



zfit: scalable model fitting in Python using TensorFlow

Master thesis of
Jonas Eschle

Supervised by
Prof. Nicola Serra
Dr. Albert Puig Navarro

Abstract

Fitting a model to data is an essential part in most High Energy Physics analyses. Several frameworks to perform this action exist in C++, but no powerful enough counterpart exists in Python, a language recently becoming more and more popular for analyses. With the recent success of deep learning, frameworks in Python such as TensorFlow came up, offering a high level interface for efficient, parallelised computing on modern architectures. In this thesis, `zfit`, a library for model fitting in HEP implemented in pure Python and based on top of TensorFlow, is presented. It offers a well structured model fitting workflow allowing to build composite models from a variety of shapes. A high level of customisation is possible due to well specified interfaces and convenient base classes allowing to easily replace any part in the workflow with a custom implementation. Together with the flexibility and scalability of TensorFlow, `zfit` extends its functionality well beyond what current model fitting libraries offer. An overview over the current status of model fitting libraries and the HEP requirements will be discussed followed by the structure of `zfit` and its implementation. Finally, examples which quantify the performance and demonstrate the feasibility of `zfit` for a whole range of real world applications are shown and an additional library for phasespace generation, `phasespace`, is introduced.

Contents

1	Introduction	1
2	Model fitting	5
2.1	Maximum Likelihood	5
2.2	Requirements	6
2.3	Existing libraries	7
2.3.1	General fitting	7
2.3.2	HEP specific	8
3	zfit introduction	10
3.1	TensorFlow backend	12
4	zfit implementation	17
4.1	Spaces and Dimensions	18
4.1.1	Limits	18
4.2	Data handling	19
4.3	Model	20
4.3.1	Parametrization	21
4.3.2	Implementing a custom PDF	22
4.3.3	Sampling	25
4.3.4	Extended PDFs	26
4.4	Loss	26
4.5	Minimisation	28
4.5.1	Different optimisations	28
4.6	Results and uncertainties	29
4.6.1	Parameter uncertainties	30
5	Performance	31
5.1	Gaussian models	31
5.2	Angular analysis	34
6	Beyond standard fitting	36
6.1	Amplitude fits	37
6.2	phasespace	38
6.3	Dalitz implementation	39
7	Conclusion and outlook	44
A	Likelihood	47
B	Backend	50
B.1	HPC and paradigms	50
B.2	Working with TensorFlow	52
B.2.1	Caching	53

C	Implementation	54
C.1	Spaces definition	54
C.2	General limits	55
C.3	Data formats	56
C.4	Data batching	57
C.5	Dependency management	57
C.6	Base Model	58
C.6.1	Public methods	58
C.6.2	Hooks	59
C.6.3	Norm range handling	59
C.6.4	Multiple limits handling	60
C.6.5	Most efficient method	60
C.6.6	Functors	61
C.7	Sampling techniques	61
C.8	Loss defined	63
D	Performance studies	63
D.1	Hardware specification	63
D.2	Profiling TensorFlow	64
D.3	Additional profiling	65
	References	65

1 Introduction

The Standard Model (SM) of Particle Physics describes the most fundamental particles in the universe and their interactions. According to it, all matter is made up of fermions: quarks and leptons. They appear in different flavours and generations as depicted in Fig. 1. There are additionally four gauge bosons that allow the particles to interact via their exchange: the photon is the electromagnetic force carrier and couples to particles with an electric charge, such as the electron and the quarks. W and Z bosons are the carrier of the weak force and couple to all fermions. Unlike other forces, the strength of the strong interaction increases with the distance of two particles. As a consequence, quarks do not appear alone in nature. Instead, they form composite particles existing of multiple quarks, the hadrons. While there are hundreds of hadrons, nearly all of them have lifetimes under nano seconds and decay to lighter particles. The only stable particles are the well known neutron, proton, electron and neutrino, which together make up the visible matter in our universe. Finally, all particles in the SM (except neutrinos) acquire mass through their interaction with the Higgs Field by the exchange of Higgs Bosons.

With the recent discovery of the Higgs boson, the last missing piece of the SM has been found. It provides a complete description of nearly all observations. And yet it does not seem to be the final answer since there are phenomena that remain unexplained, such as the existence of dark matter that interacts gravitationally and significantly determines the dynamics of galaxies does not appear in the SM, and the fact that neutrinos have mass since they oscillate does not coincide with the predictions of the SM. With larger amounts of data collected, more precise measurements need to be made in order to look for further inconsistencies of the SM that can guide us to a new theory.

In the scientific context, what is called an observations is in fact an answer extracted from nature by asking the right question and using statistics to analyse the response of the data. A question in this sense is an experimental setup and a scientific hypothesis is a proposed explanation for an observed phenomena which can be tested. Different methods can be used to verify a hypothesis, but all of them make use of a test statistic that needs a single value to quantify their agreement with the observations. Given strong enough evidence, the null hypothesis may be rejected in favour of an alternate hypothesis. As an example, this can be used for the discovery of a new particle where the background only hypothesis acts as the null and the alternate is the background and signal combined.

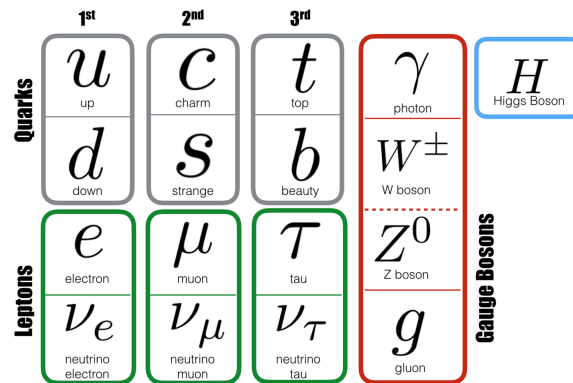


Figure 1: The particles of the SM.

33 As already mentioned, observations from experiments are needed. To study the
34 fundamental particles of the SM, high enough energies are required to produce them.
35 High Energy Physics (HEP) experiments all over the world accelerate light particles
36 such as electrons or protons and let them collide. The Large Hadron Collider (LHC) at
37 CERN accelerates protons to energies up to 6.5 TeV¹, which is currently the frontier in
38 high energies. Around the collider, there are four large experiments: the general purpose
39 detectors ATLAS and CMS, ALICE, an experiment specialized on lead interactions and
40 LHCb, a detector focused on the study of heavy flavour decays. These experiments are
41 situated at collision points around the LHC where 40 million collisions occur per second.
42 Due to the high concentration of energy on collision, heavier particles are created and
43 decay immediately to lighter particles. The experiments measure the tracks and properties
44 of the decay products that pass through the different detectors. The raw readout of those
45 is forwarded to a computer farm, where events of interest are marked and kept whereby
46 the rest of the events are not recorded. This reduces the stream of data to a frequency
47 that allows it to be stored on persistent storage. From there it can be retrieved and used
48 for further, offline analysis.

49 Performing a full analysis to measure physics observables from the data involves several
50 steps. This includes among others cleaning the samples by applying selection criteria,
51 reweighting to correct for systematic effects or creating new features that better describe
52 the event. The sample can then be used to directly infer unknown parameters by using
53 physically motivated models and performing a fit to the data. All of these analysis steps
54 require convenient, reliable and fast libraries together with enough computing resources.
55 To accomplish this, a lot of code is being written and stacked upon each other. To cope
56 with the ever increasing amounts of data, both in real-time event filtering as well as offline
57 analysis, it is mandatory to keep the computing level all in all at the state of the art.

58 Computing is still a comparably young, fast moving field. Hundreds of general
59 programming languages exist, most of them do not stay for a long time or only exist
60 in a specialist community. Even longer lived languages have both advantages as well as
61 shortcomings. This leads to different fields adopting a few, or even one, main languages
62 that specifically well serve their purpose. The field of scientific computing mainly involves
63 either simulation of systems or data analysis. In both cases languages that are fast
64 to do number crunching are required. Among the most popular languages for heavy
65 computations are Fortran and C/C++. The former is over six decades old and still used
66 up to these days. It was designed for numerical processing and contains optimizations that
67 still outperform other languages. C and its superset C++, the most popular language in
68 HEP, date back about four decades and are built for more general usage than Fortran
69 is. Although being fast and a powerful general programming language, it does not offer
70 the convenient abstractions that scripting languages offer. It allows but also requires to
71 manually handle certain resources, such as memory allocation, and is limited in terms of
72 flexibility because it is a statically compiled language. While the latter feature allows for
73 high performant execution, interpreted scripting languages such as Python offer additional
74 comfort and flexibility. While the execution of actual Python code can be significantly
75 slower than comparable static compiled languages when it comes to pure number crunching,
76 the huge Python package ecosystem offers a lot of libraries that implement time consuming
77 mathematical operations in a more efficient language, such as Fortran or C++.

¹Natural units with $\hbar = c = 1$ are used throughout.

78 makes Python a higher level library that abstracts away the handling of computation
79 demanding operations through external function calls. Together with an especially clean
80 syntax style, Python code is expressive and natural to read, giving in only a small penalty
81 for the performance from the overhead of the external calls.

82 This combination and a fast growing open source-community has established Python
83 as the most popular language for data analysis. With the recent advances in Machine
84 Learning and the rising popularity of Big Data analysis among industry leaders, the size and
85 quality of the scientific Python ecosystem has made a huge leap forward. Topics like deep
86 learning, which require highly optimised code due to the abundance of vectorised matrix
87 multiplications, have lead to the appearance of frameworks designed for this kind of massive,
88 parallel computations and supported by large companies such as Google’s TensorFlow [1]
89 (TF) or Facebook’s PyTorch [2]. With large economical interests coming into play, these
90 frameworks also focus on the efficient use of specialized hardware, such as Graphical
91 Processing Units (GPU) which are by design optimized for vectorized computations and
92 therefore fit the need of the Deep Learning community. These frameworks are optimized
93 both in terms of performance as well as in ease of use, dealing with the burden of
94 incorporating the parallelization.

95 Next to all these developments there is also a trend inside the HEP community to move
96 towards a more Python-oriented software stack. In recent surveys [3], its usage surpasses
97 C++ in collaborations such as CMS. The existence of the scientific Python ecosystem
98 offers the possibility of sharing some of the effort with the data analysis industry and
99 open-source community, allowing to perform a significant number of the analysis steps in
100 HEP within Python out of the box. This leaves to the HEP community only the burden
101 of developing the field-specific tools required to fit into the ecosystem. While some of the
102 existing frameworks in C++ offer Python bindings, they are usually not well integrated
103 with the Python language and the whole ecosystem. Several model fitting libraries in pure
104 Python have already been developed in order to fill parts of this gap, though none of them
105 offers the complete feature set that would be desired for HEP analysis and are hard to
106 extend. Therefore, while large advances have been made on this front, a viable alternative
107 to the existing, mature model fitting libraries in C++ is still missing. Nonetheless, some of
108 these libraries have proved the feasibility of using deep learning frameworks as computing
109 backends for model fitting.

110 Summarizing, HEP has

- 111 • the need for scalable, flexible model fitting;
- 112 • a strong movement towards Python with its huge data analysis ecosystem;
- 113 • the lack of a sufficiently strong model fitting library in pure Python.

114 Furthermore, modern high performance computing frameworks from deep learning
115 arose and their feasibility for computational backends in model fitting was demonstrated
116 in several projects.

117 With these ideas in mind, the `zfit` package has been developed with the goal to provide
118 this need by creating a pure Python based library built on a deep learning framework.
119 This requires the formalisation of the fitting procedure, the establishment of a stable API
120 and the usage of current knowledge from similar libraries. The following Section 2 will

121 expand on model fitting in HEP incensing a discussion on already existing libraries. With
122 this knowledge, the usage of `zfit` and its basic concepts, the formalisation of the model
123 fitting workflow, and the choice of the backend and its capabilities are outlined in Sec. 3.
124 The individual components of `zfit` will be discussed in more detail in Sec. 4. Afterwards,
125 the performance and scalability is evaluated with examples in Sec. 5. The extension of
126 `zfit` to more than the default model fitting is discussed in Sec. 6 by using its capabilities
127 to implement an amplitude fit. Lastly a brief overview of the future plans of the `zfit`
128 library and its ecosystem are given in Sec. 7.

129 2 Model fitting

130 In HEP, observations can be quantified by mathematical models which originate from
131 hypotheses or theories and make assumptions about the underlying behaviour of nature.
132 Often, these models have free parameters that we want to measure. A single model may
133 describes only parts of the observations and combinations and compositions of models may
134 be needed to build a model that describes the full data sample. Creating these models in
135 a convenient and correct way and finding the values of the parameters to maximise the
136 agreement with respect to the data is what “model fitting” refers to.

137 2.1 Maximum Likelihood

138 At the very heart of model fitting is the need to quantify the agreement, or rather the
139 disagreement, of a model with the data. This function of the parameters and data is
140 known as the loss. It is the very definition of the problem and mathematically fully defines
141 the solution. In HEP analysis losses are mostly based on the likelihood of the model under
142 the data, whereby the model typically depends on free parameters. In the following, an
143 introduction to the method of maximum likelihood is given. A more detailed explanation
144 and derivation can be found in Appendix A.

145 A likelihood can be defined by the following: given a model parametrised by θ and a
146 dataset x , the likelihood describes the odds that an event happened under θ

$$\mathcal{L}(\theta) = P(x|\theta). \quad (1)$$

147 The likelihood as shown in Eq. 1 is the quantity to be maximised in order to achieve
148 the maximal $P(\theta|x)$. To build this likelihood, we need the model $f_\theta(x)$ to be a probability
149 density function (PDF), i.e. it’s normalised to 1. Especially in HEP, it is often the case
150 that the PDF is zero outside of certain boundaries, for example because points outside a
151 specified domain are removed, in which case

$$\int_l^u f_\theta(x) dx = 1, \quad (2)$$

152 where l and u define the lower and upper boundaries of the domain, respectively. This
153 also extends to higher dimensions. It follows directly that any function $g_\theta(x)$ ² can be
154 normalised and therefore used as a PDF $f_\theta(x)$

$$f_\theta(x) = \frac{g_\theta(x)}{\int_l^u g_\theta(x) dx}. \quad (3)$$

155 A likelihood can be a product of likelihoods of independent events

$$\mathcal{L} = \prod_i \mathcal{L}_i, \quad (4)$$

156 and therefore the likelihood of dataset x can be written as the joint probability of each
157 event

$$\mathcal{L}(x|\theta) = \prod_i f_\theta(x_i),$$

²This is about the small subset of modelling functions in physics *without* pretending mathematical correctness in a general way. This includes functions $f : \mathbb{R}^n \mapsto \mathbb{R}$ that are positive, l^1 and (piecewise) C^1 .

158 with x_i a single event from the dataset x .

159 The calculation of $\mathcal{L}(x|\theta)$ involves the product of many small numbers, which is not
160 possible to perform using a normal computer given its limited precision. To solve this issue,
161 a log transformation can be applied. In addition, the log-likelihood is usually negated,
162 thus changing the target of finding the maximum to finding a minimum and ending up
163 with a negative log likelihood (NLL).

164 A maximum likelihood estimate using the transformation above is therefore given by
165 finding the minimum of the NLL

$$NLL = - \sum_i \ln(f(\theta|x_i)) \quad (5)$$

166 This maximises therefore the agreement between data and model, i.e. the *probability*
167 *of the model given the data*.

168 As seen in Eq. 4, the combination of likelihoods is quite versatile and not only limited
169 to a model shape matching the data shape. Often, a combination of several of the following
170 likelihoods is built

171 **Simultaneous** Multiple models can share parameters. To fit them simultaneously to
172 different datasets, their likelihoods can be combined (summed).

173 **Extended** While a PDF is normalised, we can add an absolute scale as an additional
174 term to the likelihood to reflect the number of events contained in this model. Given
175 the data, we know the number of events and can add a Poisson term to account for
176 them.

177 **Prior** For some parameters, a prior distribution is known. This describes the knowledge
178 obtained from other measurements and influences the likelihood if the parameters
179 spread is in the same order of magnitude as the sensitivity of the fit to the parameter.
180 A prior, or constraint, is a probability depending directly on the parameter value
181 and can also be added to the likelihood.

182 Regardless of the complexity of the model, we end up with a single number, the loss, that
183 can be used to compare the agreement between different models or parametrizations and
184 the data. When fitting a model, the loss is minimised by adjusting the parameters. While
185 the absolute value of the loss is usually not important, the ratio of losses from different
186 models can often be useful in further statistic tests.

187 2.2 Requirements

188 Some features are crucial in order to implement a model fitting library. An important part
189 of model fitting is the model building itself, but a library should also offer a convenient,
190 transparent creation of the loss and the minimisation. Especially in HEP, the following
191 features are essential:

- 192 • PDFs are by definition normalised over a certain range. In most other libraries and
193 fields, the domain is assumed to be $(-\infty, \infty)$. In HEP, this is basically never the
194 case and a finite normalisation range is used.

- 195 • Fits in HEP are often more than one-dimensional. The framework should therefore
196 naturally extend to higher dimensions.
- 197 • Building and combining models from basic shapes like Gaussian or exponential
198 functions only suffices for simpler cases, but this is often not enough to build more
199 complicated or specific models. Therefore, a convenient way to implement custom
200 models has to be provided.
- 201 • Reasonable scaling with the data size and the model complexity is a key criteria. This
202 is often especially hard to achieve in combination with the ability of specifying custom
203 models, since the latter usually requires to have the parallelization implemented by
204 the user.
- 205 • While the minimisation of the loss yields an optimal value for each parameter, it is
206 crucial in HEP to also know the uncertainty of the value. This requires the library
207 to have a transparent way of handling the parameters and their uncertainties as
208 well as to provide the flexibility to perform an advanced statistics treatment.

209 2.3 Existing libraries

210 Model fitting itself is nothing new. In fact there are already a lot of model fitting libraries
211 available. Some of these libraries are also written in Python and cover a similar scope
212 as `zfit`. Building a new fitting library from scratch sounds therefore like reinventing
213 the wheel and should be avoided if not necessary. But as already discussed in Sec. 1,
214 fundamental changes in the computing architecture are leading to vectorized paradigms.
215 Additionally, the needs in HEP for larger and more complicated while still flexible fits
216 require to keep up with the state-of-the-art in computing. And this sometimes requires a
217 re-invention.

218 However, it is an imperative to make sure that no existing library already fulfils the
219 needs or can be extended to. And even if concluding that a new library is the way to go,
220 as much as possible should be learned and taken from any existing library in order to
221 reinvent as few as necessary. In the following an overview of already existing libraries is
222 given.

223 2.3.1 General fitting

224 Fitting models to data is a task that is performed in a variety of fields independent of HEP.
225 Different general fitting libraries exist in Python, but they often contain functionality
226 not actually needed in HEP, such as mean, variance, survival function, and lack central
227 features like a custom normalisation range or the extension to more than one dimension.

- 228 • Scipy [4] is the go-to library for scientific calculations in Python and provides an
229 extensive toolbox for statistical and numerical methods. There is a module with
230 distributions that have proven to be stable and work well. Downsides of the package
231 include a non-optimized implementation in terms of parallelisation and lack of
232 support for composite models.
- 233 • `lmfit` [5] shares a lot of its design in terms of naming and concept to `zfit`. It is
234 built for model fitting, has parameters, minimisers, fit results and more. It lacks

235 more advanced features like the possibility of normalisation ranges for PDFs or good
236 scalability, since it is built on top of numpy, a fast numerical library in Python, and
237 scipy, which strongly limits the ability for massive parallelisation.

238 • TensorFlow Probability [6] provides a library for statistical reasoning. Its focus
239 is on analytical functions and only marginally extends to numerical and Monte
240 Carlo methods, which limits its application to analytically integrable functions.
241 Interestingly, it contains a lot of of features that can be used inside or together with
242 `zfit`, such as Bayesian inference with MCMC sampler and analytic functions with
243 integrals already implemented in TF.

244 2.3.2 HEP specific

245 A wide range of specialised fitters exist in HEP. The overview here is limited to general
246 purpose fitters which can be used from Python.

247 `ROOFIT` [7] is the de-facto standard tool for fitting in HEP. Models are built using
248 classes and provide automatic normalisation and integration. There is support for
249 binned as well as unbinned fits. `ROOFIT` itself extends beyond that and offers also
250 an extensive plotting and statistics module. While the library has proven itself in
251 numerous analyses over the years, and the model building part of `zfit` is actually
252 inspired by the core of `ROOFIT`, there are several shortcomings which are meant to
253 be addressed with `zfit`:

- 254 • `ROOFIT` is not a native Python library but can only be accessed through the
255 Python bindings to `ROOT`. Since `ROOFIT` manages its own memory in C++
256 and Python uses a garbage collection as well, this can lead to memory leaks
257 and completely undefined behaviour.
- 258 • Since the Python interface is barely a wrapper around the C++ classes, it does
259 not integrate well to the scientific Python stack.
- 260 • In terms of flexibility, `ROOFIT` offers the possibility to be extended up to a
261 certain degree with custom classes in pure C++. But especially when used
262 from Python, it does not provide a convenient way to define custom PDFs.
- 263 • While there are improvements in the pipeline, it is not natively optimized to
264 run vectorized on multiple cores or even accelerators like GPUs.
- 265 • Since the usage requires `ROOT`, the installation and setup is typically not
266 lightweight.

267 `probfitt` [8] is a fitting library written in pure Python that mainly uses Cython to perform
268 the heavy computations. This is a limitation in terms of performance and custom
269 PDF implementations that makes a possible extension hard. Since it does provide
270 limited features only, a large extension would be needed together with a major
271 conceptual overhaul to be able to include new features.

272 `pyhf` [9] is a re-implementation of `HistFactory` from `ROOT` in Python. It makes use of
273 TensorFlow and other libraries including PyTorch and Numpy as a backend. It is
274 purely designed to do binned template fits and does not extend its functionality
275 beyond that point.

276 The CMS Combine Tool [10] contains a subpart that implements template fits in TF. It
277 does not extend its functionality further and is currently not available as a stand-
278 alone package. Several useful parts like likelihood profiling or a minimiser in pure
279 TF have been implemented there.

280 TensorFlow Analysis [11] is a library with a simple, functional approach to built the loss
281 with TF and use Minuit [12] directly inside³ to find the minimum. It offers a lot of
282 physics content to create a model. While the lightweight approach comes with a
283 lot of flexibility, the library also leaves quite some work to the user. For example it
284 does not offer anything close to model composition with automatic normalisation.
285 Notably, in its current state, the library lacks Python 3 support. However, its
286 importance has to be stressed since it demonstrated the feasibility of using TF for
287 unbinned likelihood fits with complex models and was a major inspiration for the
288 development of `zfit`.

289 **TensorProb** is a model fitting library in Python that uses TF as the backend. In general
290 it was built with a similar goal in mind as `zfit`, providing a model fitting library in
291 Python using TF, but using more an experimental approach. It offers models that
292 can also integrate and sample. The content is based on older TF versions and the
293 library is strongly limited in functionality. Most importantly though, the project
294 never grew out of its experimental status and has been discontinued. It recommends
295 now to use `zfit` instead.

296 While the discussed model fitting libraries have different strengths and weaknesses,
297 no single one fully fulfil the needs of HEP. However it is worth pointing out that their
298 demonstration of concepts, designs and even certain functionality that can be used directly
299 with `zfit` are essential pieces in the development of `zfit`.

³This also requires to have the ROOT package installed.

3 zfit introduction

`zfit` was created in order to fill the gap of a model fitting library in pure Python for HEP. We will now have a look at it, how it is structured and supposed to fill this gap. Model fitting as implemented in `zfit` is split into five essential parts. To introduce them and `zfit` itself, an example with a sum of a Gaussian and an exponential PDF will be implemented. This example can be thought of as a fit to an invariant mass distribution of a signal and a background component.

Let's assume we are interested in an observable x within a range from 5 to 10. In `zfit` this is expressed with a `Space` defining our domain

```
limits = zfit.Space(obs="x", limits=(5, 10))
```

`zfit` can handle data from a variety of different sources. In this case, we load the data `data_np` from a numpy array

```
data = zfit.Data.from_numpy(array=data_np, obs=limits)
```

Since the data was specified with the `limits` as its observables, it will be cut automatically to be only within the `limits` range. In this context, we can think of the observable x as of the column of the data frame.

Next the model needs to be built. We create the Gaussian PDF with two free parameters, `mu` and `sigma`. Using 7 and 1.5 as initial values, respectively, this is done as

```
mu = zfit.Parameter("mu", 7)
sigma = zfit.Parameter("sigma", 1.5)
```

Creating the Gaussian in the observable x and using the parameters from before as

```
gauss = zfit.pdf.Gauss(obs=limits, mu=mu, sigma=sigma)
```

Equivalently the exponential PDF is created. A fixed value of -0.1 is used for the exponential parameter λ as in $e^{\lambda x}$ and can directly be given to the PDF⁴

```
exponential = zfit.pdf.Exponential(obs=limits, lambda_=-0.1)
```

To build the sum, an additional free parameter is used to describe the fraction of the first PDF. It is initialised with 0.5 and limited between 0 and 1

```
frac = zfit.Parameter("fraction", 0.5, 0, 1)
model = zfit.pdf.SumPDF(pdf=[gauss, exponential], frac=frac)
```

Now that the model is built, we can define the loss by combining it with the data. Here an unbinned NLL will be used

⁴Alternatively, a `Parameter` with the argument `floating` set to `False` can be created.

```
329 nll = zfit.loss.UnbinnedNLL(model=model, data=data)
```

330 which needs to be minimised in order to find the optimal parameters. To achieve this goal,
331 a minimiser such as Minuit is needed. Once it is created, we use its `minimize` method in
332 order to minimise the previously built loss

```
333 minimizer = zfit.minimize.MinuitMinimizer()  
334 result = minimizer.minimize(nll)
```

334 The outcome of this minimisation is stored in a `FitResult` object. Whether the
335 convergence was successful can be checked with

```
336 has_converged = result.converged
```

337 The free parameters of the model are updated in-place with the values obtained from
338 the minimisation. This implies that the shape of the model has changed now, since it
339 depends on the parameters. While a parameter can change again, the `result` stores the
340 values from the minimisation as immutable numbers. They can be accessed like⁵

```
341 mu_value = result.params[mu]["value"]
```

342 The value of the parameter is incomplete without an estimate of its uncertainty. For
343 an accurate estimation, we can use `error`, an advanced method that takes all correlations
344 among the parameters into account

```
345 errors = result.error()
```

346 This simple yet complete example demonstrates how model fitting in `zfit` works. All
347 parts contain a lot more functionality than just seen, but the structure of the workflow as
348 shown in Fig. 2 into five independent parts remains the same, no matter how complicated
349 a fit may be.

350 **Model building** The construction of models is the core of `zfit` and involves Functions
351 and PDFs. The difference between them is that the latter is normalised to one
352 over a certain domain. Building the model includes a set of convenient base classes
353 that allow to easily create a custom model as explained in Sec. 4.3. Furthermore,
354 composed models involving sum, products and more are available.

355 **Data** Any kind of data needs to be loaded and transferred into a well defined `zfit` format.
356 The `Data` class takes care of this and offers several formats to load from, which then
357 can be used by models. The aim here is to provide a simple way of loading data
358 from different formats into `zfit` and applying some cuts.

⁵Notice that `mu`, the parameter itself, and not `"mu"`, the name of it, is used as the key.

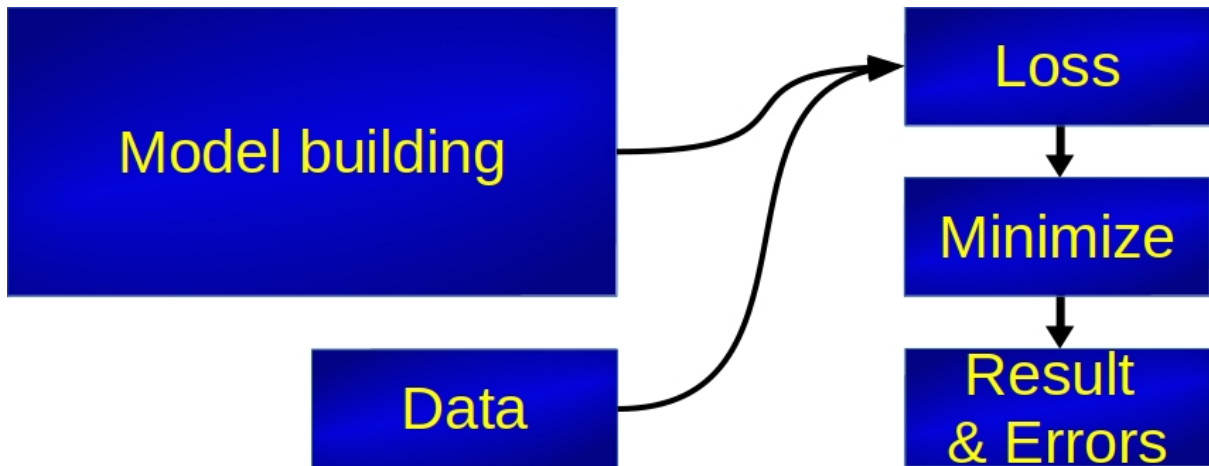


Figure 2: Fitting workflow in `zfit`. Model building is the largest part. Models combined with data can be used to create a loss. A minimiser finds the optimal values and returns them as a result. Estimations on the parameters uncertainties can then be made.

359 **Loss** This is the core definition of the problem. It uses the model and data objects to
 360 calculate a single number that quantifies the discrepancy from the model and the
 361 data. Typically, a binned or unbinned NLL or a χ^2 is used, but `zfit` offers the
 362 freedom to implement any desired loss that is not available in a straightforward way.
 363 From this step onward, it is irrelevant what data or models are *actually* used. Only
 364 the number and the gradients with respect to the models parameters matter.

365 **Minimisation** Given a loss, the minimiser minimises its value with respect to the free
 366 parameters of the models. In `zfit`, several algorithms are implemented by wrapping
 367 existing minimisation libraries.

368 **Result and Errors** After each minimisation, a `FitResult` object is created. It stores
 369 all the information about the minimisation process and allows ,amongst other things,
 370 to check if the convergence was successful. The result also includes the parameters
 371 and their values at the minimum. Furthermore, the loss and the minimiser itself are
 372 also stored in the result. Using both of them, an estimation on the uncertainty of
 373 the parameters can be made. For this purpose, some simple algorithms are provided
 374 by `zfit`, but any more sophisticated uncertainty estimation can be made by using
 375 the objects made available by the `FitResult`.

376 This formalisation is a powerful approach: the separation of the model fitting into
 377 these fine building blocks allows to improve and maintain the individual parts almost
 378 independently. Most importantly, it reveals a surprising similarity to the field of deep
 379 learning: apart from the last step, the workflow is *exactly* the same. Using a deep learning
 380 framework as the backend for a model fitting library therefore seems like an obvious choice
 381 to consider.

382 3.1 TensorFlow backend

383 Deep learning has recently gained a lot of attention as it has been quite successful as a tool
 384 in big data analysis and predictive statistics. In its core, the idea is to extract correlations

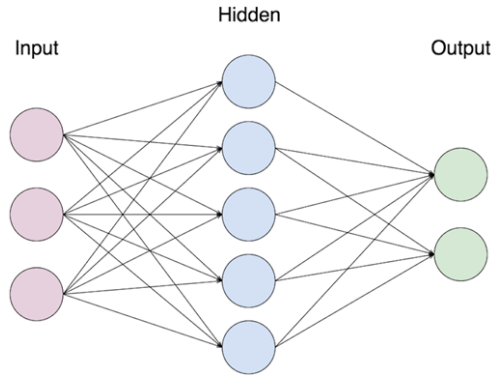


Figure 3: A schematic view of a DNN function. Input takes the data and has the dimension of an event. Each end-point of a line is multiplied by a weight, a free parameter, and added with the other lines. At a node, a non-linear function is then applied.

385 from data by training a neural network on them in order to make accurate predictions on
 386 unknown samples. This can be the classification of images, the prediction of stock markets
 387 etc. More interestingly, the typical deep learning workflow can be summarised also by
 388 Fig. 2 by replacing “model” with “neural network”,⁶ “minimisation” with “training”
 389 and removing the last block “Result & Error”. Moreover, deep learning and HEP model
 390 fitting both use large data samples and build complex models. This similarity inspired
 391 the implementation of `zfit` with a deep learning framework as the backend.

392 While we have just shown how the two workflows look incredibly similar at first
 393 glance, there are some hidden, crucial differences. Knowing them is essential in order to
 394 understand the advantages but also limitations of this approach. In the following we will
 395 have a simplified at the core of deep learning and compare then to model fitting.

- 396 • A Deep Neural Network (DNN) is simply a function $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$. The input space
 397 is the number of observables of a single event. The output of the function, the
 398 prediction, is notably m dimensional. Contrary to this, a model with PDFs outputs
 399 into the one dimensional space of a normalised probability. DNN outputs, if used
 400 for classification, correspond to pseudo-probabilities. While they are not normalised,
 401 there exists a monotonic transformation to a probability.

402 *Consequence:* normalisation over a certain range is specific to model fitting and no
 403 explicit tool, e.g. for numerical integration, exists in the deep learning frameworks.

- 404 • In model fitting, the composition of shapes is motivated by previous knowledge
 405 and an underlying theory is built. There is often a meaning behind each part and
 406 the shape of the model is restricted to a specific problem. This specific shape,
 407 coming from assumptions and previous knowledge, is what keeps the number of
 408 parameters low: typically, no more than a dozen for simple fits and a maximum of a
 409 few hundreds for the most complicated fits are used. Contrary to that, the structure
 410 of a DNN is basically agnostic to the problem and depends mostly on its complexity.
 411 It therefore incorporates a minimum of pre-knowledge and assumptions about the
 412 correlations in the data. This huge versatility is what makes deep learning such a

⁶Neural networks in deep learning are also called models. The term “model” will here solely be used for “classical” model fitting as in `zfit`.

413 successful field but comes at a price: a great number of free parameters is needed,
414 starting at thousands for really simple DNNs, typically being around hundreds of
415 thousands and going to tens of millions. DNNs are a structure consisting of layers
416 with nodes as shown in Fig. 3. Each of the nodes adds an additional parameter for
417 every incoming connection, resulting in a large number of parameters.

418 *Consequence:* While dozens of DNN building libraries like Keras⁷, PyTorch and
419 more offer great capabilities in building DNNs, their abstractions into layers are of
420 no use to model fitting.

- 421 • A DNN is essentially matrix multiplications scaled by the free parameters, with
422 an additional simple, non-linear activation function applied. In model fitting, the
423 shape can have an arbitrary complexity and contain a whole range of elementary
424 functions. Furthermore control flow elements and complex number are often used
425 as well. Not only is the function itself more complicated, also the dependency of
426 parameters can be highly non-trivial.

427 There is an additional difference in the precision of the floating point operations. In
428 model fitting, the precision required is higher, because the values of likelihoods and
429 the changes are larger compared with neural networks often having values varying
430 between -1 and 1 . Also the Quasi-Newton methods as described below build an
431 approximate second order derivatives which needs a high enough precision.

432 *Consequence:* While both fields do heavy computations, the focus of the optimizations
433 in a fitting library is slightly different and requires for example to always explicitly
434 specify float64 as data type.

- 435 • Minimising a loss is a non-trivial task. Algorithms usually start at a certain point
436 and use local information to make forward steps. The gradient and sometimes
437 higher order derivatives, usually up to the second order, are used to help finding
438 the minimum. In particular, which order is usable strongly depends on the number
439 of parameters: the Hessian matrix of n parameters has n^2 entries rendering its
440 calculation unfeasible for more than a few hundred of parameters; this restricts the
441 minimisation of DNNs to only use the first order derivatives at the cost of more
442 required steps.

443 *Consequence:* On one side, minimisers designed for DNNs and optimised to work
444 with the framework are not suitable for model fitting. On the other side the analytic
445 gradients that are provided by th frameworks for their minimisers can be easily
446 extended to higher orders and come in very handy for model fitting minimisers.

- 447 • Fitting a model has the goal to find the parameters for which the model matches
448 the data best. In terms of a loss function, this is equivalent to finding its *global*
449 minimum. Being stuck in a local minimum is a problem and requires careful
450 treatment. Contrary in DNNs the global minimum is not found but also not desired.
451 The DNN has to approximate an arbitrary data sample *good enough* and a local
452 minimum is usually found, it is in fact preferred over the global minimum since, due
453 to the high degrees of freedom of a DNN, this typically entails a huge over-fit⁸ and
454 a bad generalization.

⁷Keras is an API specification only, a reference implementation exists.

⁸Basically remembering every noise in the data.

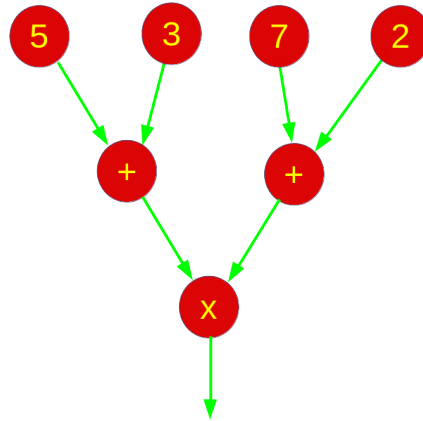


Figure 4: Example of a graph representing $result = (5 + 3) * (7 + 2)$.

455 *Consequence:* While finding the global minimum is crucial in model fitting, deep
 456 learning is interested to find a local minimum.

- 457 • Finding the global minimum in the model fitting case involves the evaluation of the
 458 loss over the whole data sample at every step. For the training of a DNN, variations
 459 of a technique called stochastic gradient descent (SGD) are used: they work by only
 460 evaluating a small mini-batch of the data, typically around 32 events, and then
 461 taking a step towards the negative gradient direction.

462 *Consequence:* Both fields are pushing the limits of handling big data and samples
 463 too large for the memory are common. Deep learning is optimized to loop through
 464 a data set with very small batch sizes to take a minimisation step but to do that
 465 millions of times. Model fitting needs to loop through the whole data sample *once*
 466 for a single step but no more then a few thousand times.

467 As seen, while there are differences, the core of the problem is still the same: build a
 468 complicated model, use data to get a value from it and tune parameters to optimize the
 469 loss.

470 To accomplish this task efficiently, deep learning frameworks use a declarative paradigm
 471 by building a computation graph as seen in Fig. 4. This allows to perform optimisations
 472 and define the parallelisation *before* the actual execution of the computation. Furthermore,
 473 it also allows to get an analytic expression for the gradient by consecutively applying the
 474 chain rule within the graph.

475 To implement model fitting in `zfit`, the TensorFlow library was chosen, restricting
 476 the implementation to static graphs as explained in more detail in Appendix B.2. The
 477 main motivation for this decision comes from the fact that models in model fitting *usually*
 478 don't change their logic but are rather built once and then minimised. While the same is
 479 true for most DNNs, more advanced fields in deep learning like reinforcement learning can
 480 heavily rely on dynamic models. The main advantages of a static graph are the additional,
 481 potential speedup and the immutability, which leads to less unexpected behaviour. The
 482 restriction that *anything built will remain like that and not change* allows for obviously
 483 more efficient optimizations compared to a graph where any part could change anytime
 484 and a re-analysis of the graph is required.

485 Working with graphs leads to difficulties and unexpected behaviours in comparison
486 with more traditional, non-graph based code, such as the one used in ROOFIT. In `zfit`,
487 most of these complications are hidden from the user and the library offers a similar user
488 experience as the one found in other model fitting libraries. This requires some extra care
489 taking behind the scenes, as explained in Appendix B.2.

490 4 `zfit` implementation

491 *Things should not be easy or hard. But consistent.*

492 In Sec. 3 an example of a complete fit with `zfit` was shown. While this was a rather
493 simple example, `zfit` generalizes to more complicated, custom model or higher dimensional
494 fits in a consistent manner. The guiding principles that facilitate this generalization are

495 **One thing** The library is meant to do one thing: model fitting. Additional capabilities
496 like advanced plotting, statistics etc. are intentionally left to other libraries.

497 **Simplicity** As often stated, easy cases should be easy or at least straightforward and not
498 hard. For simple, one dimensional models and fits, no extended knowledge about
499 advanced concepts or pitfalls should be necessary. Making mistakes can only happen
500 when clearly leaving the known grounds.

501 **Consistency** When going from the simple to the more complex case, it should not be
502 necessary to learn a whole new set of rules and behaviour. Instead, enough generality
503 should be contained also in the simple case, even if it means to slightly heighten
504 the knowledge required for it. The overall difficulty should gradually increase, and
505 expert knowledge should only be required for expert problems.

506 **Flexibility** No library can ever cover everything, so special cases will need to be imple-
507 mented specially, but most elements of the library should nonetheless be usable and
508 only the special part has to be implemented specially.

509 These requirements can be conflicting. Especially flexibility and simplicity are hard to
510 reconcile. A conflict arises with the potential of creating bugs: since not every case is the
511 most complicated one, not every user will know about all peculiarities of the library and
512 may uses the freedom in the wrong way. To avoid mistakes due to the lack of knowledge
513 of the library but also incorporate the flexibility, `zfit` follows the pythonic way, a term
514 used to describe the philosophy behind Python and its usage.

515 While the Python language itself, as any other language, consist of syntax specifications,
516 a non-neglectable part of Python consists of its philosophy, the Zen of Python. A commonly
517 used phrase in this context is the expression “*We’re all adults here*”. It refers for example
518 to the absence of enforced private class attributes in Python. In its essence it means that
519 anyone should be able to change anything, which offers a great flexibility.

520 A complementary strong guidance is given by “*There should be one – and preferably
521 only one – obvious way to do it*”. This encourages to discuss problems appearing in the
522 Python community and finding one “best” way to solve a specific problem. And even if
523 there is no “best” solution, at least a convention should be agreed on. Combining these
524 two ideas, `zfit` provides a lot of freedom and allows to implement any part of the fit flow.
525 To counter unintentional mistakes, clear guidelines and examples on *how* to accomplish
526 something specific are provided. This is especially important for the simpler cases, in
527 which one clear way is shown.

528 In the following, we will go through a few essential pieces of the library that allow to
529 extend the workflow naturally to more complex fits and discuss them in more detail.

530 4.1 Spaces and Dimensions

531 The extension to more than one dimension when fitting a model contains some ambiguities.
532 The data, limits and the models all have corresponding axes or columns that should
533 match. To deal with that in a simple, yet consistent manner, `zfit` has a `Space`. The
534 responsibility of this class is to define and handle the axes or dimensions and limits. To
535 understand the concept and the `Space` itself, three definitions need to be made:

536 **Observables** The observables are the *named axes of a coordinate system*. A single
537 observable is a string and a list of strings acts as several observables describing a
538 higher dimensional system. In the context of data this is equivalent to columns
539 in a data frame. Observables allow a named, *inter-object* identification of axes.
540 Therefore, working with observables allows to work independently of the underlying
541 data ordering.

542 **Axes** Axes are integers and are used to *enumerate the axes of a coordinate system*. This
543 corresponds to indices of an array and provides the fully order-based mapping
544 necessary for *intra-object* manipulations. For example, a `Data` with three columns
545 has three axes, 0, 1, and 2, which can though be reordered so that the corresponding
546 observables match the order of some other observables.

547 **Limits** A description of boundaries that can be used to define any kind of limits of the
548 axes. Currently only rectangular limits are supported but arbitrary shaped limits
549 will be provided in the future.

550 A `Space` can be initialized with observables and limits to define a domain. When it's
551 assigned to an object, it automatically connects the axes of the object with the observables
552 from the `Space`. More details on the implementation and use cases as well as additional
553 functionality for dimensional handling can be found in Appendix C.1.

554 4.1.1 Limits

555 Limits are used in many instances, be it in sampling limits, integration or data limits. A
556 `Space` not only defines the observables but typically also has limits associated with it. In
557 one-dimensional fits, limits as seen in the example in Sec. 3 are needed. Simple limits
558 consist of a tuple for each observable with lower and upper limits. For example a `Space`
559 in one observable `x` from `-5` to `3` can be created like

```
560 limits = zfit.Space(obs="x", limits=(-5, 3))
```

561 This is the simplest way of specifying limits and rather a special case. For anything
562 more sophisticated, such as multiple limits or multiple observables, either a composition
563 of `Spaces` or the more general format as explained in Appendix C.2 has to be used. For
564 example when blinding a region, a `Space` with multiple limits can be used.

565 While the general format is fully specified independent of a `Space` and can therefore
566 be useful programmatically, a `Space` with multiple limits can be built unambiguously
567 from `Spaces` with simple limits by adding them, either through a dedicated `zfit` function
568 or using the addition operator in Python. In this way, multiple limits can be created
569 through simple composition and without the need of using the more general format.

570 As an example, let's assume a `Space` should be created with the observables `x`, `y` in
571 the two domains l_01 and l_23

$$l_01 = \{(x, y) | x_0 < x < x_1, y_0 < y < y_1\}$$
$$l_23 = \{(x, y) | x_2 < x < x_3, y_2 < y < y_3\}.$$

572 We start out creating the domains by specifying the limits in the `x` observable

```
573 limit_x_01 = zfit.Space(obs="x", limits=(x0, x1))
limit_x_23 = zfit.Space(obs="x", limits=(x2, x3))
limits_x = limit_x_01 + limit_x_23
```

574 Equivalently `limits_y` can be composed. Since going to higher dimensions is unam-
575 biguous with two limits in each space, this can be done using the multiplication operator
576 in Python or the function `combine`.⁹

```
577 limits_xy = limits_x * limits_y
limits_yx = limits_y * limits_x
```

578 The difference between `limits_xy` and `limits_yx` is the order of the observables. In
579 the first case, it's `["x", "y"]` while for the latter it's `["y", "x"]`. In order to ensure
580 consistency, if the two `Spaces` already have observables in common, the limits in this
581 observables have to be the same.¹⁰ Reordering the `Space` is possible as well as extracting
582 a subspace, a `Space` only defined in subset of the dimensions. More details can be found
583 in Appendix C.1.

584 4.2 Data handling

585 Data to fit can come from different sources but it should be handled uniformly inside `zfit`.
586 To ease this, the `Data` object is responsible for loading, ordering and simple preprocessing
587 of data, which can have weights assigned to it. Furthermore, this abstraction layer
588 with `Data` potentially allows for more advanced use cases such as batched, out-of-core
589 computations of the likelihood.

590 The `Data` class supports a variety of data files and structures. While adding additional
591 loading capabilities is not difficult, the focus is on the following formats: the default HEP
592 format `ROOT`¹¹, Numpy arrays and Pandas DataFrames, and pure Tensors. More details
593 can be found in Appendix C.3

594 Each `Data` has a `Space` with observables it is defined in. This assigns an observable
595 to each column of the data, so the observables here act like columns from spreadsheets
596 or DataFrames. This allows to retrieve a subset or different ordering of the `Data` by
597 specifying the observables explicitly in the method that returns the data as a Tensor.

598 Once instantiated, a `Data` object appears like a lightweight wrap of the Tensor class
599 and can be directly used as such. It is possible therefore to simply operate on a `Data`

⁹If a different number of limits were defined, an error would be thrown.

¹⁰This exact behaviour of the multiplication and observables is the same if models are multiplied.

¹¹Even though `ROOT` files are supported, the `ROOT` library is not needed thanks to the `uproot` package as explained in Appendix C.3.

600 object with any operation that would also accept a pure Tensor. While this is convenient
601 for certain contexts where the correct ordering of the data is guaranteed such as inside a
602 model, the preferred way of using the `Data` is to access the columns by names using the
603 `unstack_x` method. The `Data` class can also handle data generated on the fly and not
604 fitting into memory, see C.4.

605 `Data` objects can be ordered in-place as opposed to `Space`, the reordering of which
606 returns a new instance. This is heavily used together with a context manager inside
607 models if a `Data` is given as an argument in order to match the order of its observables
608 with the order of the models observables.

609 4.3 Model

610 Building models is the core competence of `zfit`. As seen in the example in Sec. 3, this can
611 be done in a simple manner by using already implemented models and possibly combining
612 them, but the models can also be completely custom built. Within `zfit`, there are two
613 basic types of models to cover most cases: Functions and Probability Density Functions
614 (PDF).

615 The basic features of a model include

- 616 • Each model is defined inside a `Space`. Its dimension are “observables”, simple string
617 identifiers as previously discussed.
- 618 • A model implements a function that returns its value depending on some data. This
619 is either `pdf` for a PDF or `func` for a Func.
- 620 • Full as well as partial integration over a model is possible. This includes numerical
621 as well as analytical integration, if available.
- 622 • Generating a sample following the models shape using numerical or analytical
623 methods, the latter only if available.

624 The value function as well as the integration and sampling are implemented to return
625 pure Tensors. Depending on the task, higher-level methods providing either a more
626 intuitive, imperative behaviour or a significantly more performant execution for certain
627 cases such as repetitive pseudo-experiments are also implemented.

628 The main differences between a PDF and a Function concern the normalisation and
629 the output dimensionality. This leads to a few subtle differences.

630 **PDF** A PDF $f(x)$ is only well-defined with a given normalisation range. This defines the
631 normalisation constant so that, with the limits from *lower* to *upper*, the integral
632 over the PDF equals to one as

$$\sum_i \int_{l_i}^{u_i} f(x) dx = 1 \quad (6)$$

633 with i indexing all the limits that make up the domain, l_i and u_i are the lower
634 respectively upper limit.

635 A PDF object has three special attributes, which are

- 636 • the probability density function `pdf`. It returns always a rank one Tensor, i.e.
637 a simple vector, with the length number of data points.
- 638 • the probability density function *without* the normalisation constant,
639 `unnormalised_pdf`. In cases where only the shape is needed, using this function
640 is less expensive, especially if the normalisation has to be computed numerically.
- 641 • the limits that define the normalisation constant of Eq. 6. They can be set
642 using the `set_norm_range` method.

643 **Function** A Function is in a way more simple and general than a PDF. It takes the same
644 values as a PDF but returns something with dimensionality \mathbb{R}^m . It can be used to
645 transform values or to use it as a building block for more complex expressions.

646 It has the method `func` to evaluate its value and no other special attributes.

647 4.3.1 Parametrization

648 Models can be parametrized by `Parameter` objects which can be used in the implementation
649 of the shape function like any other Tensor-like object when building the model. In the
650 following we will first have a look at the `Parameter` itself. Afterwards we will see how
651 they are exactly used with a model.

652 There is a distinction between dependent and independent parameters as only the latter
653 can be changed directly and have limits while the former are any arbitrary combination
654 of them. An independent `Parameter`

- 655 • has a name with purely descriptive purpose;
- 656 • has an initial value;
- 657 • maybe has lower and upper limits;
- 658 • is either floating or not, independent of whether limits were specified;
- 659 • has a step size indicating the order of magnitude of the parameter which can be
660 given. Otherwise, it is automatically inferred. A well chosen step size improves the
661 minimization process and can be critical, mostly in the absence of limits and with a
662 weak dependence on the model, so that large steps will be required to change the
663 model.

664 Currently, the shape of parameters is implicitly restricted to a scalar.¹² As a conse-
665 quence, a parameter cannot simply be treated like data as it is possible in `ROOT`. A
666 function `Parameter` which depends on the data itself will most likely be available in the
667 future.

668 The name of a parameter, as any other name for a single object in `zfit`, has purely
669 descriptive character. There is purposely¹³ no direct way provided to access parameters

¹²This comes from the fact that the PDF will have different sized data as input which is not controlled by the user, such as when doing numerical integration. Any parameter therefore has to be able to broadcast seamlessly.

¹³Contrary to the `ROOT` framework. If the need ever arises, adding this as an additional feature is relatively simple.

670 by name: instead of using names the actual parameter object is passed around. This has
671 the advantage of avoiding double bookkeeping the parameters, since then a reference on
672 an object as well as on a string would be needed. The name is mandatory for parameters,
673 as opposed to other objects in `zfit`, since matching the value of a parameter to the right
674 name is a critical task during and after a minimization, when reading off the right value.

675 Composed parameters are dependent parameters. In general they can be any Tensor,
676 that is a result of any kind of operation. They can depend on zero, one or more
677 independent parameters, as composed parameters are arbitrary functions; therefore,
678 operations such as shifting and scaling are included as a trivial subset. For more details on
679 the actual implementation and the dependency management see Appendix C.5. Composed
680 parameters can neither be floating nor have limits currently, since the independent
681 parameters a composed parameter may depends on can have arbitrary relations and have
682 to be restricted themselves. In order to restrict a parameter, arbitrary constraints can be
683 given to the loss instead.

684 Every model can depend on multiple parameters, both dependent and independent.
685 Each parameter that parametrises the model has a name specific to the model and is
686 given on instantiation. For example a `Gauss` has parameters named `mu` and `sigma` as in
687 the example in Sec. 3. They are stored in a mapping attribute named `params`.

```
688 mu_param = gauss.params["mu"]  
sigma_param = gauss.params["sigma"]
```

689 Notice that this is not contradictory to the statement above that single objects *cannot*
690 be accessed by *their* name. Each object has a unique identifier, the name, but objects
691 can have names for their constituents that are *not* unique, like `mu`. This simply describes
692 a part of the PDF and any `Gauss` will have a parameter `mu`.

693 4.3.2 Implementing a custom PDF

694 An essential feature of `zfit` is the ability to simply create custom models. There is a large
695 freedom in building models from the `BaseModel` class, since it takes care of most boilerplate
696 and has well defined entry points than can be customized. The full implementation details
697 and possibilities for customizations are described in Appendix C.6.

698 For the most common use cases though, there exists a simple way of creating a custom
699 model. The `ZPDF`, basically a more user friendly wrapper around `BasePDF`, can be used
700 as a base class in these simple cases. The following function will be used as the PDF
701 shape

$$702 f(x, y) = a \cdot x^2 + b \cdot y^4. \quad (7)$$

703 To implement this function, `_unnormalized_pdf` has to be overwritten. For the vast
704 majority of custom models, this is the only method to be overwritten. Changing other
705 methods, especially `_pdf`, is an advanced feature and only needed in special cases. For
706 more details on the customization and the possibility of hooking into the calls see Appendix
707 C.6.

708 This is a two dimensional PDF with two parameters. To implement it in `zfit`, a new
class is created

```

class EvenPolyPDF(zfit.pdf.ZPDF):
    """Implementation of  $f(x, y) = a*x^2 + b*y^4$ """
    _PARAMS = ['a', 'b']
    _N_OBS = 2 # since two dimensional

709 def _unnormalized_pdf(x):
    xdata, ydata = x.unstack_x()
    a = self.params['a']
    b = self.params['b']
    return a * xdata ** 2 + b * ydata ** 4

```

710 Note that we only need the *shape* of the function but do not need to take care of the
711 normalisation, as numerical methods are already implemented in the base class.

712 We can see the advantage of the preprocessing that is done by the base class, especially
713 the reordering of the data `x`. It is a `Data` object and calling the `unstack_x` method returns
714 a list¹⁴ containing a Tensor for each column sorted according to the models observables.
715 The length of the list has to coincide with the specified `_N_OBS`, the number of observables.
716 It is not mandatory to specify this field in a Model and is sometimes not possible to know
717 previously, in which case it can simply be left away.

718 The naming of the parametrization of the function is defined with `_PARAMS`. These
719 exact names have to be used when creating an instance of the model. The parameters are
720 then stored in the `params` dictionary and extracting them is usually the first step inside
721 `_unnormalized_pdf`.

722 Documentation plays an important role here: it defines the name of the parameters
723 that have to be used and what they represent in the function, but it is also crucial to
724 communicate the ordering of the data. In this case, a user can see that the first observable
725 and the corresponding column in the data will be used as x in the function.

726 This PDF is already complete and works out of the box. We can create an instance
727 and use its methods. As an example, the observables of the instance will be called "xobs"
728 and "yobs" with the normalisation range going from 0 to 10 and from 0 to 5, respectively.

```

param_a = zfit.Parameter("a", 1.)
param_b = zfit.Parameter("b", 2.)
x_obs = zfit.Space("xobs", limits=(0, 10))
y_obs = zfit.Space("yobs", limits=(0, 5))
obs = x_obs * y_obs
poly_model = EvenPolyPDF(a=param_a, b=param_b, obs=obs)

729 prob = poly_model.pdf(...) # some data needed

x_limits = zfit.Space("xobs", limits=(3, 5))
y_limits = zfit.Space("yobs", limits=(1, 3))
integral_limits = x_limits * y_limits
integral = poly_model.integrate(integral_limits)

sample = poly_model.sample(n=1000)

```

¹⁴Or a single Tensor for the one dimensional case if not specified differently in the arguments.

730

Using the `integrate` method, we obtain the integral i of our normalized PDF f_{norm}

$$i = \int_3^5 dx \int_1^3 dy f_{normed}(x, y) \quad (8)$$

$$\stackrel{(3)}{=} \frac{\int_3^5 dx \int_1^3 dy f(x, y)}{\int_0^{10} dx \int_0^5 dy f(x, y)} \quad (9)$$

731

with x and y corresponding to the observables \mathbf{x} and \mathbf{y} , respectively, and $f(x, y)$ is the unnormalised PDF as defined in Eq. 7. Using Eq. 3 we end up with a precise formulation of what is *actually* executed in `zfit`. In other words, the integral was calculated over the limits 3 to 5 and 1 to 3 respectively and normalised over the normalisation range `obs`. The latter is defined on initialisation and taken as the default normalisation range. As mentioned before, the limits could also be multiplied in a different order resulting in `limits` having the order (“yobs”, “xobs”). The model takes care internally (see also Appendix C.6.1) that the right limits are at the right place.

732

733

734

735

736

737

738

The `sample` method is used to generate 1000 points from the model. With the instance created, also probability densities `prob` for points can be calculated by calling the method with some `Data` object called `data` as follows

742

```
probs = poly_model.pdf(data)
```

743

744

745

It is worth noting that none of the operations are executed yet and what is returned by the methods are Tensors, as described in Sec. 3.1. To run the actual computations, it is necessary to call `zfit.run(...)` on any Tensor. Running

746

```
probs_np = zfit.run(probs_np)
integral_np = zfit.run(integral)
```

747

returns an actual numpy array for `probs_np` and a Python float for `integral_np`.

748

749

750

751

752

753

For further, advanced customization of the PDF, methods can be overridden as described in Appendix C.6.1. However, in most use cases there exist better ways: for example, to add an analytic integral to the model, overwriting `_analytic_integrate` should be the last resort. Instead, integrals defined over a specific range only or the whole range and over all dimensions or just partially can be registered with a model. A priority attribute allows to specify preferences on one over other methods.

754

755

756

A typical use case for this feature is a special integral that is known exactly, for example the value of the integral over the full space of a Gaussian shaped model¹⁵ An integral over both dimensions is registered but partial integrals could be added as well

$$\int_{x_0}^{x_1} \int_{y_0}^{y_1} f(x, y) dx dy = (1/3 \cdot a \cdot x^3 + 1/5 \cdot b \cdot y^5) \Big|_{x_0}^{x_1} \Big|_{y_0}^{y_1}. \quad (10)$$

757

758

759

760

In case a full integral is requested but no analytic integral over all dimensions is available, the fallback of `integrate` looks for partial integrals. If available, it uses them to numerically integrate over the remaining dimensions instead of using the unnormalised PDF.

761

The implementation of the integral with `zfit` is done by registering it to the PDF

¹⁵As already hinted in Eq. 8, it is important to note that the integral *to be implemented* is over the unnormalised PDF as implemented in the `_unnormalized_pdf` method.

```

def integral_full(x, limits, norm_range, params, model):
    lower, upper = limits.limit1d
    a = params['a']
    b = params['b']

    lower = ztf.convert_to_tensor(lower)
    upper = ztf.convert_to_tensor(upper)
    def indef_integral(limit):
        return 1 / 3 * a ** 3 + 1 / 5 * b * y ** 5

    return indef_integral(upper) - indef_integral(lower)

lower = ((zfit.Space.ANY_LOWER, zfit.Space.ANY_LOWER),)
upper = ((zfit.Space.ANY_UPPER, zfit.Space.ANY_UPPER),)
limits = zfit.Space.from_axes(axes=(0, 1), limits=(lower, upper))

EvenPolyPDF.register_analytic_integral(func=integral_full,
                                       limits=limits)

```

762

763 Since the observables that will be assigned to each axis are unknown, the `Space` has
764 to be defined with axes, not with observables. This follows the principle of using axes
765 inside the model as explained in 4.1. Instead of `ANY_LOWER` and `ANY_UPPER`, it could
766 also be defined over a specific domain by using plain Python floats. The PDF will now
767 automatically use this integral if possible. Otherwise, as for example when creating a
768 partial integral, it will silently fall back to numerical integration.

769 This example demonstrates how to implement a shape that does only depend on the
770 data and parameters. Models, that also depend on other models, such as the `SumPDF`
771 from Sec. 3, require to be built from a `Functor`. This implies as a minor change only
772 an additional argument to the base class, the depending models, in order to track the
773 dependencies correctly. This is described in more details in Appendix C.6.6.

774 4.3.3 Sampling

775 Sampling from a model can be done in two different ways to cover two distinct use
776 cases. On one hand, the method `sample` returns a pure Tensor, which behaves as any
777 other sampling algorithm in TF, for example `tf.random.uniform`. This can be useful
778 for specific and mostly advanced cases, but the overall behaviour can be unintuitive for
779 inexperienced users, as discussed in Appendix B.2.

780 On the other hand, sampling from a model is typically used to perform toy studies,
781 which consists of the generation of events according to the model and afterwards a fit of the
782 model with randomly initialized parameters to the sampled data. With `create_sampler`,
783 a `Data`-like object which handles the correct storage and the sampling from the model is
784 created. As any other `Data`, it can be used to build a loss. The actual data is sampled by
785 invoking `resample` and stays unchanged until it is invoked again. This sampler keeps a
786 dependency on the original model and uses it to sample, keeping the original parameter
787 values it was created with. If desired, any parameter that has changed in between will be
788 at that new value when `resample` is called, effectively changing the sampling shape.¹⁶

¹⁶Just to be clear: the default behaviour is that the sampler samples from the exact distribution that it was created with *in terms of parameter values*.

789 To illustrate the use of the toy sampler, let's look at a toy study example. Note that
 790 the actual (large) *objects* needed such as the model, minimiser, loss etc. are created
 791 *outside* of the loop and only the necessary *methods* are called *inside*. In general, this
 792 removes boilerplate code from the additional object creation but is especially important
 793 when working with a graph based backend. We assume to have previously built a model
 794 already and a minimizer as for example shown in Sec. 3, and we generate 1000 events
 795 according to our model.

```

sampler = model.create_sampler(n=1000)
nll = zfit.loss.UnbinnedNLL(model, sampler)
for _ in range(n_toys):
    sampler.resample() # now the sampling happens
    for param in model.get_dependents(only_floating=True):
796     param.set_value(np.random.normal()) # any initial value
    result = minimizer.minimize(nll)

    if result.converged: # check if the fit was successful
        ... # safe results in a list or similar

```

797 The sampling is implemented with the accept-reject method that works with arbitrary
 798 shapes. If an analytic inverse integral is registered, this will be used, providing a more
 799 efficient way of sampling. More details about the different techniques and the use of
 800 importance sampling are described in Appendix C.7.

801 4.3.4 Extended PDFs

802 In addition to the shape, a PDF can carry more information, namely a yield, in which
 803 case we refer to it as an Extended PDF. The yield is a scale that describes the magnitude
 804 of the PDF, typically reflecting the number of events in the data sample. It can be used
 805 to multiply the output of `pdf`, the normal probability density function, which results in a
 806 number probability. This implies that when multiplying the integral with the yield, the
 807 number of events is retrieved instead of the probability over a certain range.¹⁷ To count
 808 for example how many signal particles are in the sample, a composite PDF existing of
 809 a background and a signal component, both extended, can be fitted to the data. The
 810 `SumPDF` used to build this composition does not need a fraction if the PDFs are already
 811 extended, but uses the yields, normalised to one, as fractions. In this case, integrating the
 812 PDF that represents the signal shape over the whole range returns the number of signal
 813 events.

814 To create an extended PDF, `create_extended` can be used. This method returns a
 815 copy of the PDF and adds a yield to it.

816 4.4 Loss

817 As discussed in Sec. 2, the loss is the core of the problem specification, since it describes
 818 the disagreement between a model and the corresponding data. Additionally, it can
 819 contain constraints which help further specify the problem.

¹⁷ *Currently*, the integral of an extended PDF is multiplied by the yield *by default* but this behaviour is meant to change in a future version

820 Having the loss as an independent part of the whole workflow is a crucial design feature
 821 in `zfit`: it is the connection between the model, data, and their relation on one side and
 822 the minimisation process on the other side. Having the loss as an extra step accomplishes
 823 decoupling the former from the latter: a minimiser can take a loss and minimise it *without*
 824 knowing anything about the underlying models, data or the actual definition of the loss.
 825 Therefore, it is important that the loss knows everything that is needed for a minimisation,
 826 as listed in detail in Appendix C.8.

827 Basic loss implementations like the `UnbinnedNLL` use a PDF and data to calculate the
 828 loss according to Eq. 5. With extended PDFs, an additional term is taken into account
 829 if using `ExtendedUnbinnedNLL` instead, as derived in Eq. 26. Furthermore, additional
 830 terms can be added to express prior knowledge about parameters as in Eq. 29. To keep
 831 the flexibility, *any* Tensor can be added as a constraint to the loss with the method
 832 `add_constraint`. Alternatively, a custom constraint can be implemented by using the
 833 base class `BaseConstraint`. `zfit` implements the most often used constraints to improve
 834 usability: in the example from Sec. 3 a Gaussian constraint for the parameter μ could be
 835 applied by adding the following line after the creation of the `nll`.

```
mu = ... # parameter
nll = zfit.loss.UnbinnedNLL(gauss, data)
836 mu_constr = zfit.constraint.nll_gaussian(params=mu, mu=6.8, sigma=0.4)
nll.add_constraint(mu_constr)
```

837 Typically, some parameters are shared between different fits to different data samples.
 838 These can be obtained through a simultaneous fit of all the datasets by creating multiple
 839 PDFs with some of the `Parameter` objects being the same. Since this corresponds to
 840 a simple addition of the losses as seen in Eq. 25, `zfit` allows to perform precisely this
 841 operation. As an example we create two Gaussians and assume to already have their data.
 842 Here the μ is shared while the σ is not. Limits and the data are created as in the example
 843 of Sec. 3:

```
mu = zfit.Parameter("mu", 7)
sigma1 = zfit.Parameter("sigma", 1.1)
sigma2 = zfit.Parameter("sigma", 1.5)

844 gauss1 = zfit.pdf.Gauss(mu=mu, sigma=sigma1, obs=limits)
gauss2 = zfit.pdf.Gauss(mu=mu, sigma=sigma2, obs=limits)

nll1 = zfit.loss.UnbinnedNLL(gauss1, data1)
nll2 = zfit.loss.UnbinnedNLL(gauss2, data2)
simul_nll = nll1 + nll2
```

845 Alternatively, a list of models and their corresponding data can be given to create the
 846 loss

```
847 simul_nll = zfit.loss.UnbinnedNLL([gauss1, gauss2], [data1, data2])
```

848 There is a special loss available in `zfit` that gives a flexibility rarely found in other
 849 fitting packages: the `SimpleLoss`. This object lightly wraps any Tensor, which allows
 850 to build *any* kind of loss. No dependency on the data structure or the model layout is

851 required and it is completely up to the user. This allows to create not yet implemented
852 losses, such as binned ones, and allows other libraries which build loss functions with
853 TF to simply hook in with this mechanism. Because the `SimpleLoss` can then be used
854 with the next steps such as the minimisation and error estimation, an other library can
855 therefore use the whole available tools in `zfit` as well as any library that builds on top of
856 it.

857 4.5 Minimisation

858 The optimisation of functions is a large topic by itself and a lot of implementations of
859 different algorithms exist. They usually need a function that returns a value, such as
860 the loss, which depends on the parameter values that they use as arguments. Since the
861 computation of the loss is efficiently implemented in `zfit` thanks to TF and this usually
862 is the heavy part of the minimisation, in practice any minimiser can be wrapped easily.
863 In `zfit` several minimiser algorithms are implemented by wrapping existing libraries and
864 giving them a common API. An important part of the design is that the creation of a
865 minimiser object does neither execute any minimiser function nor tie itself to a specific loss,
866 it's simply the configuration of the minimiser. This implies that the minimiser is stateless,
867 and to actually invoke it, the `minimise` method needs to be called. The information about
868 the minimisation procedure and the parameter values are collected in a `FitResult` and
869 returned by the minimiser.

870 4.5.1 Different optimisations

871 While a whole variety of algorithms exists, not all are equally feasible to be used for
872 model fitting as done in `zfit`. There are various distinctions that influence the choice of
873 a certain minimiser.

874 **Derivative** Some optimisers use the derivatives and others don't. In general, using the
875 derivative provides a huge advantage since it tells about the local shape of the
876 function. This requires though that the function to minimise be continuous, which
877 is not always the case. For model fitting, functions are continuous and, thanks to
878 TF, `zfit` uses an analytic expression for them.

879 **Global/Local** Some optimisers are better in finding a global minimum by doing a
880 variation of a large grid search. Others focus on a local minimum by using a starting
881 point and going along a path to the next minimum. The latter is more accurate and
882 faster *if* a good initial parameter estimation is given, so that the minimiser does
883 not start far from the true minimum. In model fitting, a reasonable estimate can
884 often be made.

885 **Dimensionality** The number of parameters that have to be tweaked in order to minimise
886 the loss has an impact on the strategy. While the Hessian matrix can help greatly
887 with the minimisation, it has n_{params}^2 entries. This renders its usage impossible for
888 hundreds of thousands of parameters as used in deep learning. For model fitting
889 with a maximum of a few hundreds and typically around a dozen of parameters,
890 using the Hessian is feasible.

891 In HEP model fitting, mostly local, second order derivative minimisers are used. An
892 important algorithm is the Newton method of minimisation: leaving the mathematical
893 details away, the algorithm performs a second order Taylor approximation of the function
894 using the exact Hessian. Assuming that the target is a saddle point with the derivate
895 equal to zero, the equation is solved and the algorithm jumps to the estimated minimum.
896 Since this solution is only the true minimum if the second order approximation were exact,
897 this step is repeated until some convergence criteria are fulfilled.

898 Newton’s method is often not directly used since the computation of the Hessian and
899 its inverse can be computationally expensive. Instead an approximation of the inverse
900 Hessian is calculated and updated on every minimisation step, giving rise to a family
901 of methods called Quasi-Newton methods. Since these updates bring some ambiguity
902 in higher dimensions, different methods use different assumptions on the updates. A
903 prominent example is the BFGS algorithm and its low memory variant L-BFGS, which
904 only stores the most recent steps. L-BFGS is also the method implemented in Minuit, the
905 most common used minimiser in frameworks like ROOFIT or ROOT.

906 Another important point is the stopping criteria. While Minuit has its own stopping
907 criteria and in general good criteria have to be found with experience, this tuning for
908 other minimisers is still ongoing work in `zfit`. It is though simply possible to define and
909 change this criteria for a minimiser.

910 A special mention has to be made about TF optimisers. As discussed above, for
911 deep learning first order algorithms are used relying on the so called “gradient descent”
912 technique. This algorithm simply follows the negative gradient iteratively. Advanced
913 variations alter the step size based on heuristics but are still largely inferior for low
914 dimensional problems in comparison with Quasi-Newton methods. A wrapper for any TF
915 optimiser is available in `zfit`, an example is provided using the implementation of the
916 Adam [13] optimiser. Simple benchmarks though yield an order of magnitude increase
917 in the number of minimization steps compared to Quasi-Newton methods, they are not
918 competitive. to the

919 4.6 Results and uncertainties

920 After every minimisation, the information about the process as well as the results are
921 returned as a `FitResult`. This object does not only store information but is also capable
922 of performing uncertainty estimations of the parameters.

923 The most important information that it collects is:

- 924 • Information about the parameters, including the values at the function minimum,
925 information about their limits and uncertainties calculated using different methods.
- 926 • General information about the minimisation itself, including
 - 927 – A flag that indicates whether the minimisation was successful or not.
 - 928 – The minimum of the loss function that was determined by the minimiser.
 - 929 – An estimation made by the minimiser of how far away the minimum value is
930 from the actual true minimum.
 - 931 – All the additional information produced by a minimiser. This can highly vary
932 depending on which minimiser was used.

- 933 • The instance of the minimiser that was used to perform minimisation. Since
934 minimisers are stateless, no information is stored in it.¹⁸
- 935 • The instance of the loss that was minimised. Since a loss keeps references to the
936 model and data it was built with, it is possible to thereby retrieve all information
937 regarding this minimisation.

938 As the last of five steps in the minimisation workflow, the `FitResult` serves again
939 as an additional abstraction layer. A lot of different statistical quantities such as limits,
940 confidence intervals and more can be calculated using the result, loss and minimiser. Since
941 they are all bundled together in the `FitResult`, no other object is required for advanced
942 statistical treatment of fit results.

943 **4.6.1 Parameter uncertainties**

944 The values of the parameters at the minimum are important, but they are meaning-
945 less without an uncertainty estimate. Therefore, the `FitResult` provides two ways of
946 calculating it:

- 947 • For a fast, approximative and symmetric estimate, the `hesse` method can be invoked.
948 It provides an estimation based on the Hessian matrix, assuming a second order
949 approximation of the loss around the minimum value of the parameter.
- 950 • If there are high non-linearities in the loss and the parameter correlations, the
951 actual uncertainties differ strongly from what `hesse` returns. Good estimates can
952 be retrieved by creating a profile: fixing the parameter at a certain value and run
953 a complete minimisation. The calculation is invoked by using the `error` method.
954 If the Minuit minimiser was used to perform the minimisation, the `minos` method
955 can be invoked in this way. Currently, no other error estimation is implemented,
956 but any custom statistical method can be easily applied thanks to the information
957 contained in the fit result.

958 To only calculate the uncertainties with respect to specific parameters, the desired
959 parameters can be given as arguments to `hesse` and `error`. The results for each parameter
960 are cached and won't be recomputed on an additional call.

¹⁸Ideally. Some minimisers like `iminuit` have a state and can be accessed like this. A copy of the actual minimiser is stored in the `FitResult`

5 Performance

In the previous Sections, the structure and logic of model fitting in `zfit` was elaborated. Apart from the functionality, model fitting involves a lot of number crunching and thus state-of-the-art performance is a hard requirement.

The performance of code depends on several factors and can often be divided into a serial and a parallel part. The efficient implementation of parallelised parts of code are as important as deciding *when* to actually run in parallel, as discussed in Appendix D.2. While this part is mostly left to TF, the fact that some pieces of code simply *are not* parallelisable, and therefore cause a bottleneck, depends on the nature of the problem itself. In real code, the two are not clearly separated and a mixture of both influence the actual performance.

In this Section we will focus on quantifying the performance of `zfit` by measuring the execution time of the whole process, mixing together the performance of TF and `zfit` with the bottlenecks coming inherently from model fitting. While this limits the ability to draw conclusions about bottlenecks and remove them, it reflects the performance of the library as experienced in real life usage.

In the following, two studies are presented: on one hand, a more artificial but well scalable example consisting of Gaussian models, and on the other hand, a more realistic case, namely a fit to an angular distribution is performed. In both cases the performance of toy studies is measured and compared. The hardware specifications can be found in Appendix D.1.

5.1 Gaussian models

There are three quantities of interest to describe the scaling of the library: the complexity of the model, the number of free parameters and the size of the data sample used to fit. In this study, a sum of Gaussian PDFs is used as the model. The fractions are constant and the mean and width are parameters, each shifted by a different constant, and either all depending on only two parameters or scaling with the number of Gaussians.

This model is used to perform toy studies. First a sample from the model is created using a fixed, initial parameter setup. Then the parameters are randomized and a fit to the sampled data is performed. The implementation of sampling and the minimisation are two distinct steps and for fits to data, only the latter is actually invoked. Since this combines two execution time measurements making reasoning even harder and there is a known, temporary¹⁹ bottleneck in `zfit` sampling, only the execution time of the minimisation is measured. A approximate comparison of the current performance with sampling can be found in Appendix D.3.

For each setup, twenty toys are run ranging from 2 up to 20 Gaussians with sample sizes from 128 to 8 million events (except for GPU, where it goes only up to 4 million²⁰). A comparison with an implementation in ROOFIT is performed, although only ranging from 2 to 9 Gaussians.²¹ To have a fair comparison, both use the Minuit algorithm for

¹⁹This problem is resolved now, though in this thesis the old measurements are still used since the conclusions are the same.

²⁰The GPU used has a comparably small memory. Performing larger-than-memory computations and multi-GPU is still work in progress.

²¹Due to technical problems using a sum of more than 9 pdfs with ROOFIT that were overcome only

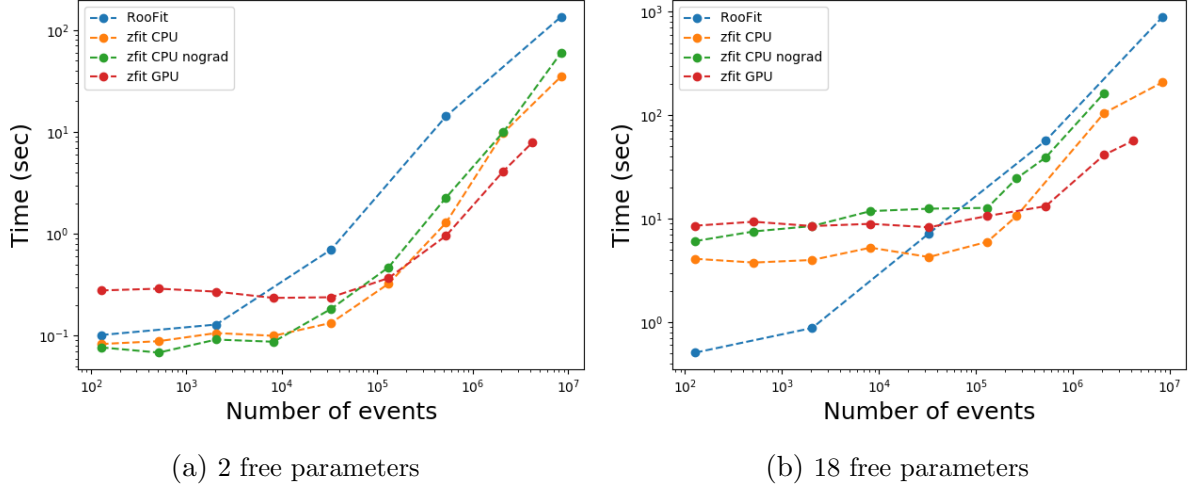


Figure 5: A sum of 9 Gaussian PDFs with shared (left) or individual (right) μ and σ . On the y axes, the time for a single fit is shown, averaged over 20 fits. It is plotted against the number of events that have been drawn per toy.

1000 minimisation.

1001 In the following, 4 cases are considered:

- 1002 • an implementation in `zfit`, labelled as “zfit CPU”, using the analytic gradient
1003 provided by TF. In this case, an initial run is done to remove the graph compile
1004 time. While not significant, this provides a more realistic estimation.
- 1005 • The same implementation as above is used but the analytic gradients provided by
1006 TF are disabled, denoted by the addition “nograd”. Instead, the Minuit minimiser
1007 calculates a numerical approximation of the gradients internally.
- 1008 • The same as “zfit CPU” but run on a GPU and labelled as “zfit GPU”.
- 1009 • An implementation in `ROOFIT` using the Python bindings with `PyROOT`. The
1010 parallelisation is done when invoking the `fitTo` method and equals the number of
1011 cores available.

1012 In Fig. 5, the time per toy is measured and plotted against the number of events
1013 used in the toy. While `zfit` on CPU outperforms `ROOFIT` in the left plot with only two
1014 free parameters, in the right plot with 18 free parameters, `ROOFIT` performs significantly
1015 faster with a low number of events. This comes from the fact that while the Minuit
1016 minimiser is used in both measurements, their are tweaked differently: the version used by
1017 `ROOFIT` is configured to better cope with a bumpy likelihood as given with a high number
1018 of parameters and a low number of events, which results in a vastly superior performance
1019 for complicated, low statistics cases while reducing the performance in simpler cases. It is
1020 to note that tweaking the minimiser is still ongoing effort in `zfit` and the performance
1021 behaviour is likely to change in the future. Different tuning cannot simply be quantified
1022 and compared, since not only the number of evaluations of the loss function, but also the

after the measurements were done.

1023 gradient and the (very expensive) Hessian computation is part of the strategy. This limits
1024 the conclusions that can be drawn from the comparison of the performance with a low
1025 number of events. Furthermore, the minimiser for 18 free parameters using `zfit` was not
1026 able to converge for more than a million of events *without the automatic gradient from*
1027 *TF*, therefore no data points are available there.²²

1028 The comparison in Fig. 5 still reveals a few important things that agree very well with
1029 the expectations.

- 1030 • As the number of events increase, the execution time of ROOFIT monotonically
1031 increases.
- 1032 • There is no advantage using parallelisation for very few calculations since the
1033 overhead of splitting and collecting the results is dominant. With increasing number
1034 of events this gets negligible and as an effect the execution time of `zfit` increases
1035 way slower than for ROOFIT. This also comes from the fact that more events mean
1036 a more stable loss shape, so the minimiser used in `zfit` performs better.
- 1037 • The GPU is highly efficient in computing thousands of events in parallel. For only
1038 a few data points though, the overhead of moving data back and forth dominates
1039 strongly, making it unfeasible for only a small number of events. A continuous
1040 decrease of the time up to 10'000 in the left plot of Fig. 5 events confirms that and
1041 even shows a drop in the computation time.

1042 As a conclusion, the speed for very small examples around 100 events of `zfit` is
1043 *marginally* slower than the corresponding ROOFIT example. For larger fits, the speedup
1044 of `zfit` is up to a factor of ten at around a million events. For more complicated fits and
1045 a small number of events, ROOFIT is an order of magnitude faster because of its currently
1046 better tuned minimiser, though there is ongoing work in `zfit`. The GPU delivers in
1047 larger examples a similar performance as the multicore setup. Finally, not using the TF
1048 gradient yields only a minor penalty for large fits and can even be faster for small ones,
1049 experiments have shown that the effect of an increased time can be larger for complicated,
1050 real use cases and helps reducing the number of steps required to take by the minimiser.

1051 In Fig.6, a comparison of what can be achieved within a certain time frame is shown.
1052 The fits scale with the number of events for different number of added Gaussians, having
1053 only two free parameters *in total*. The observation matches the intuitive behaviour of
1054 scaling with complexity and size of the sample. Compared to ROOFIT, the fitting time of
1055 `zfit` increases an order of magnitude less; note that both time scales end at 100 seconds.
1056 In this time, `zfit` fits 8 millions of events with 16 Gaussians, ROOFIT does half a million
1057 with 9 Gaussians. For this setup, 8 CPUs on a shared cluster were used, leaving slight
1058 ambiguity about the results due to the unknown configuration and actual workload on
1059 the machine. However, the order of magnitude matches with the results in Fig. 5, which
1060 were performed in a clean environment.

1061 Scaling the number of free parameters with the number of Gaussians added is shown
1062 in Fig.7. We see clearly that for a large number of parameters having a higher number
1063 of events can actually be more performant on an absolute scale, at least up to a certain
1064 threshold, since this reduces the number of steps to be taken for the minimisation.

²²This indicates numerical instabilities due to limited precision and there are ways to circumvent them which are planned to be implemented in `zfit`.

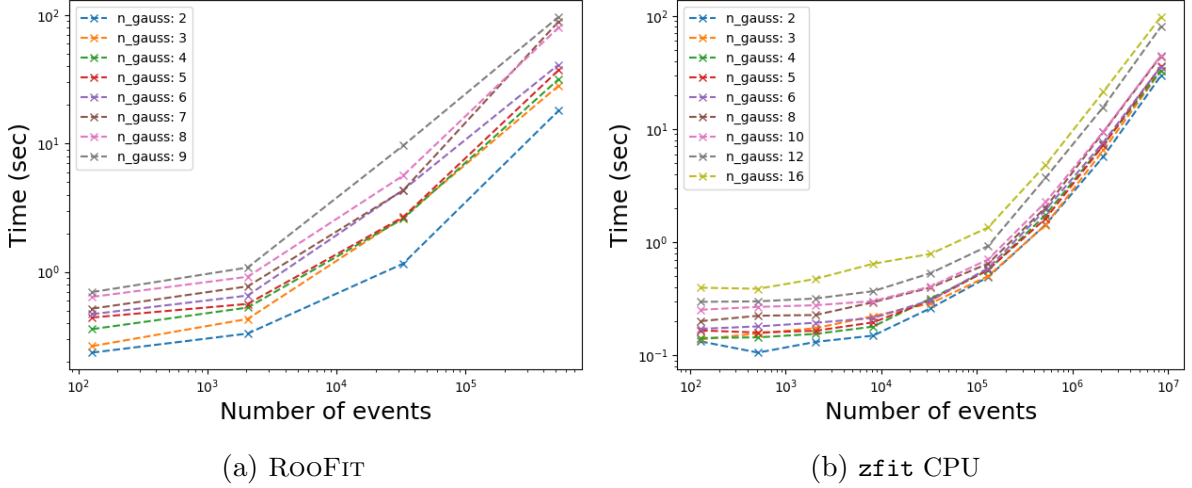


Figure 6: Measurement of the computation time with a sum of n_gauss Gaussians and in total two free parameters. Notice that the y-scale is the same for both plots but the x-axis for **zfit** goes an order of magnitude higher. Also, **zfit** sums up to 16 Gaussians whereas **ROOFIT** only goes to 9.

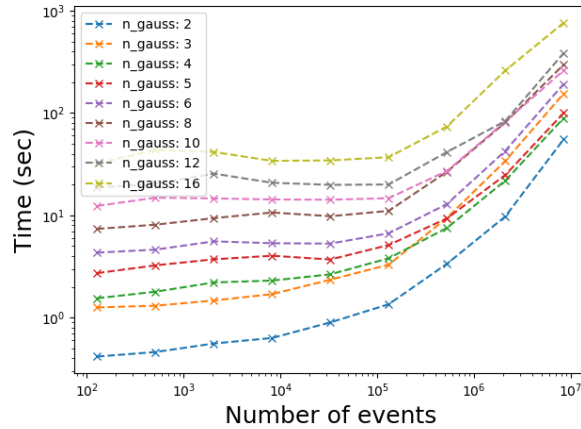


Figure 7: Time of a single toy in dependence of the number of events used. Plotted for a sum of n_gauss Gaussians and two free parameters *per Gaussian*.

1065 5.2 Angular analysis

1066 In order to study the scalability of **zfit** in a challenging, realistic setting, a toy study from
 1067 an ongoing analysis involving **zfit** is used. The example, which consists in the angular
 1068 analysis of $B^0 \rightarrow K^*(\rightarrow K^+\pi^-)\ell^+\ell^-$ with ℓ being either e or μ , is an important legacy
 1069 analysis and the fact that it was implementable in **zfit** is itself already an achievement.
 1070 The model is from the folded angular analysis. The folding to measure the P5' parameter

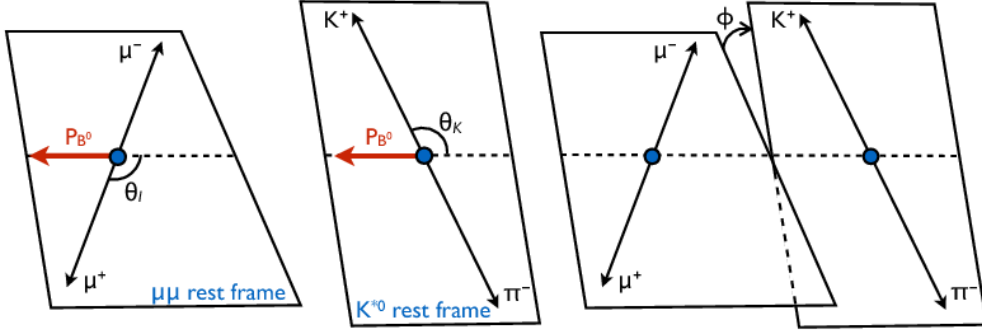


Figure 8: Schematic view of the angles of the $B^0 \rightarrow K^*(\rightarrow K^+\pi^-)\ell^+\ell^-$ decay. The image was taken from [15]

1071 as defined in [14] is implemented. The resulting model is given by

$$\begin{aligned}
 f_{angular}(\theta_K, \theta_\ell, \phi; F_L, A_T^{(2)}, P_5') &= \frac{3}{4}(1 - F_L) \sin^2 \theta_K + F_L \cos^2 \theta_K \\
 &+ \frac{1}{4}(1 - F_L) \sin^2 \theta_K \cos 2\theta_\ell \\
 &- F_L \cos^2 \theta_K \cos 2\theta_\ell \\
 &+ S_3 \sin^2 \theta_K \sin^2 \theta_\ell \cos 2\phi \\
 &+ S_5 \sin 2\theta_K \sin \theta_\ell \cos \phi
 \end{aligned} \tag{11}$$

1072 with

$$\begin{aligned}
 S_3 &= \frac{1}{2} A_T^{(2)} (1 - F_L) (1 - F_L) \\
 S_5 &= P_5' \sqrt{F_L (1 - F_L)}.
 \end{aligned}$$

1073 The model is depends on three angles, shown in Fig. 8, where

- ϕ is the angle between the plane spanned by the flight direction of the two leptons with the plane spanned by the kaon and pion,
- θ_K is the angle between the kaon and the negative flight direction of the B^0 and
- θ_ℓ is the angle between the ℓ^+ (ℓ^-) and the negative flight direction of the B^0 .

1074
1075
1076 The model is extended to four dimensions by adding a description of the B invariant
1077 mass distribution by building the product with the angular part as defined in Eq. 11. The
1078 model used for the mass shape is a combination of two Gaussians, each with a powerlaw
1079 tail.

1080 The implementation of the angular part is straight forward and follows the example in
1081 Sec. 3. Using the implemented `DoubleCB` for the mass, the four dimensional model can
1082 be built as

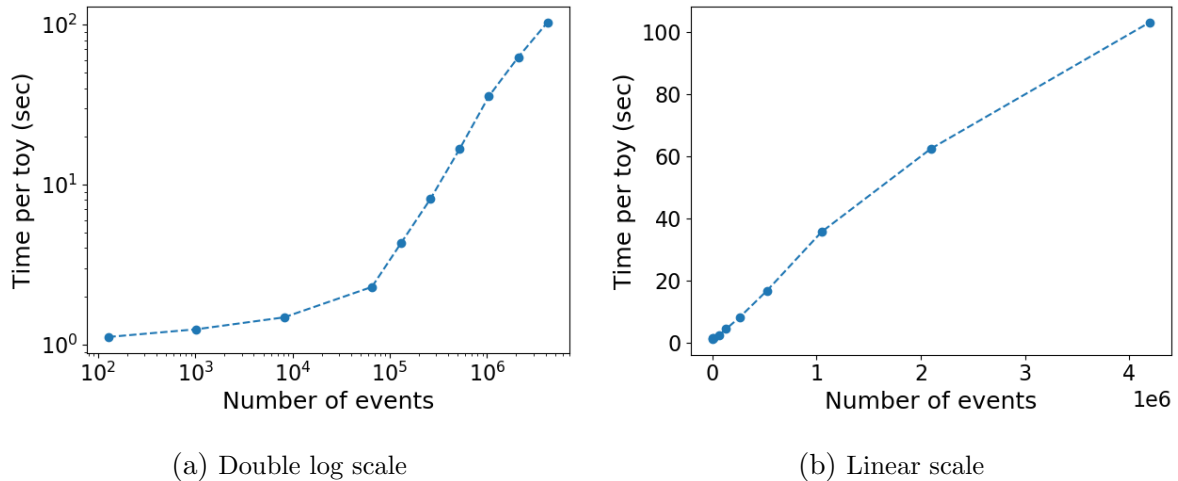


Figure 9: The figures show the time needed per toy for the four dimensional P5' folded angular distribution. 25 toys are produced for each number of events.

```

limit_thetak = zfit.Space('thetak', limits=...)
limit_thetal = zfit.Space('thetal', limits=...)
limit_phi = zfit.Space('phi', limits=...)
limit_bmass = zfit.Space('bmass', limits=...)

1083 angular_obs = limit_thetak * limit_thetal * limit_phi
angular = AngularPDF(obs=angular_obs, ...)
mass = zfit.pdf.DoubleCB(obs=limit_bmass, ...)

model = angular * mass

```

1084 With this model a set of toy studies is performed, varying in each of them the number
1085 of events while keeping the number of free parameters fixed to nine.

1086 Fig. 9 shows the performance of the toys on a shared cluster with 8 CPUs requested
1087 and we can see that the time per toy increases slightly sublinearly with the number of
1088 events. An interesting number is the average time for toys with around 1000 events,
1089 the expected number of events to be seen in the data and therefore actually used in the
1090 real toy studies, which is around one second per toy. Both go well together with the
1091 expectations of a model fitting library. The whole example demonstrates the suitability of
1092 `zfit` to be used in non-trivial, real world analyses.

1093 6 Beyond standard fitting

1094 In the previous Sections, we have seen the capabilities of `zfit` for fitting a model to data.
1095 The inherent flexibility and the powerful functionality of the core parts allow for the
1096 library to be extended beyond the usual feature set of general model fitting libraries in
1097 HEP. To perform more specialized fits, the components of `zfit` can be used as building
1098 blocks of a higher level fitter. In order to keep the size and features of `zfit` at a reasonable
1099 level, the whole project is split into multiple packages that are tightly coupled to each

1100 other. In this Section, we will have a closer look at the different packages that extend
1101 `zfit` beyond the simple model fitting we have seen so far.

1102 The main package is the core library `zfit`, as described in the thesis up to now, that
1103 provides all the fundamental building blocks and is in itself a self-contained fitting library.
1104 Content-wise, `zfit` offers a more field agnostic selection and does not contain models and
1105 tools too specific for HEP. The focus is on a stable and solid implementation together
1106 with the API and format definitions.

1107 More field specific content can be added in other libraries such as `zfit-physics`. This
1108 repository is meant to be the place for community contributions focusing on HEP-specific
1109 content. Guidelines, examples and automated tests are planned to be in place in order to
1110 lower the threshold for contributions. Furthermore, it will also contain functions dealing
1111 with physical quantities such as kinematics, which can serve as building blocks for models.
1112 Higher level interfaces that use `zfit` to build specific models are also intended to be
1113 placed in this repositories, for example a whole amplitude analysis framework currently
1114 under development which will be explained in the following Section.

1115 Additionally, the project also contains libraries that `zfit` depends on but which
1116 themselves are self-contained and which are factored out of the core- and extension
1117 libraries to standalone packages. This allows other projects than `zfit` to make use of them
1118 as well. An example is the implementation of `phasespace`, which generates four momenta
1119 of particles from a decay taking into account the correct kinematics. The package will be
1120 explained in more detail in Sec. 6.2.

1121 6.1 Amplitude fits

1122 `zfit-amplitude` is a higher level fitting library, which is currently under active develop-
1123 ment and in its early stage, built with elements from `zfit` in order to perform amplitude
1124 analyses, including Dalitz. In the following, the extension within the scope of the `zfit`
1125 project is discussed by showing an implementation example of the decay $D^0 \rightarrow K^+ \pi^- \pi^0$.
1126 As part of this, an additional library will be introduced, `phasespace`, which covers an
1127 essential part of being able to generate amplitude fits.

1128 In order to perform amplitude fits, `zfit-amplitude` contains

- 1129 • shapes such as the Breit-Wigner distribution that can be used to model resonances;
- 1130 • helper classes to build PDFs such as some that efficiently cache intermediate results
1131 specific to amplitude fits, and
- 1132 • a higher level interface to build amplitude fits in a transparent manner, similar
1133 to other specialised amplitude fitters in HEP. This removes the need of low level
1134 handling *but* keeps the flexibility to replace any part of it by a custom object.

1135 In order to perform toy studies as described in 5, an efficient way of sampling from a
1136 model is needed. While accept-reject is an universal and good working way, the efficiency
1137 can be very low in higher dimensional spaces and/or with peaky model shapes. To avoid
1138 inefficiencies, samples can be produced by using importance sampling as described in
1139 Appendix C.7, which requires a distribution approximately resembling the shape of the
1140 model. In `zfit-amplitude` the approach is to use the kinematic constituents of the decay
1141 particles. They are then transformed to the desired variables that are used in the model.

1142 This is a very general approach that allows all sorts of variables to be constructed from the
 1143 kinematics. On the downside, in order to produce the particles with realistic kinematics,
 1144 a phasespace generator is needed. While there is an implementation in ROOT called the
 1145 TGenPhasespace class, no equivalent is available in pure Python. Therefore, a library
 1146 porting the above with an extension to real world experiment kinematics and implemented
 1147 using TF has been created.

1148 More details on the `zfit-amplitude` library, especially on the higher level interface,
 1149 will be provided in a future paper and goes beyond the scope of this thesis.

1150 6.2 phasespace

1151 The kinematics of a particle are a four-dimensional tuple representing its four-momentum.
 1152 In a decay of a parent particle to lighter children particles, their kinematics are constrained
 1153 by the fundamental physical laws of momentum and energy conservation. With

$$p^{0,1,2} = p^{x,y,z} \quad p^3 = \sqrt{p^2 + m^2 c^2}$$

1154 this reads as

$$p_{parent}^\mu = \sum_i p_i^\mu \quad (12)$$

1155 with p_i being the four-momenta of all children particles. For a two-body decay $A \rightarrow BC$,
 1156 there are six free variables from the momenta of the two children. Using Eq. 12 there are
 1157 two degrees of freedom left in the decay kinematics that lead to a distribution. Furthermore,
 1158 the mass is only fixed for stable particles but is a distribution for resonances. To sample
 1159 from the phasespace within the `zfit` project a package named `phasespace` was created.
 1160 The purpose of it is to generate arbitrary decays obeying the physical constraints discussed
 1161 before and expressed in Eq. 12. Since it uses TF as the backend, it integrates seamlessly
 1162 into `zfit`.

1163 The algorithm used in `phasespace` is the Raubold and Lynch method for n-body
 1164 events as described in [16]. In principle, the algorithm builds a tree where every node has
 1165 two leaves representing a two-body decay. It calculates the available kinematic energy
 1166 that will be assigned to the particles from the difference of the parent rest mass and the
 1167 sum of all the children rest masses. The latter are either fixed or drawn from the mass
 1168 distribution in case of short-lived particles

$$E_{kin} = m_{parent} - \sum_i f_i^{mass}(m_i^{min}, m_i^{max})$$

1169 with f_i^{mass} being the mass distribution of the particle i , m_i^{min} the minimum mass recursively
 1170 determined from the children masses of particle i and m_i^{max} the available energy from the
 1171 top particle minus the other parent particles of particle i . The remaining E_{kin} is randomly
 1172 split into fractions along all the decaying particles in the tree where each fraction is the
 1173 kinetic energy for the boost of this particle.

1174 Given the mass of each particle in the tree and its kinetic energy, particles are recursively
 1175 generated starting from the top of the tree, *i.e.*, with the lightest particles. Each parent
 1176 particle randomly generates the two children in the available phasespace. Then the whole
 1177 decay tree is boosted to the parent momentum. This continues until the top particle
 1178 is reached, which is not boosted by default. However, an additional argument to the

1179 generation method allows to also boost it. This can be used to reproduce the physics
1180 happening in actual colliders.

1181 Usage

1182 The library has a main class `GenParticle` that describes the particles mass, either fixed
1183 or as a distribution, and has a method to set the children particle it decays to. This allows
1184 to build an arbitrary decay chain in an object-oriented way. As an example, the decay of
1185 $D^0 \rightarrow K^*(892)(\rightarrow K^+\pi^-)\pi^0$, which will be implemented as an amplitude fit in the next
1186 Section, would be implemented as following

```
1187 import phasespace as phsp

kplus = phsp.GenParticle("K+", mass=KPLUS_MASS)
piplus = phsp.GenParticle("Pi+", mass=PIPLUS_MASS)
piminus = phsp.GenParticle("Pi-", mass=PIMINUS_MASS)
kstar = phsp.GenParticle("K*", mass=kstar_mass_func)
dzero = phsp.GenParticle("D0", mass=DZERO_MASS)

kstar.set_children(kplus, piminus)
dzero.set_children(kstar, piplus)
```

1188 Here, `kstar_mass_func` is a function sampling from the mass distribution of a K^* .
1189 This is a custom function that can be implemented by the user.

1190 Having specified the decay chain, two methods can be used to generate the actual
1191 particles: Either `generate_tensor` which returns the four momenta as a Tensor and can
1192 be used directly inside `zfit`, or alternatively, `generate` returns the same information but
1193 as a numpy ndarray. Internally, a TF session is called and the computation is run.

```
1194 weights, particles = dzero.generate(1000) # generate 1000 particles
```

1195 The weights correspond to every single event and quantify the probability of the
1196 generated event. The returned `particles` is a dictionary containing the momentum of
1197 each particle. For example

```
1198 kstar_kinematics = particles['K*']
kstar_x = kstar_kinematics[:, 0]
kstar_mass = kstar_kinematics[:, 3]
```

1199 and the format of the kinematic is (number of events, components).

1200 As a shortcut for simple decays, a high level function `generate_decay` is available
1201 which allows to describe a decay with stacked lists of masses.

1202 6.3 Dalitz implementation

1203 As an example of an amplitude fit, a Dalitz analysis of the decay $D^0 \rightarrow K^+\pi^-\pi^0$ is
1204 implemented using some parts of `zfit-amplitude` together with `zfit` and `phasespace`.
1205 This implementation of the resonances and their shapes is done following the analysis in
1206 Ref. [17]. The amplitude is built using the isobar approach, which calculates the coherent

1207 sum of the individual contributions, the intermediate resonances. The implemented
 1208 resonances are $\rho(770)$, $\rho(1700)$, $K^{*0}(892)$, $K^{*+}(892)$, $K^{*0}(1430)$, $K^{*+}(1430)$ and $K_2^*(1430)$

1209 For amplitude analyses, dedicated fitting libraries exist and using a general purpose
 1210 fitter like `zfit` for this case is rather special. In the following, we will go through the
 1211 elements that are needed from a top-down approach and see how the problem can be
 1212 separated into smaller pieces using the functionality of `zfit`. Although unimportant
 1213 technicalities are left away and the example is a sketch only of the actual implementation,
 1214 it is going to be rather advanced and involves functionality only explained in Appendices.

1215 The observables in this case are the invariant masses of pairs of the different children
 1216 particles, which are $m_{K^+\pi^-}$, $m_{K^+\pi^0}$ and $m_{\pi^-\pi^0}$. The whole PDF f_{tot} is of the form

$$f_{tot}(m_{K^+\pi^-}, m_{K^+\pi^0}, m_{\pi^-\pi^0}) = \left\| \sum_i^{n_{amp}} c_i A_i(m_{K^+\pi^-}, m_{K^+\pi^0}, m_{\pi^-\pi^0}) \right\|^2 \quad (13)$$

1217 with c_i being the complex coefficient of each amplitude A_i . Expanding this term contains
 1218 cross ($i \neq j$) and square ($i = j$) terms b_{ij} of the form

$$b_{ij} = c_i c_j^* A_i A_j^*$$

1219 To implement this in `zfit`, a custom class is created that takes the coefficients,
 1220 amplitudes and a `Func` called `AmplitudeProduct` which calculates the above products b_{ij} .
 1221 Before we will look at these three components, we can build the global PDF as defined in
 1222 Eq. 13

```

class SumAmplitudeSquaredPDF(zfit.pdf.BaseFuncor):
    def __init__(self, obs, coefs, amps, ...)
        self._cross_terms = [AmplitudeProduct(c1, c2, a1, a2)
                               for (c1, a1), (c2, a2)
                               in combinations(coefs, amps)]

        self._square_terms = [AmplitudeProduct(coef, coef, amp, amp)
                               for coef, amp in zip(coefs, amps)]
    ...

    def _unnormalized_pdf(self, x):
        value = tf.reduce_sum([2. * amp.func(x=x) for amp in self.
                               _cross_terms]
                               + [amp.func(x=x) for amp in self._squared_terms],
                               axis=0) # over which axis to sum
        return tf.real(value) # convert it to a real number

```

1224 As the next piece, we need to have the amplitude product. Since this is complex in
 1225 general, it is again split into smaller parts that represent the real and imaginary parts by
 1226 a class `AmplitudeProductProjection`.

```

class AmplitudeProduct(zfit.func.BaseFuncorFunc):
    def __init__(self, coef1, coef2, amp1, amp2,...):
        prod_real = AmplitudeProductProjection(amp1, amp2, proj=tf.math.
            real)
        prod_imag = AmplitudeProductProjection(amp1, amp2, proj=tf.math.imag
            )
        super().__init__(funcs=[prod_real, prod_imag],
            params={'coef1': coef1, 'coef2', coef2}, ...)
    ...

    def _func(self, x)
        coef1 = self.params['coef1']
        coef2 = self.params['coef2']
        prod_real, prod_imag = self.funcs
        coeffs = coef1 * tf.conj(coef2)
        return coeffs * tf.complex(prod_real.func(x), prod_imag.func(x))

```

1227

1228 Splitting the real and imaginary parts has the advantage of keeping both of them real
1229 and therefore also their integrals. This allows to make full usage of the `zfit` integration.
1230 Since these parts are *independent* of any coefficient, this allows to cache the integral value.

```

class AmplitudeProductProjection(zfit.func.BaseFuncorFunc):
    def __init__(self, amp1: ZfitFunc, amp2: ZfitFunc, proj, ...):
        self.projector = proj
        super().__init__(funcs=[amp1, amp2], ...)
        self._cache_integral = None

    def _func(self, x):
        amp1, amp2 = self.funcs
        return self.projector(amp1.func(x) * tf.conj(amp2.func(x)))

    def _single_hook_integrate(self, limits, norm_range, name):
        integral = self._cache_integral
        if integral is None:
            integral = super()._single_hook_integrate(limits=limits,
                norm_range=norm_range,
                name=name)
            integral = zfit.run(integral) # simplified
            self._cache_integral = integral
        return integral

```

1231

1232 where we used the hook for the integration as described in Appendix C.6.

1233 All that is left now are the coefficients, which are just parameters, and the resonances.
1234 Using the Breit-Wigner function from `zfit-physics`,²³ we can build them

```

resonances = [
    'rho(770)': RelativisticBreitWigner(rho_770_plus_mass,
        obs=zfit.Space('m2pipi', limits...))
    ...
]
coeffs = [
    zfit.ComplexParameter.from_polar('c_rho770', 1.0, 0.0),
    ...
]

```

1235

²³The Breit-Wigner is currently an open merge request.

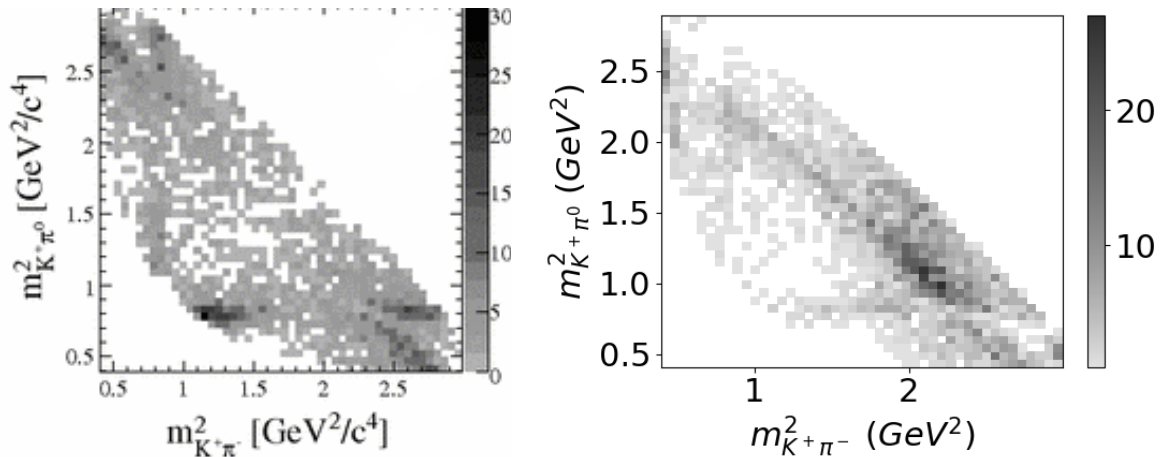


Figure 10: Dalitz plot of the invariant mass $K^+\pi^-$ and $K^+\pi^0$. The left histogram is taken from the paper while the one on the right was sampled from a distribution built with `zfit` and using `phasespace` as sampler.

1236 and combining all of the above we can finally create the whole model

```
1237 obs = zfit.Space(['m2kpi0', 'm2kpi-', 'm2pipi', limits=...])
pdf = SumAmplitudeSquaredPDF(obs=obs, coeffs=coeffs, amps=resonances,
...)
```

1238 This PDF is quite an accomplishment: while the implementation compared to all
1239 previous examples is more extensive, none of the above parts is *really* complicated to
1240 implement and uses basic `zfit` functionalities. Nonetheless, with this few lines of code,
1241 `zfit` extends its functionality into the world of amplitude fitters, where usually only
1242 dedicated tools exist.

1243 A missing part is the sampling: since the kinematics of the particles is not encoded into
1244 the PDF, it is needed inside the sampling. For that, we can create a decay as shown in
1245 examples above in Sec. 6.2 with the `phasespace` package. For simplicity, the Breit-Wigner
1246 distribution is used to model the mass shape. Building the decay and sample from it can
1247 be registered with the PDF as importance sampling.²⁴

1248 In order to compare the `zfit` implementation with the measured data, a sample is
1249 drawn from the PDF. It is to note that the right plot in Fig. ??, drawn with `zfit`, is
1250 currently not representative, since an instability in the sampling can vary the samples
1251 strongly depending on a single proposed event.²⁵ While the two distributions do not agree,
1252 it is notable that resonances are taken into account and that the `phasespace` is sampled
1253 with correct borders. This unveils another current limitation in `zfit`, that currently only
1254 rectangular limits are possible. As a future extension, arbitrary limits for `Space` will be
1255 implemented.

1256 What was shown in this example are the low level components to build an amplitude
1257 analysis using `zfit`. Most of the work above is straight forward but cumbersome and

²⁴Since this feature is not yet fully public in `zfit`, no explicit example is shown. The mechanism will be similar to the registration of an integral.

²⁵As described in C.7, weights going to zero as occurring in the `phasespace` can cause large problems.

1258 can be hidden behind a higher level interface that takes care of the right assignment of
1259 observables, resonances and compositions. Furthermore, and adding an additional layer
1260 around the resonances, the **phasespace** decays can be kept with the actual resonances
1261 which allows the phasespace to be an integral part of the amplitude. A higher level
1262 interface, which leaves the user with the specification of the resonances and coefficients
1263 but still offers the possibility to replace any part that was shown above, is currently under
1264 work in `zfit-amplitude`.

7 Conclusion and outlook

`zfit` is a versatile library that fills the gap of model fitting in Python for HEP. Built on top of the deep learning framework TensorFlow, it has shown great advantages including a remarkable speedup for parallelisation. The formalisation into five loosely coupled parts and an extensive base class for custom models extend its scope far beyond the usual feature set of HEP fitting libraries.

The project has been successful so far, and is being used in several high-impact analyses. In addition, it has been shown how the `zfit` design allows to build extremely complex analyses, namely amplitude analyses, in a reasonably simple way, unlike most general fitting libraries. However, a lot of work remains to be done on the way to establish a stable, reliable fitting library. The main features to be improved on in the near future are:

Binned fits Fits are sometimes performed in bins of data, mostly to speed up the computation. Furthermore, the shape of a model can be too complicated to be described analytically and has to be deduced from simulation or data samples by creating a template PDF. Currently, there is no native support yet for binned fits or template PDFs. This is right now under active development and will be added to `zfit` in the future.

Optimisation Model fitting can be a numbers game: in order to estimate uncertainties of parameters or to study the sensitivity of a fit with toys, a large number of repeated fits have to be performed. To keep this feasible in terms of time and computing resources, performance matters. There are currently still various places where the computation can be optimised. This includes the caching of computations and more efficient numerical integration by using advanced Monte Carlo techniques or other numerical methods.

Serialisation Models are currently built within a script. Often, a model needs to be stored and used again later on or in a modified version, which is not well achieved by just dumping the code. To actually define a model, for most cases no code is actually needed but a configuration file with the model description is sufficient. This allows to change certain parts of a model and rather inexperienced users to safely build a model. Therefore, a complete serialisation of a model into a human readable format is planned for `zfit`.

Content In HEP, there are quite a few different shapes and possibilities of combinations that models are built with in order to describe the observables correctly. This includes angles, masses, incorporating smearing effects and more. `zfit` and its extensions currently don't contain a lot of different models or losses. The essential parts are contained but more are planned to come in the future. It is expected for them to be continuously added, also depending on the needs that may arise. Furthermore, with `zfit-physics` a repository especially created for content and simple community contributions is available.

Large scale Fits in HEP can be large, both in terms of data as well as in the complexity of the fitting model. With future experiment upgrades an increased amount of data is expected and a fitting library has to cope with that. Complex models and more precise measurements also increase the need for a reliable normalisation, achieved

1308 by a higher number of random samples drawn for the integration. While scalability
1309 to medium scales is already available with `zfit`, the software should not be the
1310 limit in terms of scaling, the computing infrastructure should be. This requires that
1311 on-the-fly normalisation computation can be performed. The extension to huge data
1312 samples with out-of-core computations and to use multiple nodes as well as GPUs
1313 is also a requirement. TF supports this quite well, it was designed for that, but the
1314 explicit implementation inside `zfit` is not yet there.

1315 Most of these current shortcomings have been foreseen and make it into the idea of
1316 `zfit` to become a stable library; a clean implementation with a minimal maintenance
1317 effort is preferred over quantitative content. Additionally, the flexibility and available
1318 base classes allow the user to add these features on top of `zfit` as they are required.

1319 Another future challenge is provided by a significant change of the backend. TensorFlow
1320 2.0 is currently in beta stage and expected to appear somewhere during Summer/Fall
1321 of 2019. A complete restructuring of the library is expected, including a lot of clean up.
1322 While a lot will remain the same, some work will be needed to adjust `zfit` to it.

1323 In summary, `zfit` started not only filling an open gap in the HEP Python ecosystem
1324 but also extends its functionality through the formalisation and flexibility far beyond of
1325 what traditional fitting frameworks are able to do. While still under heavy development,
1326 the current library is already well suitable for a diversity of simple to advanced analysis.

1327 Acknowledgements

1328 I would like to first express my gratitude to Professor Nicola Serra for letting me do this
1329 thesis in his group and being very supportive and trusting overall.

1330 A huge thanks goes to Dr. Albert Puig Navarro, who was not only a great supervisor
1331 but also a core member of the `zfit` project, a co-author, code reviewer, advertiser as well
1332 as main author of `phasespace` and `zfit-amplitude`, both libraries that are very closely
1333 interconnected with `zfit`. Without all the discussions about the large and small details
1334 of the library and use cases, the pure coding and reviewing contributions, this project
1335 would not be on a comparable level and may never grew over a small collection of scripts.

1336 I would like to greatly thank Dr. Rafael Silva Coutinho who is also part of the `zfit`
1337 core team and part of bringing the project to life. His contributions in discussions with a
1338 more user sided view as well as extensive usage with real use cases is, next to code and
1339 documentation contributions, a great support in designing and creating the library.

1340 There are a lot of members in my group that I would like to thank as well for a variety
1341 of things. Be it for discussions about the development of `zfit` and code snippets I like to
1342 specially thank Dr. Abhijit Mathad, Dr. Julian Garcias and Dr. Oliver Lantwin. For
1343 the usage and trial of the library including new features where I like to thank Davide
1344 Lancierini, Sascha Liechti and Martina Ferrillo. And last but not least for the idea of the
1345 library, the technical and user sided experience with an already existing tool and helpful
1346 advices I would like to thank Dr. Andrea Mauri and Michele Atzeni.

1347 A specially thank goes to Prof. Anton Poluektov for several discussions and whose
1348 library `TensorFlow Analysis` was a major inspiration to build `zfit`.

1349 Furthermore, I would also like to thank the scikit-hep community for various design
1350 discussions, especially Dr. Eduardo Rodrigues and Chris Burr.

1351 A great thanks goes to the University of Zurich and CERN for offering the opportunities
1352 to do this kind of research by providing the adequate resources.

1353 Many thanks also go to Google and the TensorFlow development team for open-sourcing
1354 the library and thereby allow libraries like `zfit` to be built.

1355 Last but not least I'd like to thank two persons not only professionally for their support
1356 but also privately for various things, too many as they could even remotely be summarised
1357 here: Simone Steinbrüchel and Patrik Eschle

1358 Appendices

1359 A Likelihood

1360 A reasonable loss is to have a quantity that expresses the *probability of the model given*
1361 *the data*, which can then be maximised. A form of Bayes theorem about probability can
1362 be used in the setting of a model parametrised under θ and data x to express the above

$$P(\theta|x) = \frac{P(x|\theta)P(\theta)}{P(x)} \quad (14)$$

1363 It states that *the probability of θ under our observed data equals to the likelihood of*
1364 *finding our data given θ times the prior of θ over the prior of our data.* While $P(\text{data}|\theta)$
1365 is a probability, fixing the data and varying θ is called the likelihood of θ and denoted as
1366 $\mathcal{L}(\theta)$. If there is no previous knowledge about θ , we can assume a uniform prior, this is
1367 usually the case. The same is done for the data, whose marginal probability acts as a
1368 normalising constant. This means they both introduce constants into our term since they
1369 do not depend on our parametrization θ .

1370 Note that the following derivations rely on a discrete probability distribution as opposed
1371 to a continuous one. This simplifies the argumentation. It can be shown that the results
1372 are equivalently valid for the latter, therefore no strict distinction is made. We start out
1373 with a form of Bayes theorem, which states that the combined probability of two events
1374 are independent of their order.

$$P(A \cap B) = P(B \cap A) \quad (15)$$

1375 The probability of two events to happen can be rewritten in terms of the conditional
1376 $P(A|B)$, read as *the probability of A given B*, and the marginal probability $P(B)$ as

$$P(A \cap B) = P(A|B)P(B) \quad (16)$$

1377 Combining 15 with 16 and rearranging we get

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)} \quad (17)$$

1378 that is usually in this form famously known as Bayes theorem. It states that *the*
1379 *probability of A under B equals the likelihood of B given A times the prior of A over the*
1380 *prior of B.* While $P(B|A)$ is a probability, fixing B and varying A is called the likelihood
1381 of A and denoted as $\mathcal{L}(A)$. If there is no previous knowledge about A, we can assume
1382 a uniform prior. The same is done for B, whose marginal probability usually acts as a
1383 normalizing constant. This means they both introduce constants into our term and are
1384 therefore not of interest later on. They will be left away from now on

$$\mathcal{L}(B) = P(A|B) \quad (18)$$

1385 Given a hypothesis H_0 and a dataset x , the likelihood is the probability that an event
1386 happened under a certain hypothesis

$$\mathcal{L}(H_0) = P(x|H_0). \quad (19)$$

1387 While the probability itself is normalized over x , a likelihood is not. Notice that
 1388 the likelihood of H_0 is a function in H_0 , namely the probability of x under H_0 *with* H_0
 1389 *changing*. So the difference between likelihood and probability is the distinction of what
 1390 is the free parameter and what is the parametrization. We assume now our hypothesis is
 1391 described by a model parametrised with θ , a PDF as defined in 2.

1392 There is actually a distinction between parameters of interest (POI) and nuisance
 1393 parameters when speaking of parametrisation. The latter describe parts of the models
 1394 shape but are not of direct interest in the sense that they do not appear in the hypothesis
 1395 but rather describe well known effects like the width of some smearing. Since they also
 1396 have to be inferred in the same way, we restrict ourselves in the derivation of the likelihood
 1397 to models with exclusively POI. The distinction becomes apparent later in the hypothesis
 1398 testing, which goes beyond the scope of this thesis.

1399 While the likelihood denotes the odds of the model parametrised with θ under the
 1400 observations x , it is, as seen above, equal to the probability density of finding x given the
 1401 parametrisation θ

$$\mathcal{L}(\theta|x) = f_\theta(x) \quad (20)$$

1402 Since the likelihood under data x is the product of likelihoods under each independent
 1403 data point, we can write the likelihood as a product of probability densities as in 4

1404 Using this expression, we can get a maximum likelihood estimate of our parametrisation
 1405 θ

$$\hat{\theta} = \operatorname{argmax}(\mathcal{L}(\theta|x)) \quad (21)$$

1406 In practice, finding the maximum is done using numerical methods. Eq. 4 involves the
 1407 product of many small numbers <1 , not feasible for most computers given their limited
 1408 precision on numbers. The monotony of the logarithm can be used to transform the
 1409 probability densities to log densities whereby the multiplication becomes an addition

$$\operatorname{argmax}(f(\theta|x)) = \operatorname{argmax}(\ln(f(\theta|x))) \quad (22)$$

1410 It is often more convenient to find a minimum instead of a maximum, so that our
 1411 likelihood takes the form as in 5

1412 which is our NLL. Therefore, our likelihood estimation becomes

$$\hat{\theta} = \operatorname{argmin}\left(-\sum_i^n \ln(f(\theta|x_i))\right) \quad (23)$$

1413 This estimation is valid for n measured data points with weight one. We can generalize
 1414 this expression, including event weights. They introduce a power factor but can be taken
 1415 out of the logarithm to have a simple multiplication

$$\hat{\theta} = \operatorname{argmin}\left(-\sum_i^n w_i \cdot \ln(f(\theta|x_i))\right) \quad (24)$$

1416 with w_i being the weight of the i^{th} event. This generalization is not only useful when
 1417 using data points directly to build an unbinned NLL but allows straightforward to bin

1418 the data in the first place and use the bin height as the event weights by choosing an
 1419 appropriate x_i for each bin. Binning can speed up the computation significantly but loses
 1420 slight precision²⁶.

1421 We considered now the likelihood of a model fit to a dataset. The same parameters
 1422 often occur in more than one model, for example the invariant mass of a mother particle.
 1423 To take this into account, a simultaneous fit can be performed, which is simply the
 1424 multiplication of several likelihoods.

$$\mathcal{L}_{f(x)}(\theta|data_0, data_1, \dots, data_n) = \prod_i \mathcal{L}(\theta_i, data_i) \quad (25)$$

1425 where θ_i is a subset of parameters of θ used for the model fit to $data_i$.

1426 The likelihood can be multiplied by constraint terms, most notably an extended
 1427 likelihood estimator (EML) and parameter constraints.

1428 If a parameter changes both shape and the overall normalization of the pdf, an EML
 1429 fit is superior. This situation typically arises if there is a sum of several shapes as
 1430 when adding the background and the signal shape. So for example, taking the model
 1431 $N_{sig} * PDF_{sig} + N_{bkg} * PDF_{bkg}$ and assuming a Poisson distribution of the number of
 1432 events in the data, we can multiply the likelihood by a term

$$\mathcal{L}_{extended} = poiss(N_{tot}, N_{data}) \quad (26)$$

$$= N_{data}^{N_{tot}} \frac{e^{-N_{data}}}{N_{tot}!} \quad (27)$$

1433 where $N_{tot} = N_{sig} + N_{bkg}$. Therefore the EML looks like

$$\hat{\theta} = argmin\left(-\sum_i \ln(f_i(\theta|data_i)) - \sum_i \ln(poiss(N_i, N_{data_i}))\right) \quad (28)$$

1434 where $data_i$ are, as before, different datasets. The weights are inside $data_i$ and taken
 1435 into account as described above. The total number of events is, in general, the sum of the
 1436 weights of all the events.

1437 For certain parameters, prior knowledge is available. If the order of magnitude of the
 1438 knowledge uncertainty is within the expected fitting sensitivity, a constraint term can be
 1439 used to incorporate this. This is nothing else than an additional term the likelihood is
 1440 multiplied with

$$\mathcal{L}(\theta) = \mathcal{L}_{unconstrained} \prod_i f_{constr_i}(\theta) \quad (29)$$

$$= \mathcal{L}_{unconstrained} \mathcal{L}_{constr} \quad (30)$$

1441 as an example, a parameter θ_i that is Gaussian constraint with μ_{constr} and σ_{constr}
 1442 looks like this

²⁶Currently, only unbinned losses are implemented in `zfit` but the extension to binned is already tested and will be there in the future.

$$constr_i = Gauss(\theta_i; \mu_{constr}, \sigma_{constr}) \quad (31)$$

1443 Combining equations 25, 26 and 29 yields for the likelihood in general

$$\mathcal{L} = \mathcal{L}_{f(x)} \cdot \mathcal{L}_{extended} \cdot \mathcal{L}_{constr} \quad (32)$$

1444 The absolute magnitude of the likelihood itself does not have any specific meaning,
 1445 but ratios and scanning over the parameter space of θ allow for statistical interpretation.
 1446 This however goes beyond the scope of `zfit` and is intentionally left to other packages
 1447 like `lauztat` [18].

1448 As mentioned in the beginning, a likelihood as described in 32 is not the only possibility
 1449 to define a loss, though the most common used one. A prominent example is the χ^2 loss,
 1450 which is equivalent to a likelihood in the limiting case if each point is normally distributed.

1451 B Backend

1452 Choosing the right computational backend for model fitting is crucial. In the following
 1453 Sections, the different paradigms and backend designs are introduced. Furthermore, the
 1454 TF library and how it is wrapped inside `zfit` is explained.

1455 B.1 HPC and paradigms

1456 High performance computing (HPC) solves an entirely different problem than high level
 1457 programming languages do. While the flexibility offered by Python is nearly unlimited
 1458 and an incredible strong and comfortable feature, it is only possible due to the dynamic
 1459 nature of the language. For the interpreter of the language that actually translates and
 1460 executes the code to machine level, nearly no assumptions can be made about any object.
 1461 Neither about their size or type nor about their future in the code. This means that no
 1462 previous optimization is possible. For HPC this is the crucial key to success: the more is
 1463 known *previously* of the actual computation, the better the performance will be. The less
 1464 flexibility is available, the more assumptions can be made. This way the layout of arrays,
 1465 parallelisation of computation, caching and more can be efficiently achieved. Since less
 1466 flexibility is not an overall desired characteristics, a good compromise has to be found.
 1467 The distinction between two fundamental paradigm are of importance

1468 **Imperative** Most code is run imperatively: a statement is executed when the line is
 1469 hit. Python, C++ and many more work that way. The advantage is that the
 1470 state changes as the code runs. An addition of two numbers for example will be
 1471 executed right when the command appears. The disadvantage is that it's impossible
 1472 to optimize over more than one step, since the next line of code is not known yet.

1473 **Declarative/Graph** When a statement is supposed to be executed, it is *actually* not
 1474 yet run but somehow remembered. It runs either explicitly when asked for it or
 1475 implicitly if it is needed as a dependency. The latter could also be described as
 1476 “lazy evaluation”. The disadvantage is that the state of the code may looks quite
 1477 undetermined since for example an addition of two numbers is not yet actually

1478 executed when the line is hit. Instead, an object that actually calculates the
1479 computation is usually returned. The paradigm can also be described as a “graph”
1480 based approach, since an execution graph is built. In the special case of HPC, this
1481 can be a computational graph displaying data input, operations and outputs. While
1482 a declarative paradigm is better in terms of raw performance and optimizations in
1483 general, it is also more limited in functionality and cannot offer the flexibility that
1484 the imperative paradigm offers. The workflow is divided into *compilation* or *graph*
1485 *construction* time and *run* time.

1486 In reality, a mixture of this approaches is usually applied on different levels which
1487 some languages and frameworks tending strongly toward one or the other. For example,
1488 any ahead-of-time compilation as used in a C/C++ compiler is kind of a declarative step.
1489 But the language C/C++ itself is imperative. For deep learning frameworks, there is a
1490 strong tendency towards a declarative paradigm, more specific a graph is built. To ease
1491 the difficulty of debugging and making experimentation simpler, an immediate execution
1492 of a graph can make a graph based approach behave imperatively. Some frameworks
1493 therefore offer a mix with this immediate execution. For the maximal performance, a
1494 graph based approach will though be superior for decent complex problems. There are
1495 three major advantages of using a graph.

1496 **Distributed computing** Moore’s law states since over 40 years that the processing
1497 power, which originally equals to the clock speed, of CPUs will double approximately
1498 every 1.5 years. While there is no theoretical foundation for this statement, reality
1499 holds up to it. Since the beginning of the 2000’s, CPUs have reached a clock speed
1500 around gigahertz and started hitting a limit: the thermal dissipation of a CPU goes
1501 with the third power of the clock speed. As a consequence most CPUs today don’t
1502 have a higher clock rate in order to stay efficient. Nonetheless, CPU power keeps
1503 increasing at a similar rate by using more than one CPU unit in parallel. Multicore
1504 systems became common on almost every machine. It is an imperative therefore
1505 that any numerical intensive library makes a good use of distributed computing.
1506 This includes shared memory machines to large scale clusters with multiple nodes.
1507 Moreover, with the recent success of ever growing DNNs and frameworks built to
1508 make easy use of them, vectorial computing devices like GPUs can be used easily as
1509 an alternative to CPUs.

1510 **Mathematical optimization** Since there is not only one step for each computation
1511 available but the whole computation at once, the possibility for mathematical
1512 optimization arise. Most notable, an analytic expression for the derivative is
1513 available through automatic differentiation. This subsequently applies the chain
1514 rule to any operation.

1515 **Caching** A throughout analysis of the computation graph allows for various computa-
1516 tional optimizations. Common sub-expression elimination identifies if the same
1517 operation is executed in several places and substitutes the nodes for a single compu-
1518 tation.

1519 The disadvantage of building a graph first is the overhead to build, analyse and store
1520 the graph. And the additional indirectness, since the main principle is to build the

1521 computation graph once and run it several times. This implies a non-intuitive behaviour
1522 for users coming from an imperative background.

1523 While the graph based approach seems very feasible, there are two fundamental
1524 different assumptions that need be made with graphs.

1525 **Static** A static graph is built once and *cannot* be altered. So if we have for example
1526 a complicated function taking as input a uniform value generator, we build the
1527 graph once and execute it several times. In order to change the input to a normal
1528 distributed random generator, this requires to build the whole graph again with one
1529 operation, the random input generator, changed.

1530 **Dynamic** Dynamic graphs are mutable and can be used if there are a lot of actual
1531 modifications to the graph. Any part can be changed, its not restricted to append-
1532 only.

1533 As an example, TensorFlow offers a static graph while PyTorch uses a dynamic graph.

1534 B.2 Working with TensorFlow

1535 Graph based approaches imply another level of indirectness and need some wrapping in
1536 order to avoid unexpected or inefficient behaviour. In `zfit` this is mostly hidden, exposing
1537 some of it but also removing the obstacles for most use cases.

1538 As an example, we'll look at a task of several random number generations in a row. In
1539 the imperative style this would look like this, where numpy is used to generate a uniform
1540 distribution

```
1541 for _ in range(n_samples):  
    sample = np.random.uniform(...)  
    # do something with the sample
```

1542 In a declarative approach as with TF, the same is achieved by *first* building the actual
1543 operation of the graph and *then* executing the computation

```
1544 sample_op = tf.random.uniform(...)  
for _ in range(n_samples):  
    sample = zfit.run(sample_op)  
    # do something with the sample
```

1545 where `zfit.run(...)` is just the command to execute a certain part of the graph.
1546 The object `sample` is in both cases a numpy array. Note that in the declarative approach,
1547 we created the operation only once but *execute* it multiple times. This indirection can be
1548 surprising for the unaware and a typical mistake is to implement it this way

```
1549 # BAD EXAMPLE  
for _ in range(n_samples):  
    sample_op = tf.random.uniform(...)  
    sample = zfit.run(sample_op)  
    # do something with the sample
```


1550 which just fills up the graph with additional operations that actually all do the same
1551 thing. In order to hide this inconvenience to the user, a caching system is in place in
1552 `zfit` to prevent an accidental re-creation of the operation. In the above example, we can
1553 think of wrapping a function `uniform` (*fictive example*), create the operation on the first
1554 call and use the cached operation afterwards. In pseudo code, this would look like the
1555 following

```
cached_op = None
def uniform(...):
    global cache_op # make "cache_op" assignable
1556     if cached_op is None:
        cached_op = tf.random.uniform(...) # create op
    return cache_op
```

1557 This function can then be used a mixed stile

```
for _ in range(n_samples):
    sample_op = uniform(...) # as created above
1558     sample = zfit.run(sample_op)
    # do something with the sample
```

1559 This hides some of the difficulties. While this example of just random number
1560 generation may seem artificial, a prominent example is the loss that is typically built
1561 using models and data. Calling `value` builds the operation, adds it to the graph and
1562 stores it, so that in subsequent calls, the stored operation is returned instead of a new
1563 one created. While this is very convenient for the user combining the good of two worlds,
1564 it comes with an additional burden for `zfit` of having to cache the operations efficiently.

1565 B.2.1 Caching

1566 In HPC, sometimes certain parts of a computation do not need to be recomputed and
1567 storing the results in a cache for later reuse can improve the performance significantly.
1568 Since `zfit` uses a declarative graph based approach, there are two different caches: For
1569 the operations built on construction time as discussed above and only appearing due
1570 to the graph based approach. Furthermore, intermediate calculations can be cached in
1571 general, which corresponds to the objects built on run time. The latter is commonly used
1572 also in imperative approaches.

1573 As seen before, creating operations with a declarative approach adds each of them
1574 to a global graph, a collection of operations. And since these are only instructions and
1575 not actual computations, there is no need to rebuild the operations and stitching them
1576 together again but rather re-use the previously built instruction.

1577 Therefore `zfit` has a possibility to cache Tensors after they have been built for the
1578 first time inside the object that the Tensor was created in. Since the graph is static and
1579 cannot change, an even small modification to a part of it requires a complete rebuild²⁷.
1580 This can be as less as removing an single, additional term in the model. Any object that
1581 may need to perform such a modification is therefore registered within the caching object
1582 and can notify the cacher in order to invalidate its cache and rebuild any Tensor.

²⁷There are ways of changing the *value* but not the logic, the computation flow.

1583 Numerical results of computations inside the graph can remain the same during several
1584 executions of an operation which can therefore be cached. Opposed to the construction
1585 time caching, numerical caches do also invalidate if a number inside the graph changes,
1586 not only its structure. Caching in a declarative approach is not straightforward since this
1587 implies to replace a node, which can be a large subgraph by itself, with a single value.
1588 Since the graph cannot be changed, this would require a complete rebuild of the graph,
1589 rendering the caching way less efficient. There are two ways to circumvent this problem

1590 **feed_dict** The TF sessions `run` method offers the possibility to overwrite specific nodes
1591 with a value using a `feed_dict`. This means that cached values can simply be stored
1592 and overridden on runtime. At the current stage, `zfit` is not designed to always
1593 control the execution but leaves the freedom to the user to invoke the `Session.run`
1594 method themselves.

1595 **Variable** Since all TF object so far are just operations but not numbers, storing a value
1596 between the runtimes requires a special object. A TF `Variable` can be used for
1597 this case. Its purpose is to store a value but also act as a node in the graph with a
1598 read operation that return the current value. Its value can be changed between the
1599 runtimes in case it has to be updated. This though can have side effects on other
1600 object that still use this cache. The disadvantage here is that the initialization of
1601 the `Variable` can be tricky if it is created inside a control flow such as `tf.while` or
1602 `tf.cond`.

1603 In `zfit` currently a `feed_dict` independent implementation is being tested, but that is
1604 likely to change in the future. One main problem with this kind of caching is the gradient.
1605 If a value is not supposed to change, then its simple since there is no gradient. But if
1606 there is a value that *can* change, the gradients value also needs to be cached. A likely
1607 solution is to wrap the ordinary `Tensor` class from TF and provide a caching `Tensor`, that
1608 automatically takes care of gradient caching as well. Also the `feed_dict` seems like an
1609 efficient solutions, at least for cases without a gradient.

1610 C Implementation

1611 C.1 Spaces definition

1612 A `Space` is either initialized through observables *or* axes and *maybe* also has limits. It *can*
1613 have both, observables and axes, which means there is an order-based, bijective mapping
1614 defined between the two. In general, a `Space` is immutable and adding or changing
1615 anything will always return a copy.

1616 When a user creates a `Space`, observables are used and define with that the coordinate
1617 system. Once a dimensional object, as a model or data, is created with a `Space`, the
1618 order of the observables matters. Since the `Space` at this point only has observables and
1619 does not yet have any axes, when the dimensional object is instantiated, the axis are
1620 created by filling up from zero to $n_{obs} - 1$. This step is crucial and defines the mapping
1621 of internal axes of this dimensional object to the externally used observables. In other
1622 words, every dimensional object has *implicitly* defined axes by counting up to $n_{obs} - 1$,
1623 assigning observables creates a mapping by basically enumerating the observables.

1624 For example we assume a model was instantiated with a `Space` consisting of some
1625 observables and data with the same observables but in a different order. Now the
1626 assignment of observables to the model and the data columns are fixed, therefore it is
1627 well defined how the data has to be reordered if it is passed to the model.

1628 While we used data and a model in the example above, the same is true for limits that
1629 can be used to specify the bounds of integration and more. Since limits can be part of
1630 a `Space`, the reordering is done automatically if the order of the observables or axes is
1631 changed, not in-place but in the returned copy of it.

1632 To help with the accounting of dimensions, `Space` can return a copy of itself with
1633 differently ordered observables. Internally of the `Space`, the axes, if given, and the limits,
1634 if given, are reordered as well. This is crucial in input preprocessing for any dimensional
1635 object since with that each `Space` is ordered in the same way as the observables of that
1636 object.

1637 A subspace can be created from a `Space`. This is a subset of the dimensions of the
1638 `Space`. It can be used for example if a model is composed of lower dimensional models.
1639 This is often the case for functors such as a product PDF as described in Sec. C.6.6.

1640 C.2 General limits

1641 Simple limits are just tuples. But a more general format is needed to express multiple
1642 limits and higher dimensions in a straight forward way. Multiple limits are technically
1643 done with multiple tuples of lower and upper limits for each observable. The `Space` is
1644 handled as *one* domain. So the area of a `Space` is the sum of all the areas of each simple
1645 limit, the integral over a `Space` is the sum of integrals over each simple limit and so on.
1646 Multiple limits are defined separately and are not built from the projections of all limits.
1647 The format is therefore to specify the lower limits and the upper limits in each dimension
1648 as a tuple. Multiple limits contain multiple lower limit tuples and multiple upper limit
1649 tuples.

1650 As an example, a `Space` is created with the observables `x`, `y` and the two limits l_1
1651 and l_2

$$l_1 = (x_0, y_0) \text{ to } (x_1, y_1)$$
$$l_2 = (x_2, y_2) \text{ to } (x_3, y_3)$$

1652 to write the limits in the right order, the upper and lower limits have to be separated
1653 and concatenated, so that the lower and upper limits look like this

$$lower = ((x_0, y_0), (x_2, y_2))$$
$$upper = ((x_1, y_1), (x_3, y_3))$$

1654 By definition of the format, *lower* and *upper* have to have the same length. The
1655 limits for the `Space` is the tuple $(lower, upper)$. Creating this `Space` with `zfit` is done as
1656 follows

1657

```

lower = ((x0, y0), (x2, y2))
upper = ((x1, y1), (x3, y3))
1658 limits = (lower, upper)
multiple_limits = zfit.Space(obs=["x", "y"], limits=limits)

```

1659 The same format is returned by the property `Space.limits`. This is a quite general
1660 format that covers the needs for rectangular shaped limits. However, more advanced
1661 shapes may be necessary, see also 6.1, which will most probably be provided in the future.

1662 Since the order of observable matters and `limits_xy` and `limits_yx` as used in Sec.
1663 4.1 define the same domain (apart from the order of the axis), they can be converted into
1664 each other using the `Space.with_obs` method

```

1665 limits_xy_resorted = limits_yx.with_obs(limits_xy.obs)
limits_xy == limits_xy_resorted # -> True

```

1666 Notice that this created a new `Space` and left `limits_yx` untouched. This method
1667 can also be used to only select a subspace

```

1668 limits_x == limits_xy.with_obs("x") # -> True

```

1669 and therefore go to lower dimensions again.

1670 C.3 Data formats

1671 Currently, the following formats can be read by `Data`.

1672 **ROOT** The standard file format used in HEP analysis. It efficiently stores data samples
1673 and is the native format of the ROOT library. Due to the recent development of
1674 `uproot` [19], it is possible to load these files in Python *without* an installation of
1675 ROOT.

1676 **Pandas DataFrame** Pandas [20] is the most extensive data container in Python used
1677 for data analysis. It provides DataFrames that offer an extensive set of data analysis
1678 tools going from plotting to feature creation and selections. It is the de-facto
1679 standard in Python and has the ability to load from a variety of data formats
1680 including hdf5, csv and more. The possibility to load them directly into `zfit` is
1681 therefore a powerful feature because it allows to do any preprocessing in DataFrames.
1682 `Data` can also be converted *to* DataFrames, which allows to load for example from
1683 ROOT files into `Data`, convert to a DataFrame, apply some preprocessing steps and
1684 then load again into `Data`.

1685 **Numpy** Numpy [21] is the standard computing library in Python that has been around
1686 since a long time. Several libraries, including TF, are inspired by its API and
1687 behaviour design. The numpy arrays are the default way to handle any vectorized
1688 data and are also returned by TF as a result of computations.

1689 **Tensors** `Data` can also take a pure Tensor as input. While this may seem at first glance
1690 the obvious thing to do, it is trickier: a Tensor is, compared to the other data types
1691 *not* fixed per se, since it is only an instruction to compute a certain quantity. While
1692 constants for example behave straight forward and will always return the same, a
1693 `Data` initialized with a random Tensor will produce different data every time it is
1694 called. Therefore, special care has to be taken for this case, from the developer as
1695 well as from the user site.

1696 C.4 Data batching

1697 Small datasets are internally simply converted to a Tensor and attached to the graph.
1698 Large datasets though, which either exceed the memory limit of the computing device
1699 or the limit of the graph size (which is 2 GB), need to make use of batched out-of-core
1700 computation. TF has a data handling class `Dataset`, which provides a performant way to
1701 do batched computations. It incorporates the loading of the batches from disk into the
1702 whole graph as several operations. This allows the runtime to split the execution in order
1703 to asynchronously load a batch of the data and run the graph of an already loaded data
1704 batch.

1705 Another way of on-the-fly computation can be more generally be done with Tensors,
1706 since they can be used to instantiate `Data`. For example, `Data` has a subclass `Sampler`,
1707 which is specialised on this. It allows to evaluate the Tensor and store its value. This
1708 way, the Tensor is only re-evaluated when requested. The `Sampler` acts for example as
1709 the returned `Data` when sampling from a model. This data depends on the model but can
1710 be used like a normal data for example to construct a loss.

1711 C.5 Dependency management

1712 The graph built by TF can be fully accessed, the parent operations of any operation is
1713 available. This enables to detect any dependency by walking along the graph. Or as
1714 TF does internally, to create the gradient. Using this in general can be risky though
1715 since for example caching with a `Variable` changes how the *actual* graph looks like.
1716 Furthermore, it is also time consuming on a larger scale. Inside `zfit`, it is used as an
1717 additional feature to figure out dependencies automatically of certain subgraphs. There
1718 is one type of independents that can change as also described in 4.3.1, `Parameter`, that
1719 other objects can depend on, directly or indirectly. To have a dependency structure that
1720 is independent of the graph, `zfit` has a `BaseDependentsMixin`. A subclass implements
1721 a method of returning the dependents for itself. This can then be done recurrently up
1722 to the independent parameters (see also 4.3.1) that return themselves. Any major base
1723 class implements the appropriate functions but requires for example to have a distinction
1724 between a model that depends only on parameters or also on other models. Both of them
1725 are `Dependents` but the model has to be aware to not only extract dependents from the
1726 parameters but also from the models.

1727 To get the dependents from any object, `get_dependents` can be used, which returns
1728 a set of independent `Parameters`. This is fundamentally different from the `params` that
1729 each model has. The former will return all the *independent* parameters that the model
1730 depends on. This can be any number of *independent* parameters. The latter returns
1731 the exact parameters that are used in the function defining the shape of the model. For

1732 fitting, when tuning the parameters, the `get_dependents` should be used, since the actual
1733 changeable parameters matter. When reading off a value from a model, like the mean
1734 by accessing `mu`, the `params` has to be used. Using the latter for fitting can easily result
1735 in an error if not all of the parameters are independents, since the value of dependents
1736 (including constant) parameters *cannot be set*.

1737 C.6 Base Model

1738 From all the classes, the `BaseModel` and therefore also its subclasses, `BaseFunc` and
1739 `BasePDF`, contain the most logic. The implementation has a few peculiarities that will
1740 be highlighted. It is meant to be used as a base to implement custom models providing
1741 flexibility but ease of use at the same time, as discussed in Sec. 4. Therefore in this class
1742 a structure is provided where everything can be directly controlled *but does not have*
1743 *to be*. The class takes care of anything that is *unambiguous* but maybe cumbersome to
1744 do while leaving the full control to the user. The intended usage as a base class leaves
1745 it to the user to implement himself any method he wants to control directly. Namely
1746 the class provides the guarantee that any of the main methods can be overridden by
1747 changing the implementation of `_method`, the same method name but with a leading
1748 underscore. Furthermore, any direct control can always be given back to the class by
1749 simply raising a `NotImplementedError`. The class acts as if the overridden method was
1750 never called. Furthermore, the base class takes care of some unambiguous handling of
1751 arguments (see below e.g. C.6.3, C.6.4). It is mandatory to decorate each `_method`
1752 with the `supports()` decorator. This allows to specify if the method can handle certain
1753 things, like normalization ranges. By default, this filters the more advanced arguments
1754 and handles them automatically. It is meant to provide a way to still allow specific
1755 implementations and workarounds.

1756 C.6.1 Public methods

1757 The internal logic of the methods `pdf` (and to some extent `func` and `unnormalized_pdf`),
1758 `sample` and all the `integrate` have a very similar layout. Their logic is split into a public
1759 part and more internal methods. We'll start with the outermost method, the public
1760 methods, and follow the subsequent method calls. There is a strict order on what will be
1761 called after which method, we follow the same order here.

1762 A public method is a method starting without an underscore. It serves as the entry
1763 point to a model, asking it to do something. It is supposed to provide a clean API to the
1764 user as well as appropriate documentation of the method. The functional responsibility
1765 of the method is to clean the input which mostly means to take care of the ordering
1766 of dimensions and automatically convert certain input to a more general format. The
1767 following cleaning is done to the input

1768 **limits** Limits for `norm_range` or for sampling and integration can be given as arguments
1769 to certain methods. First, the limits are automatically converted to a `Space` *if it's*
1770 *save to do so*. Namely a simple limit consisting of a tuple in a one dimensional case
1771 is converted, anything else raises an error. Second, the method takes care of sorting
1772 the space in the right way. Since each model has a `Space` with observables assigned
1773 to it with a given order, the given or auto converted `Space` is sorted accordingly.
1774 This assures that the limits are internally in the right order.

1775 **data** As for limits, the data is first converted to the right format, a `Data` *if possible*. This
1776 will convert any input data that coincides with the dimensionality of the model.
1777 While this is convenient since it allows to directly feed numpy arrays and Tensors to a
1778 model, it relies on the correct order of both the data as well as the model observables.
1779 Since this can silently lead to mistakes by using the wrong order, probably in the
1780 future this won't be possible anymore, given that a simple conversion to `Data` can
1781 always be made.

1782 As a second step, the `Data` is sorted according to the observables of the model,
1783 just like the limits. This ordering is done using context managers that revert the
1784 reordering once the data exits the method.

1785 The ordering is a crucial element to allow the direct usage of any object inside a model
1786 while having the matching ordering guaranteed.

1787 C.6.2 Hooks

1788 After the public method, two hooks follow. Assuming we have a public method named
1789 `_method`, the first hook is named `_single_hook_method`. Its purpose is to be directly
1790 called from the public method *or* if the model itself needs to call one of it's own methods.
1791 The second hook called by the first is `_hook_method` and is used for repeated calls by
1792 the very same method that was called. This is useful if for example a scaling factor is
1793 supposed to be applied at the end of the calculations. If a method calls itself recursively,
1794 the scaling is supposed to be applied just once at the very end. The general idea of hooks
1795 is to provide a convenient way to change the behaviour of a method *without* altering
1796 the public method and its input cleaning. It provides the user the possibility to directly
1797 change any input before it is passed further down or any output right before it is returned.
1798 So any kind of advanced, model specific pre- or post-processing or setting of internal flags
1799 can be handled here. Some implementations inside `zfit` make use of this. For example
1800 the implementation of the exponential shape uses the exponential-shift trick²⁸. This
1801 requires to determine the shift before the actual computation method is called. Whatever
1802 modification is applied in a hook, it is important that any hook always calls it's parent
1803 method to allow stacking hooks.

1804 C.6.3 Norm range handling

1805 Following the hooks, `_norm_method` is invoked. It only exists if a normalization range is
1806 used inside the function. This methods responsibility is to automatically take care
1807 of the *normalization logic* if the underlying function does not handle it itself. A
1808 `NormRangeNotImplementedError` can be raised by any deeper nested function which
1809 is caught here. For example in the case of an integral, the method first calls the subse-
1810 quent method and if it catches an error, it splits into two calls: it calculates two integrals,
1811 each one *without* a normalization range, one over the limits to be integrated over and one
1812 over the normalization range. The first is then divided be the second, which is the very
1813 definition of a normalized integral.

²⁸A normalized exponential shape is invariant under translation. This can be used to increase the numerical stability by avoiding large number calculations and keep it around zero.

1814 To illustrate this behaviour in pseudo code, a function `integrate` that takes the limits
1815 and the normalization range as arguments is assumed to exist. `False` as argument to the
1816 normalization range means the calculated integral is unnormalized.

```
try:  
    integral = integrate(limits, norm_range) # pseudo method that  
                                             integrates  
except NormRangeNotImplementedError:  
1817    integral_unnormalized = integrate(limits, False)  
    normalization = integrate(norm_range, False) # integrate over  
                                                  norm_range  
    integral = integral_unnormalized / normalization
```

1818 This allows for a user to implement only a simple integral when overwriting or
1819 registering without the need to worry about its normalization. For most integrals this
1820 would anyway end in two times calling the integration function itself, basically what was
1821 done above. The default behaviour is that the normalization range will be automatically
1822 handled, as described in Appendix C.6. But it still leaves room for the possibility to
1823 implement a method that handles the integral as well as the normalization of it. There are
1824 special cases where this can be achieved with less computations than two times calculating
1825 the integral, such as in the case where the normalization range equals the limits.

1826 C.6.4 Multiple limits handling

1827 Next in the call sequence the method `_limits_method` is invoked. It has the responsibility
1828 of handling multiple limits which are described in 4.1.1. Equivalently to the way `norm_range`
1829 is handled in C.6.3, multiple limits are caught here as well. Consecutive methods are called
1830 inside a `try-except` block in order to catch any `MultipleLimitsNotImplementedError`.
1831 If handling multiple limits is well defined, `_limits_method` is supposed to take care of it.
1832 As an example, the implementation of integration with a limit of n limits is to split them
1833 into n independent `Spaces`, each with only one limit, and call the following methods with
1834 this new `Space`. It is in then a simple matter of summing the results.

1835 Given some multiple limits as a `Space` to be integrated over by invoking the pseudo
1836 `integrate` function, the implementation looks like this

```
try:  
    integral = integrate(limits)  
except MultipleLimitsNotImplementedError:  
1837    integral = 0  
    for limit in limits.iter_limits():  
        integral += integrate(limit)
```

1838 C.6.5 Most efficient method

1839 In `_call_method`, the actual functions are invoked. There are usually several choices for
1840 which function to invoke depending on availability and efficiency. The order chosen there-
1841 fore starts with the most efficient implementation. If that raises a `NotImplementedError`,
1842 the next method is called.

- 1843 • First the `_method` is called and returned if successful. This is a method that by
1844 default raises a `NotImplementedError` but can be overwritten by the user. It is
1845 guaranteed to be executed in this case making the public `method` call seem like a
1846 call to `_method`, module input cleaning and limit handling. This asserts the full
1847 level of flexibility in a model: any major function can completely be overwritten by
1848 this procedure.
- 1849 • If no explicit method is implemented, closely related alternatives are invoked. For
1850 example, if `log_pdf` was called and `_log_pdf` is not overwritten, `_pdf` is invoked
1851 and if implemented, its logarithm is returned. As an other example, if `integrate`
1852 was called and `_integrate` is not implemented, an analytic integral calculation
1853 is performed. If that is not available and also raises an error, more expensive
1854 alternatives are called.
- 1855 • The freedom of `_call_method` is to try simple, but maybe failing alternatives to
1856 return the desired. When all of the above fails, the `_fallback_method` is invoked.
1857 This is the last resort and may be quite a complex function. The calculations that
1858 it returns may be expensive but the method is guaranteed to work. It must not fail
1859 if the `Model` is implemented correctly.

1860 C.6.6 Functors

1861 To implement a function just depending on data, the normal base class as described
1862 above is sufficient. But a model can also depend on other models. Creating compositions,
1863 as for example the `SumPDF` from Sec. 3, requires an additional tweak regarding the
1864 dependencies: a functor does not only depend on its own parameters but also on the sub
1865 models parameters associated with it. This is automatically taken care of by having a
1866 functor base class. Compared to normal models, they often don't need to define their
1867 observables but are inferred from the sub models.

1868 C.7 Sampling techniques

1869 Implementing a sampling from an arbitrary distribution is not a straight forward task.
1870 The generation of uniformly distributed numbers is comparably simple and is used as
1871 a basis for any more advanced sampling technique. Two major ways are often used to
1872 sample from a distribution. If the inverse analytic integral function is known, then this
1873 can be used to transform a uniform distribution to the desired distribution by treating it
1874 as the sample on y . The inverse returns values proportional to the target distribution.
1875 This method is very fast and efficient since every drawn event is used. The downside
1876 is that it requires the inverse analytic integral, which is not available for most shapes,
1877 especially custom ones.

1878 For a model where only the shape and no integral is known, the accept-reject method
1879 can be used. Thereby, samples are randomly generated in the model domain and evaluated.
1880 A random number is drawn between 0 and the maximum of the target shape for each
1881 event. If the value returned by the model is larger than the random number, the event is
1882 accepted. Otherwise it is rejected. The technique is illustrated in Fig. 11.

1883 This can be very inefficient though for peaky distributions since all red values are lost.
1884 An increase in efficiency can be achieved by sampling from a distribution that follows

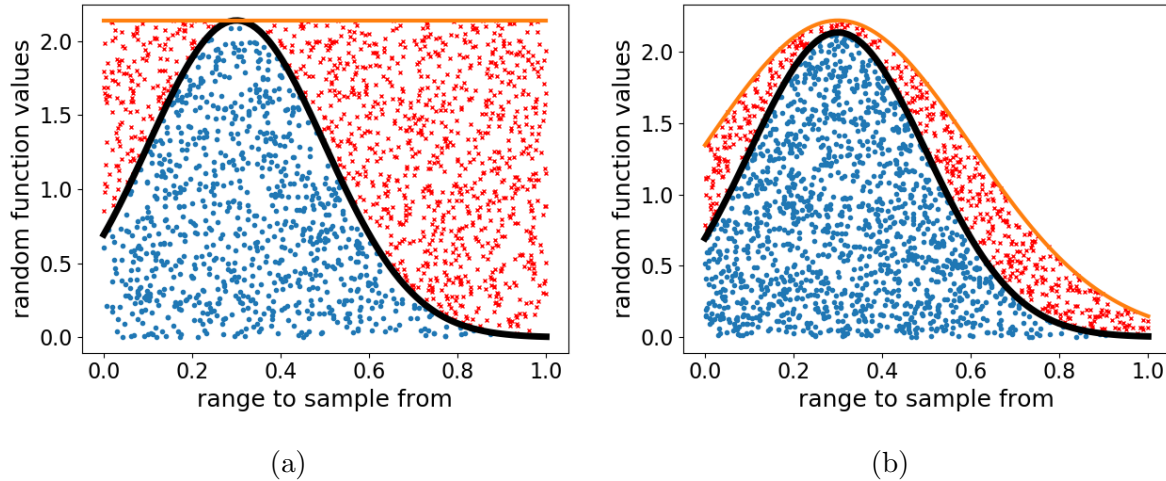


Figure 11: Visualization of the accept reject method. Proposed events are randomly sampled in the valid range. In a) a uniformly sampled y value and in b) a Gaussian shaped y is used to either accept or reject them. The black like is the true shape of the model. The orange line represents the distribution the y were drawn from. Blue values are accepted, red are rejected.

1885 better the target shape, so called “importance sampling”, as shown in 11 on the right. It
 1886 needs to take a little bit more into account, namely

1887 **target** This is the distribution we want to approximate.

1888 **probability** The target probability is the function value of the target shape at a given x.
 1889 While called probability, this does not have to be normalized to anything but be
 1890 proportional to a real probability.

1891 **proposal** The proposed sample is the events that will be either accepted or rejected.
 1892 They are drawn from the proposal distribution.

1893 **weights** This is the probability (or proportional to it) of an event from the proposal
 1894 being drawn from the proposal distribution.

1895 **rnd** A random number drawn uniformly between 0 and 1 to decide whether to accept or
 1896 reject an event.

1897 To approximate the target, a sample is drawn. To accept or reject samples, the
 1898 following is checked

$$accept = prob_i < weight_i \cdot rnd \tag{33}$$

1899 This is only unbiased if $all(prob_i < weight_i)$. Otherwise the distribution will be
 1900 misshaped as shown in Fig. 12

1901 Therefore to be unbiased, the weights have to be scaled enough. This can though lead
 1902 to problems if the weight is significantly lower (even if only in a single point) than the
 1903 target probability. Since this requires a large rescaling, the sampling of the rest of the
 1904 target gets rendered inefficient. Therefore it is important that the proposal distribution
 1905 matches the target reasonably well. Most importantly the maximum and minimum of the
 1906 ratios of the two distributions should be as close as possible.

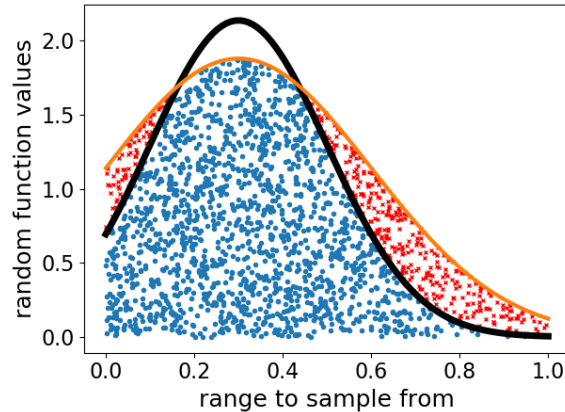


Figure 12: Importance sampling with a wrong scaled weight. The sampled Gaussian distribution (blue) is cut on the top and does not resemble the correct shape.

C.8 Loss defined

The following parts are provided by a loss to enable minimisation

value The actual value of a loss is needed since it is the desired object to be minimised. The method `value` returns the Tensor that can be run with `zfit.run(...)`. This Tensor contains typically the heavy computations.

parameters The loss depends on one or more parameter that is floating. They are associated with models and change the shape thereof. As for a model, all dependents can be retrieved by using `get_dependents`. Changing their values affects the value of the loss.

gradients Calling `gradients` returns the gradients of the loss with respect to the parameters. Gradients provide a helpful tool for the minimisation algorithm.

D Performance studies

D.1 Hardware specification

The measurements are either performed on CPUs only or on an additional GPU. The following hardware and software stack has been used for Fig. 5.

CPU 12 core Intel i7 8850H with 2.60 GHz, 6 cores are virtual using hyper-threading. The available shared memory is 32 GB RAM, which was never even half-way filled. While hyper-threading can be very useful for applications where the bottleneck is not the actual computation by the CPU, for HPC this is often not the case. As the experiments have shown, there is only a minor difference between using 6 or 12 cores. Therefore, only 6 cores, the physical ones, are used in order to quantify the speedup correctly.

1929 **GPU** Mobile Nvidia P1000 with 4GB RAM. It contains the same processing unit as the
1930 consumer GTX 1050 series but is for professional usage and performs more efficient
1931 float64 computations.

1932 It is notable that the price of the GPU and the CPUs are roughly the same, which
1933 allows for some kind of comparison between them.

1934 For Fig. 6, a cluster server with varying hardware was used. Eight cores were requested
1935 for the studies, though the workload of other jobs and the CPU type may have an impact
1936 on the results.

1937 The tests were performed with the TensorFlow version 1.13. It was pre-built by
1938 anaconda and uses the MKL library. There is also a version with the Eigen library
1939 available, tests revealed differing performances for different tasks, around a factor of two
1940 in time. For the GPU version, CUDA 10.0 with cudnn was used.

1941 D.2 Profiling TensorFlow

1942 Code consists of parallelised and serialised parts. While the speed of the former scales²⁹
1943 with the number of cores, the latter does not. The total execution time t is given by

$$t = \sum_i^{n_s} t_s^{(i)} + \sum_i^{n_p} (t_p^{(i)} / n_{cpu} + t_o^{(i)}). \quad (34)$$

1944 where n_s and n_p are the number of *serial* and *parallel* parts, respectively. $t_s^{(i)}$ refers to the
1945 execution time of the i th serial part, the $t_p^{(i)}$ for the parallel parts if executed serial and
1946 $t_o^{(i)}$ denotes the overhead that is needed for each parallel execution.

1947 The serial part consists of

- 1948 • Reading in data from disk.
- 1949 • Setup code such as building a model in `zfit`.
- 1950 • Global operations such as reductions on all values. For example determining whether
1951 a stopping criteria such as the sum of all gradients has gone below threshold is a
1952 serial operation.

1953 while the parallel time $t_{parallel}$ contains usually the heavy computations: evaluating a
1954 function on data whereby the data can be split amongst the cores. The overhead for the
1955 parallel execution time includes

- 1956 • the overhead of creating a new thread for the parallel execution.
- 1957 • the time to move data between the CPUs or even to the GPU.

1958 The serial time t_{serial} consists of

- 1959 • Bottlenecks in I/O or moving data

1960 In order to achieve maximum performance and minimize t

²⁹Ideally. In reality, cores

	1 CPU	6 CPU	GPU
1 x problem	1.0 sec	0.27 sec	0.093 sec
12 x problem	13.4 sec	3.3 sec	1.0 sec

Table 1: Execution time measurement of a loss-like function execution. The complexity of the problem is scaled by n times adding the same loss again to the reduce function.

- 1961 • there should be as few serial code execution time as possible. This is though heavily
1962 limited by the logic and a *certain* amount will always be there.
- 1963 • a minimum of splitting into serial and parallel parts should occur, since each add a
1964 constant t_o term.

1965 This two points are often heavily conflicting and end up with the simple equation to
1966 describe when to parallelise

$$t_s^{(i)} - t_o^{(i)} > t_p^{(i)} / n_{cpu}$$

1967 which reveals that even for large n_{cpu} , the overhead can be the decisive term. Furthermore,
1968 whether it is suitable to execute a piece of code serial or parallel depends on the n_{cpu} .
1969 Together with the difficulty of predicting the overhead time, this leaves *just the decision*
1970 of whether a perfectly parallelisable piece of code actually should be run in parallel as a
1971 heuristic problem. TF uses as a strategy to find the optimal parallelisation to run a small
1972 simulation of the graph, thereby determining the overhead and the number of cores.

1973 Since TF actually executes the computations, any execution time measurement highly
1974 reflects the performance of TF for this task. As TF itself is under active development, the
1975 performance in general is expected to improve in the future.

1976 To get a reasonable estimate of what TF is capable of and somewhat avoid potential
1977 bottlenecks from `zfit`, a dummy test function similar to a loss was written in pure TF.
1978 The function creates three times one million of random numbers and does a few operations
1979 on them before adding and reducing them to a single number. This is added to the
1980 previous calculation in a loop 100 times. There are no I/O bottlenecks and, while not as
1981 an optimal example for TF, it seems reasonable to what can be expected in model fitting.

1982
1983 We can see that the speedup is roughly a factor of 2/3 per core compared to the ideal
1984 case of 1. For example, the time from 1 CPU to 6 CPUs could be expected to decrease by
1985 a factor of 1/6 but does by 1/4 instead. While there are ways of building more efficient
1986 code, the example was chosen to reflect an arbitrarily, non-optimized implementation as
1987 expected to be found in `zfit`, mostly with custom models.

1988 D.3 Additional profiling

1989 Performance studies have been conducted, not shown in Sec. 5. They are displayed here.

1990 References

- 1991 [1] M. Abadi *et al.*, *TensorFlow: Large-scale machine learning on heterogeneous systems*,
1992 2015. Software available from tensorflow.org.

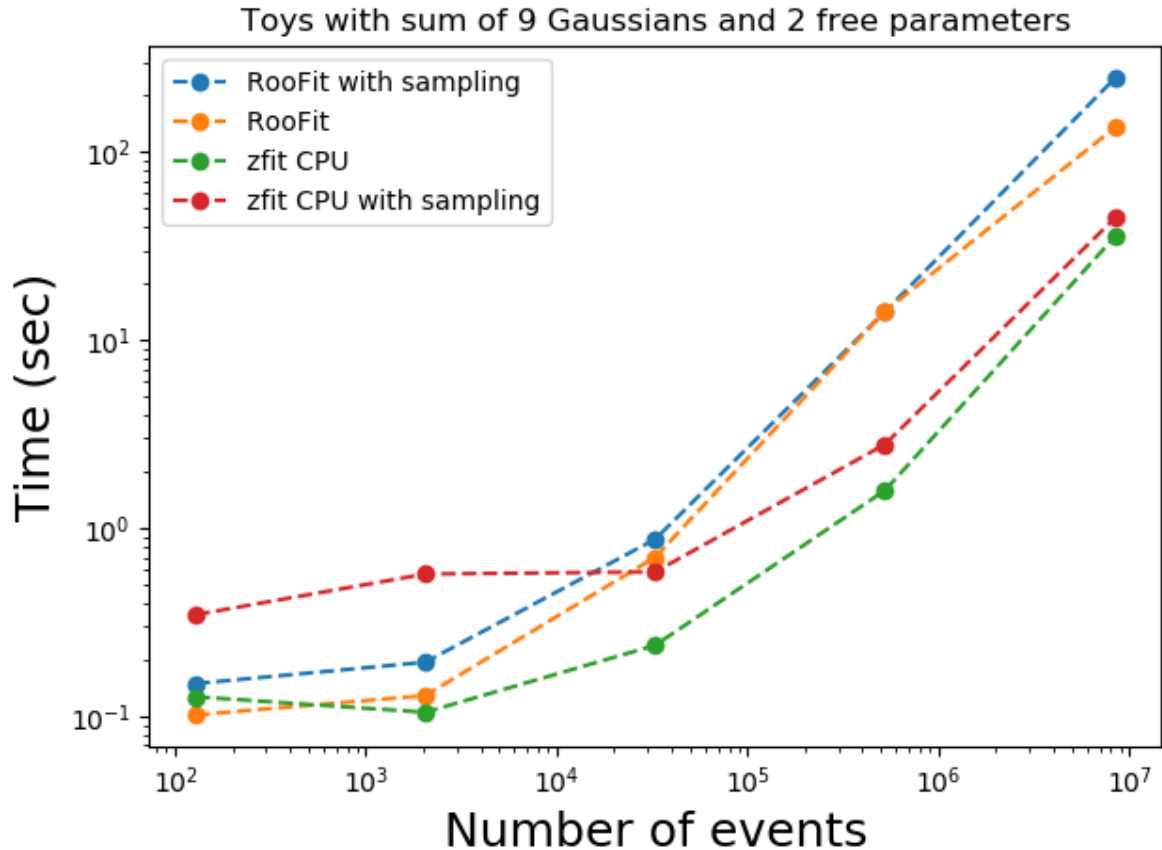


Figure 13: Full toy study with sum of 9 Gaussians and 2 free parameters. We can see that `zfit` s temporary bottleneck in sampling causes an extraordinary increase in execution time mostly for low number of events, but the conclusions and the overall scaling behaviour is still the same as described in Sec. 5.1.

- 1993 [2] A. Paszke *et al.*, *Automatic differentiation in pytorch*, .
- 1994 [3] J. Pivarski, *Non-fork repositories of github users who forked cmssw/aliphysics*,
1995 [https://github.com/jpivarski/2019-06-10-usatlas-argonne-python/blob/](https://github.com/jpivarski/2019-06-10-usatlas-argonne-python/blob/master/01-why-python-in-hep.ipynb)
1996 [master/01-why-python-in-hep.ipynb](https://github.com/jpivarski/2019-06-10-usatlas-argonne-python/blob/master/01-why-python-in-hep.ipynb), 2019.
- 1997 [4] T. E. Oliphant, *Python for scientific computing*, *Computing in Science Engineering*
1998 **9** (2007) 10.
- 1999 [5] M. Newville *et al.*, *lmfit/lmfit-py 0.9.13*, 2019. doi: 10.5281/zenodo.2620617.
- 2000 [6] J. V. Dillon *et al.*, *Tensorflow distributions*, *CoRR* **abs/1711.10604** (2017)
2001 [arXiv:1711.10604](https://arxiv.org/abs/1711.10604).
- 2002 [7] W. Verkerke and D. P. Kirkby, *The RooFit toolkit for data modeling*, *eConf* **C0303241**
2003 (2003) MOLT007, [arXiv:physics/0306116](https://arxiv.org/abs/physics/0306116), [186(2003)].
- 2004 [8] P. Ongmongkolkul *et al.*, *scikit-hep/probfit: 1.1.0*, 2018. doi: 10.5281/zenodo.1477853.
- 2005 [9] Lukas *et al.*, *diana-hep/pyhf v0.0.15*, 2018. doi: 10.5281/zenodo.1464139.

- 2006 [10] J. Bendavid, *Higgsanalysis-combinedlimit*, <https://github.com/bendavid/HiggsAnalysis-CombinedLimit/tree/tensorflowfit>, 2016.
- 2007
- 2008 [11] A. Poluektov, *Tensorflow analysis*, <https://gitlab.cern.ch/poluekt/TensorFlowAnalysis/>, 2017.
- 2009
- 2010 [12] F. James and M. Roos, *Minuit: A System for Function Minimization and Analysis of the Parameter Errors and Correlations*, *Comput. Phys. Commun.* **10** (1975) 343.
- 2011
- 2012 [13] D. P. Kingma and J. Ba, *Adam: A method for stochastic optimization*, 2014.
- 2013 [14] LHCb, R. Aaij *et al.*, *Angular analysis of the $B^0 \rightarrow K^{*0} \mu^+ \mu^-$ decay using 3 fb^{-1} of integrated luminosity*, *JHEP* **02** (2016) 104, [arXiv:1512.04442](https://arxiv.org/abs/1512.04442).
- 2014
- 2015 [15]
- 2016 [16] F. James, *Monte-Carlo phase space*, .
- 2017 [17] BABAR Collaboration, B. Aubert *et al.*, *Measurement of $D^0-\bar{d}^0$ mixing from a time-dependent amplitude analysis of $D^0 \rightarrow K^+ \pi^- \pi^0$ decays*, *Phys. Rev. Lett.* **103** (2009) 211801.
- 2018
- 2019
- 2020 [18] M. Marinangeli, *marinang/lauztat: v1.1.2*, 2019. doi: 10.5281/zenodo.2648147.
- 2021 [19] J. Pivarski *et al.*, *scikit-hep/uproot: 3.6.3*, 2019. doi: 10.5281/zenodo.3239529.
- 2022 [20] W. McKinney, *Data Structures for Statistical Computing in Python*, 2010.
- 2023 [21] T. Oliphant, *NumPy: A guide to NumPy*, USA: Trelgol Publishing, 2006-. [Online; accessed |today|].
- 2024