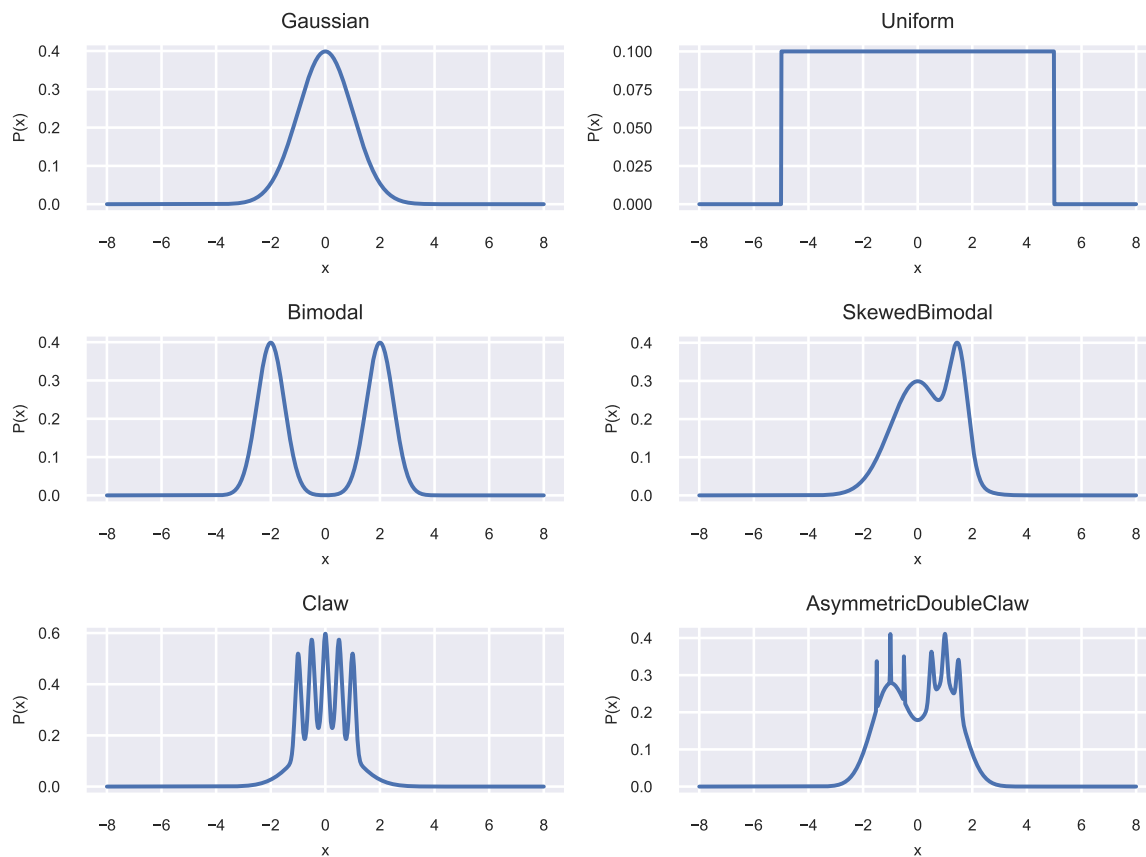
Basic implementation against Binned and FFT implementations

```
figure, axes = runner.plot_estimations(restricted_estimations,
    ↳ distributions_to_evaluate, 1e3, ['ZfitExact',
    ↳ 'ZfitExactwithAdaptiveBandwidth', 'ZfitBinned'])
```

figure



4.1.1 Runtime



4.1.2 Accuracy

4.2 New implementation against KDEpy

4.2.1 Runtime

4.2.2 Accuracy

4.3 New implementation run with GPU support

4.3.1 Runtime

4.3.2 Accuracy

5 Summary

In summary we can conclude that...

References

¹CERN (n.d.).

² M. Rosenblatt, (1956).

³ T. Duong, An Introduction to Kernel Density Estimation (2001).

⁴ M. Lerner, Histograms and Kernel Density Estimation | Biophysics and Beer (2013).

⁵ K.S. Cranmer, (2000).

⁶ J. Eschle, A.P. Navarro, R.S. Coutinho, and N. Serra, (2019).

⁷ J.V. Dillon, I. Langmore, D. Tran, E. Brevdo, S. Vasudevan, D. Moore, B. Patton, A. Alemi, M.D. Hoffman, and R.A. Saurous, CoRR **abs/1711.10604**, (2017).

⁸ R.P. Brent, The Computer Journal **14**, 422 (1971).

⁹ T. Odland, Comparison - KDEpy 1.0.5 Documentation (n.d.).

¹⁰ A. Gramacki, in *Nonparametric Kernel Density Estimation and Its Computational Aspects* (Springer, 2018), pp. 85–118.

¹¹ T. Odland, GitHub Repository (2020).

¹² Z.I. Botev, J.F. Grotowski, D.P. Kroese, and others, The Annals of Statistics **38**, 2916 (2010).

¹³ D.P. Kroese, T. Taimre, and Z.I. Botev, *Handbook of Monte Carlo Methods* (John Wiley & Sons, 2013).

¹⁴ S.J. Sheather and M.C. Jones, Journal of the Royal Statistical Society: Series B (Methodological) **53**, 683 (1991).

¹⁵ M.P. Wand and M.C. Jones, *Kernel Smoothing* (Crc Press, 1994).

List of Tables

List of Figures

1	Vectors \vec{c} and \vec{k}	8
---	---	---

Listings