

# Numerical Computation of $p$ -values with *myFitter*

M. WIEBUSCH\*

*Institute for Theoretical Particle Physics,  
Karlsruhe Institute of Technology (KIT), D-76128 Karlsruhe, Germany*

## Abstract

Likelihood ratio tests are a widely used method in global analyses in particle physics. The computation of the statistical significance ( $p$ -value) of these tests is usually done with a simple formula that relies on Wilks' theorem. There are, however, many realistic situations where Wilks' theorem does not apply. In particular, no simple formula exists for the comparison of models that are not *nested*, in the sense that one model can be obtained from the other by fixing some of its parameters. In this paper I present methods for efficient *numerical* computations of  $p$ -values, which work for both nested and non-nested models and do not rely on additional approximations. These methods have been implemented in a publicly available C++ framework for maximum likelihood fits called *myFitter* and have recently been applied in a global analysis of the Standard Model with a fourth generation of fermions.

---

\*email: wiebusch@particle.uni-karlsruhe.de

# 1 Introduction

Even though the LHC experiments have, so far, not found any clear signs for physics beyond the Standard Model (SM) they already put strong constraints on the favourite SM extensions of many theorists. For example, the SM with a (perturbative) sequential fourth generation of fermions (SM4) has recently been excluded at the  $5\sigma$  level by a combination of Higgs and electroweak precision data [1] (see also [2, 3]). Other models with additional fermions or even some constrained versions of Supersymmetry may follow soon.

In this situation some thoughts should be spent on the methods and criteria by which we decide if a certain model is ruled out. A well-established technique in (frequentist) statistical analyses is the method of *likelihood ratio tests*. (For an introduction see e.g. [4] or the statistics chapter of [5].) In this method two models are compared with a test statistic constructed from the ratio of their likelihood functions. Wilks' theorem states that under certain assumptions the test statistic is distributed according to the well-known  $\chi^2$ -distribution [6]. In this case the relation between the likelihood values at the best fit points and the statistical significance ( $p$ -value) of the corresponding hypothesis test is described by the normalised lower incomplete gamma function.

There are, however, also many realistic scenarios where Wilks' theorem does not hold and the probability density function of the test statistic is not known analytically. One example is the case of likelihood ratio tests where the two models to be compared are not *nested*, meaning that one model can not be obtained from the other by fixing some of its parameters. This problem was encountered in the above-mentioned analyses of the Standard Model (SM) with a fourth generation of fermions [1–3]. In these analyses it is not possible to regard the SM with three fermion generations as a limiting case of the SM with four generations due to non-decoupling contributions of chiral fermions in electroweak precision observables and Higgs production and decay rates. Another case where analytical formulae for  $p$ -values are not reliable is the situation where some of the parameters of a model are bounded, in the sense that they are only allowed to float within a certain range. Most notably, this applies to analyses where systematic errors are treated within the *RFit* scheme [7], i.e. by introducing so-called nuisance parameters with a limited range.

When analytic formulae fail one has to resort to numerical methods, and the computation of  $p$ -values is no exception. The brute-force method is to generate a large sample of random *toy measurements* distributed according to the prediction of the null hypothesis. For each toy measurement the value of the test statistic is computed and compared to the value obtained from the actual data. With a large enough sample we can then estimate the probability that the value of the test statistic is larger than a certain number, usually chosen to be the value of the test statistic obtained from the observed data. This probability is called *statistical significance* or  $p$ -value of the test. Unfortunately, the computational cost of the required numerical simulations can

be rather high, especially when the  $p$ -value is small. In this paper I discuss some methods for improving the efficiency of numerical computations of  $p$ -values. These methods have been applied in [1, 3], where, based on the constraints from Higgs searches and electroweak precision observables, likelihood ratio tests comparing the SM with three and four fermion generations were performed. The methods are also implemented in a publicly available code called *myFitter*, which I present in this paper.

The paper is organised as follows: in Sec. 2 I describe the general mathematical setup and the definitions of the test statistics for nested and non-nested models. In Sec. 3 I sketch the derivation of Wilks' theorem and discuss its range of applicability. In Sec. 4 I explain the strategy for improving the efficiency of numerical computations of  $p$ -values. The *myFitter* framework and the implementation of the methods from Sec. 4 are discussed in Sec. 5. Performance tests of the *myFitter* code are presented in Sec. 6. I conclude in Sec. 7.

## 2 General Setup

Let  $\mathbf{X} = (X_1, \dots, X_n)$  be a set of experimental observables. In frequentist statistics we regard observables as random variables distributed according to some probability density function (PDF). A statistical model with free parameters  $\boldsymbol{\xi} = (\xi_1, \dots, \xi_k)$  is therefore described by a function  $f(\mathbf{x}, \boldsymbol{\xi})$ , which must be a PDF for any fixed value of  $\boldsymbol{\xi}$  and considered as a function of  $\mathbf{x}$  only:

$$\int d^n \mathbf{x} f(\mathbf{x}, \boldsymbol{\xi}) = 1 \quad . \quad (1)$$

The problem of statistical inference is to draw conclusions about the parameters  $\boldsymbol{\xi}$  from a given set of measurements  $\mathbf{x}$  of the observables  $\mathbf{X}$ .

In global analyses like [1–3] the observables come from many different collider experiments and the parameters  $\boldsymbol{\xi}$  to be determined are the fundamental parameters of some theory of particle physics (the SM or extensions thereof). In this situation the function  $f$  usually does not depend on the parameters  $\boldsymbol{\xi}$  directly. For example, a cross section  $\sigma$  is measured by counting the number  $N$  of events that pass certain cuts and dividing by the integrated luminosity  $\mathcal{L}$  and the selection efficiency  $\varepsilon$ . If the theory is realised with parameters  $\boldsymbol{\xi}$  we denote the predicted value of the cross section as  $\tilde{\sigma}(\boldsymbol{\xi})$ . The integrated luminosity and selection efficiency are usually constants which, to a good approximation, do not depend on the theory parameters. Since  $N$  follows a Poisson distribution with mean value  $\mathcal{L}\varepsilon\tilde{\sigma}(\boldsymbol{\xi})$  the distribution of the *measured* value of the cross section,  $\sigma = N/(\mathcal{L}\varepsilon)$ , depends on  $\boldsymbol{\xi}$  only through the *predicted* cross section  $\tilde{\sigma}(\boldsymbol{\xi})$ .

These considerations motivate us to write the function  $f$  in the following way:

$$f(\mathbf{x}, \boldsymbol{\xi}) = \exp[-\frac{1}{2}D(\tilde{\mathbf{x}}(\boldsymbol{\xi}), \mathbf{x})] \quad , \quad (2)$$

where  $\tilde{\mathbf{x}}(\boldsymbol{\xi})$  are the “predicted” values of the observables and  $D$  is a function which we shall call the *input function*. The goal of this re-writing is to cleanly separate information about the theoretical model from information about the experimental uncertainties: the theory is described by the function  $\tilde{\mathbf{x}}$ , which maps parameters to observables, and the experimental uncertainties are described by the function  $D(\tilde{\mathbf{x}}, \mathbf{x})$ , which is  $-2$  times the logarithm of the probability density for measuring values  $\mathbf{x}$  if the “true” values are  $\tilde{\mathbf{x}}$ . For example, in many situations the experimental errors are Gaussian, independent of the true values  $\tilde{\mathbf{x}}$  and described by a covariance matrix  $V$ . In this case we have

$$D(\tilde{\mathbf{x}}, \mathbf{x}) = (\tilde{\mathbf{x}} - \mathbf{x})^\top V^{-1} (\tilde{\mathbf{x}} - \mathbf{x}) + n \ln(2\pi) \quad , \quad (3)$$

which, if substituted in (2), gives the usual expression for a correlated Gaussian probability density with central value  $\tilde{\mathbf{x}}(\boldsymbol{\xi})$ . Up to a constant term,  $D(\tilde{\mathbf{x}}(\boldsymbol{\xi}), \mathbf{x})$  is simply the  $\chi^2$ -value associated with the parameters  $\boldsymbol{\xi}$ , input data  $\mathbf{x}$  and covariance matrix  $V$ .

For the methods presented in this paper, the exact definition of the functions  $\tilde{\mathbf{x}}$  and  $D$  is not important, as long as  $D$  has the following properties:

1. For any given  $\tilde{\mathbf{x}}$  the input function  $D$  must satisfy the normalisation condition

$$\int_S d^n \mathbf{x} e^{-\frac{1}{2} D(\tilde{\mathbf{x}}, \mathbf{x})} = 1 \quad . \quad (4)$$

2. For any given  $\mathbf{x}$ , and considered as a function of  $\tilde{\mathbf{x}}$ , the input function  $D(\tilde{\mathbf{x}}, \mathbf{x})$  must be bounded from below and have its unique absolute minimum at  $\tilde{\mathbf{x}} = \mathbf{x}$ .
3. For any given  $\tilde{\mathbf{x}}$ , and considered as a function of  $\mathbf{x}$ , the input function  $D(\tilde{\mathbf{x}}, \mathbf{x})$  must be bounded from below and have its unique absolute minimum at  $\mathbf{x} = \tilde{\mathbf{x}}$ .

The first property guarantees that (1) is satisfied. The second property guarantees that, if parameters  $\hat{\boldsymbol{\xi}}$  exist with  $\tilde{\mathbf{x}}(\hat{\boldsymbol{\xi}}) = \mathbf{x}$ , the maximum likelihood estimate of the parameters is indeed  $\hat{\boldsymbol{\xi}}$ . The third property can be regarded as a definition of the term “predicted value”: if the theory is realised with some parameters  $\boldsymbol{\xi}$ , the most likely outcome of a measurement of the observables  $\mathbf{x}$  should be  $\mathbf{x} = \tilde{\mathbf{x}}(\boldsymbol{\xi})$ .

Note that, without any modifications to the model, the third property does *not* hold in the presence of systematic errors, since a systematic error is an offset between the true value of an observable and its most likely measured value. This offset is the same each time the measurement is performed and does therefore not average out when the measurement is repeated many times. This results in a difference between  $\mathbf{x}$  and the maximum of the distribution of the random variables  $\mathbf{X}$ . The central idea of the RFit method [7] is that systematic errors should not be treated as errors at all, but as unknown theory parameters, so-called *nuisance parameters*, that may vary within a certain range. In a way, the presence of a systematic error means that theorists and experimentalists are simply not talking about the same quantity. Since the difference

between the two quantities can neither be modeled nor measured it has to be treated as an additional model parameter, but with a limited range of possible values. Thus, the third assumption *does* hold if systematic errors are treated within the *RFit* scheme, i.e. by introducing a nuisance parameter for each source of systematic errors.

The results of hypothesis tests for a certain theoretical model should not depend on the way we parametrise the model. To make this parametrisation-independence manifest it is convenient to define the *theory manifold* as the image of the function  $\tilde{\mathbf{x}}$ :

$$M = \{\tilde{\mathbf{x}}(\boldsymbol{\xi}) \mid \boldsymbol{\xi} \in \Omega\} \quad , \quad (5)$$

where  $\Omega \subset \mathbb{R}^k$  is the *parameter space* (i.e. the set of allowed parameter values) of the model. Different parametrisations of the same model are represented by different functions  $\tilde{\mathbf{x}}$  and parameter spaces  $\Omega$ , but always have the same theory manifold.

The general procedure for a likelihood ratio test (LRT) with *nested* models may now be described as follows: given certain experimental data  $\mathbf{x}$ , we first maximise the PDF  $f(\mathbf{x}, \boldsymbol{\xi})$  with respect to the parameters  $\boldsymbol{\xi}$ . This is equivalent to minimising the function  $D(\tilde{\mathbf{x}}, \mathbf{x})$  with respect to  $\tilde{\mathbf{x}}$  on the theory manifold  $M$ :

$$f^{\max}(\mathbf{x}) = \exp[-\frac{1}{2}D^{\min}(\mathbf{x})] \quad \text{with} \quad D^{\min}(\mathbf{x}) = \min\{D(\tilde{\mathbf{x}}, \mathbf{x}) \mid \tilde{\mathbf{x}} \in M\} \quad . \quad (6)$$

Next, we consider a *constrained* version of the model, which is usually obtained from the original model by fixing some of its parameters. However, with the notion of theory manifolds at hand, we can be more general and simply require that the theory manifold  $M_c$  of the constrained model is a subset of  $M$ :

$$M_c \subset M \quad . \quad (7)$$

Maximising the likelihood for the constrained model we get

$$f_c^{\max}(\mathbf{x}) = \exp[-\frac{1}{2}D_c^{\min}(\mathbf{x})] \quad \text{with} \quad D_c^{\min}(\mathbf{x}) = \min\{D(\tilde{\mathbf{x}}, \mathbf{x}) \mid \tilde{\mathbf{x}} \in M_c\} \quad . \quad (8)$$

Now we construct a test statistic  $S$  from the ratio of the two maximum likelihood values:

$$S(\mathbf{x}) = -2 \ln \frac{f_c^{\max}(\mathbf{x})}{f^{\max}(\mathbf{x})} = D_c^{\min}(\mathbf{x}) - D^{\min}(\mathbf{x}) \quad . \quad (9)$$

To perform the actual test, we choose a certain realisation of the constrained model as null hypothesis. Let  $\boldsymbol{\xi}_0$  be the corresponding parameters and  $\tilde{\mathbf{x}}_0 = \tilde{\mathbf{x}}(\boldsymbol{\xi}_0)$ . The statistical significance, or *p*-value, of the test is obtained by considering an ensemble of *toy measurements*  $\mathbf{x}$  distributed according to the PDF  $f(\mathbf{x}, \boldsymbol{\xi}_0)$  and computing the probability that  $S(\mathbf{x})$  is larger than some threshold value  $S_0$ :

$$p = \int d^n \mathbf{x} f(\mathbf{x}, \boldsymbol{\xi}_0) \theta(S(\mathbf{x}) - S_0) \quad , \quad (10)$$

where  $\theta$  denotes the Heavyside step-function. If, in the real experiments, the data  $\mathbf{x}_0$  was measured, one typically takes the maximum likelihood estimates for  $\mathbf{x}_0$  in

the constrained model as null hypothesis, i.e. one chooses  $\xi_0$  so that  $\tilde{\mathbf{x}}_0 \in M_c$  and  $D(\tilde{\mathbf{x}}_0, \mathbf{x}_0) = D_c^{\min}(\mathbf{x}_0)$ . Then one performs the test with  $S_0 = S(\mathbf{x}_0)$ .

Note that the definition (10) is manifestly independent of the parametrisation of the models: the factor  $f(\mathbf{x}, \xi_0)$  is fixed by the null hypothesis and the test statistic is defined in terms of the functions  $D^{\min}$  and  $D_c^{\min}$  whose definitions (6) and (8) depend on the manifolds  $M$  and  $M_c$ , but not on their parametrisation. This parametrisation-independent language allows us to easily generalise the definition (10) for the case of (a large class of) non-nested models. Consider two models with theory manifolds  $M_1$  and  $M_2$  such that  $M_1 \not\subset M_2$  and  $M_2 \not\subset M_1$ . However, we assume that the relation between the PDFs of the two models and their respective theory manifolds is still given by (2) with the *same* input function  $D$ . This usually holds for global fits in particle physics, where a model imposes certain relations between the predicted observables, but the random distribution of the measured quantities is fixed by the predicted values, irrespective of the model under consideration. In this case we can simply combine the two theories into one by joining their theory manifolds,

$$M \equiv M_1 \cup M_2 \quad \Rightarrow \quad M_1, M_2 \subset M \quad , \quad (11)$$

and do a LRT as described above, with  $M$  as the full theory and either  $M_1$  or  $M_2$  as the constrained theory. Let

$$D_1^{\min}(\mathbf{x}) = \min\{D(\tilde{\mathbf{x}}, \mathbf{x}) \mid \tilde{\mathbf{x}} \in M_1\} \quad , \quad D_2^{\min}(\mathbf{x}) = \min\{D(\tilde{\mathbf{x}}, \mathbf{x}) \mid \tilde{\mathbf{x}} \in M_2\} \quad . \quad (12)$$

Then the test statistic for testing  $M_2$  against  $M$  is

$$S_2(\mathbf{x}) = \begin{cases} D_2^{\min}(\mathbf{x}) - D_1^{\min}(\mathbf{x}) & \text{for } D_1^{\min}(\mathbf{x}) < D_2^{\min}(\mathbf{x}) \\ 0 & \text{otherwise} \end{cases} \quad (13)$$

and the test statistic for testing  $M_1$  against  $M$  is obtained by exchanging  $D_1^{\min}$  and  $D_2^{\min}$ . Assume without restriction that for the measured data  $\mathbf{x}_0$  we have

$$D_1^{\min}(\mathbf{x}_0) \leq D_2^{\min}(\mathbf{x}_0) \quad . \quad (14)$$

Then  $S_1(\mathbf{x}_0)$  is zero and the LRT for  $M_1$  (using  $S_1(\mathbf{x}_0)$  as threshold value for the test) has a  $p$ -value of 1. So, only the LRT for the model which describes the data less well (i.e. gives a bigger value for  $D_i^{\min}(\mathbf{x}_0)$ ) can have a  $p$ -value smaller than one.

### 3 Analytical Formulae for $p$ -values

In many cases the computation of  $p$ -values in LRTs is trivial due to a theorem by Wilks [6]. It states that the test statistic  $S$  from (9) follows a  $\chi^2$  distribution with  $\dim(M) - \dim(M_c)$  degrees of freedom *if* the models are nested and the maximum

likelihood estimates  $\hat{\boldsymbol{\xi}}(\mathbf{x})$  of the parameters  $\boldsymbol{\xi}$  follow a Gaussian distribution. In this case the  $p$ -value (10) is given by

$$p = 1 - P_{\nu/2}(S_0/2) \quad , \quad (15)$$

where  $\nu = \dim(M) - \dim(M_c)$  is the difference of dimensions of the theory manifolds (usually equal to the number of parameters that were fixed) and  $P$  denotes the normalised lower incomplete Gamma function.

In global analyses in particle physics Wilks' theorem is commonly used, but the validity of its underlying assumptions are rarely discussed. For PDFs of the form (2) the requirements for Wilks' theorem translate to certain assumptions about the function  $D$  and the theory manifolds  $M$  and  $M_c$ . These assumptions are:

**Gaussianity.** The function  $D(\tilde{\mathbf{x}}, \mathbf{x})$  only depends on the difference  $\tilde{\mathbf{x}} - \mathbf{x}$  and is quadratic in this difference.

**Linearity.** The theory manifolds  $M$  and  $M_c$  are *hyperplanes*.

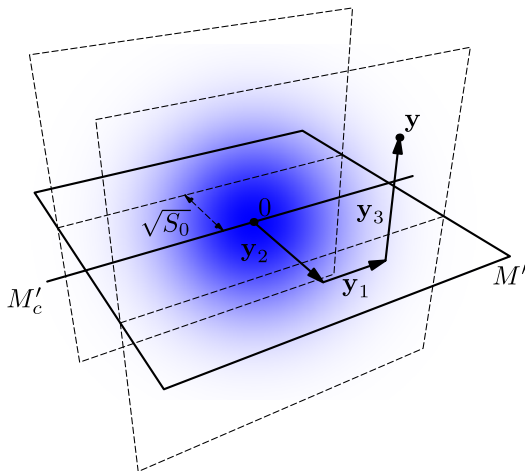
**Nestedness.** The constrained theory is a subset of the full theory:  $M_c \subset M$ .

The first assumption, combined with the properties of  $D$  discussed in Sec. 2, implies that  $D$  is of the form (3), i.e. that the experimental errors are Gaussian. The second assumption is invalid if experimental errors are large, so that the curvature of the theory manifolds can not be neglected. It also fails if some parameters of the model have upper or lower bounds, so that the corresponding manifold does not extend to infinity. The last assumption is invalid if none of the two models to be compared can be considered as a special case of the other model.

The derivation of Wilks' theorem from the assumptions above will be instructive for our discussion of numerical methods in the next section, so I will briefly sketch it here. The first step is to perform an affine-linear coordinate transformation in the space of observables, which maps  $\tilde{\mathbf{x}}_0$  (the predicted observables under the null hypothesis) to the origin and changes the PDF to an  $n$ -dimensional normal distribution. In other words, we introduce new coordinates  $\mathbf{y} \equiv \mathbf{y}(\mathbf{x})$ , so that  $\mathbf{y}(\tilde{\mathbf{x}}_0) = 0$  and

$$f(\mathbf{x}, \boldsymbol{\xi}_0) = \frac{1}{(2\pi)^{n/2}} e^{-\|\mathbf{y}\|^2/2} \quad \Rightarrow \quad D(\tilde{\mathbf{x}}_0, \mathbf{x}) = \|\mathbf{y}\|^2 + n \ln(2\pi) \quad , \quad (16)$$

The linear part of this transformation is easily constructed by diagonalising the matrix  $V$  from (3) and then scaling the new coordinates appropriately. We see that, up to a constant term, the function  $D$  is simply the squared euclidean norm of the vector  $\mathbf{y}$  (denoted as  $\|\mathbf{y}\|^2$ ). Let  $M'$  and  $M'_c$  denote the images of  $M$  and  $M_c$  under the coordinate transformation  $\mathbf{y}$ . Since  $M$  and  $M_c$  both contain  $\tilde{\mathbf{x}}_0$ , the hyperplanes  $M'$  and  $M'_c$  both contain the origin and are therefore linear subspaces. Consequently, the functions  $D^{\min}$



**Figure 1:** Derivation of Wilks’ theorem. The blue colour indicates the probability density in the transformed space of observables  $\mathbf{y}$ . It is an  $n$ -dimensional normal distribution. The  $\theta$  function in (10) vanishes in the region between the planes indicated by the thin dashed lines.

and  $D_c^{\min}$  (see Eq. 6 and 8) are simply the squared euclidean length of the component of  $\mathbf{y}$  perpendicular to  $M'$  and  $M'_c$ , respectively. For  $n = 3$ ,  $\dim(M) = 2$  and  $\dim(M_c) = 1$  the situation is depicted in Fig. 1. We may write the vector  $\mathbf{y}$  as a sum of three orthogonal vectors  $\mathbf{y}_1$ ,  $\mathbf{y}_2$  and  $\mathbf{y}_3$  with  $\mathbf{y}_1 \in M'_c$  and  $\mathbf{y}_2 \in M$ . The test statistic  $S$  is then:

$$S(\mathbf{x}) = \|\mathbf{y}_2 + \mathbf{y}_3\|^2 - \|\mathbf{y}_3\|^2 = \|\mathbf{y}_2\|^2 \quad . \quad (17)$$

In terms of the coordinates  $\mathbf{y}$  the integral from (10) becomes

$$p = \frac{1}{(2\pi)^{n/2}} \int d^n \mathbf{y} e^{-\|\mathbf{y}\|^2/2} \theta(\|\mathbf{y}_2\|^2 - S_0) \quad (18)$$

In other words, the  $p$ -value is the integral of an  $n$ -dimensional normal distribution in the region outside an (infinitely long) “hyper-cylinder” defined by  $\|\mathbf{y}_2\|^2 > S_0$ . In Fig. 1, this “cylinder” is the region between the planes indicated by dashed lines. The integral in (18) can easily be computed. The integrals over the components  $\mathbf{y}_1$  and  $\mathbf{y}_3$  are just Gaussian integrals and give an overall factor of  $(2\pi)^{(n-\nu)/2}$ . Introducing spherical coordinates in the  $\nu$ -dimensional subspace corresponding to the component  $\mathbf{y}_2$  and exploiting rotational symmetry immediately leads to (15).

## 4 Numerical Calculation of $p$ -values

In the last section we have seen how Wilks’ theorem emerges from geometric arguments if the models under consideration satisfy three assumptions, which we called *gaussianity*, *linearity* and *nestedness*. In practice, these assumptions are rarely satisfied exactly. Usually, they are only more or less valid approximations. If we do not want to rely on



these approximations we have to resort to numerical integration methods to calculate  $p$ -values. To make these methods efficient it is a good idea to take some or all of the approximations as a starting point and optimise the numerical integration for the case where they are valid. In the following, we will use Monte Carlo integration with importance sampling to compute the integral (10), and construct sampling densities which are optimal for models that satisfy gaussianity and linearity.

To compute the integral (10) with the importance sampling method, we generate a large number  $N$  of sample points  $\mathbf{x}_i$  according to some sampling distribution  $\rho$ . The integral (10) is then estimated as

$$p = \frac{1}{N} \sum_{i=1}^N \frac{f(\mathbf{x}_i, \boldsymbol{\xi}_0)}{\rho(\mathbf{x}_i)} \theta(S(\mathbf{x}_i) - S_0) \quad . \quad (19)$$

To reduce the statistical error of this estimate one has to choose the function  $\rho$  as similar as possible to the integrand, so that the terms in the sum are (ideally) all of the same size. A common approach (especially if  $f(\cdot, \boldsymbol{\xi}_0)$  is a Gaussian distribution) is to choose  $\rho(\mathbf{x}) \approx f(\mathbf{x}, \boldsymbol{\xi}_0)$  and let the numerics take care of the theta function in the integrand. For large  $p$ -values this is a viable option, but if  $p$  is small most sample points give a contribution of zero to the integrand and the numerical integration becomes very inefficient. (Remember that each evaluation of the test statistic  $S(\mathbf{x})$  requires the computation of  $D^{\min}(\mathbf{x})$  and  $D_c^{\min}(\mathbf{x})$  which, in general, has to be done by numerical minimisation.) For small  $p$ -values, the efficiency of the integration can be significantly improved by choosing a sampling density  $\rho$  which avoids the region where the theta function is zero (i.e. the region between the dashed planes in Fig. 1). Knowledge about the geometric properties of the theory manifolds can be used to construct such a sampling density. In the numerical methods proposed in this paper, we assume gaussianity and linearity for the purpose of constructing the sampling density, but make no approximations when computing the  $p$ -value.

To see how this works, let us start with the case where the nestedness assumption is still valid. For definiteness, we choose a parametrisation so that

$$M = \{\tilde{\mathbf{x}}(\boldsymbol{\xi}) \mid \boldsymbol{\xi} \in \Omega\} \quad , \quad M_c = \{\tilde{\mathbf{x}}(\boldsymbol{\xi}) \mid \boldsymbol{\xi} \in \Omega \wedge \xi_1, \dots, \xi_\nu = 0\} \quad . \quad (20)$$

Let  $\boldsymbol{\xi}_0$  denote again the parameters under the null hypothesis. Now we define the hyperplanes  $H$  and  $H_c$  as tangent planes on  $M$  and  $M_c$  at the point  $\tilde{\mathbf{x}}_0 = \tilde{\mathbf{x}}(\boldsymbol{\xi}_0)$ :

$$\begin{aligned} H &= \left\{ \tilde{\mathbf{x}}_0 + \mathbf{h} \mid \mathbf{h} \in \text{span} \left( \left. \frac{\partial \tilde{\mathbf{x}}(\boldsymbol{\xi})}{\partial \xi_1} \right|_{\boldsymbol{\xi}=\boldsymbol{\xi}_0}, \dots, \left. \frac{\partial \tilde{\mathbf{x}}(\boldsymbol{\xi})}{\partial \xi_k} \right|_{\boldsymbol{\xi}=\boldsymbol{\xi}_0} \right) \right\} \quad , \\ H_c &= \left\{ \tilde{\mathbf{x}}_0 + \mathbf{h} \mid \mathbf{h} \in \text{span} \left( \left. \frac{\partial \tilde{\mathbf{x}}(\boldsymbol{\xi})}{\partial \xi_{\nu+1}} \right|_{\boldsymbol{\xi}=\boldsymbol{\xi}_0}, \dots, \left. \frac{\partial \tilde{\mathbf{x}}(\boldsymbol{\xi})}{\partial \xi_k} \right|_{\boldsymbol{\xi}=\boldsymbol{\xi}_0} \right) \right\} \quad . \end{aligned} \quad (21)$$

By construction the function  $D(\tilde{\mathbf{x}}_0, \mathbf{x})$ , considered as a function of  $\mathbf{x}$ , has a minimum at  $\mathbf{x} = \tilde{\mathbf{x}}_0$  (see Sec. 2). Consequently, we define the matrix  $V^{-1}$  as the Hessian matrix at

that minimum:

$$(V^{-1})_{ij} = \left. \frac{\partial^2 D(\tilde{\mathbf{x}}_0, \mathbf{x})}{\partial x_i \partial x_j} \right|_{\mathbf{x}=\tilde{\mathbf{x}}_0}. \quad (22)$$

As in the derivation of Wilks' theorem, we now perform an affine-linear coordinate transformation which maps  $\tilde{\mathbf{x}}_0$  to zero and transforms  $V^{-1}$  to a unit matrix. To this end, we define

$$y_i \equiv y_i(\mathbf{x}) = \frac{1}{\sigma_i} [O(\mathbf{x} - \tilde{\mathbf{x}}_0)]_i \quad (\text{no sum over } i) \quad , \quad (23)$$

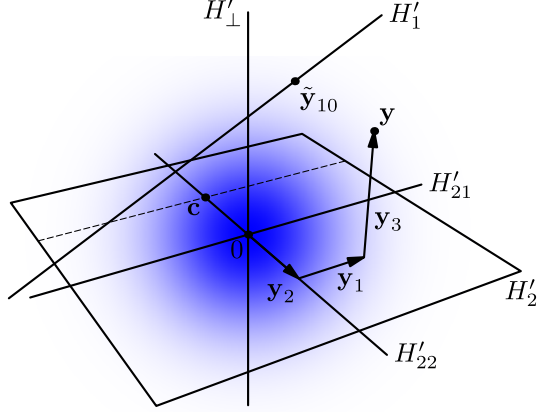
where  $O$  is an orthogonal matrix chosen so that  $OVO^T = \text{diag}(\sigma_1^2, \dots, \sigma_n^2)$  with positive eigenvalues  $\sigma_1^2, \dots, \sigma_n^2$ . Let  $H'$  and  $H'_c$  denote the images of the hyperplanes  $H$  and  $H_c$ , respectively, under this coordinate transformation. Since  $H$  and  $H_c$  contain  $\tilde{\mathbf{x}}_0$ ,  $H'$  and  $H'_c$  must contain the origin and are therefore linear subspaces of  $\mathbb{R}^n$ . Any vector  $\mathbf{y}$  may thus be decomposed into three orthogonal components  $\mathbf{y}_1$ ,  $\mathbf{y}_2$  and  $\mathbf{y}_3$  with  $\mathbf{y}_1 \in H'_c$  and  $\mathbf{y}_2 \in H$  (see Fig. 1). If the assumptions of gaussianity and linearity were satisfied exactly, the theta function in (19) would vanish if and only if  $\|\mathbf{y}_2\|^2 < S_0$  and we should not waste sample points on this region. If gaussianity and linearity are only approximations, we should be more careful and use a sampling density  $\rho$  which is small, but non-vanishing for  $\|\mathbf{y}_2\|^2 < S_0$ . A choice which can still be sampled efficiently is

$$\rho(\mathbf{x}) = J e^{-\frac{1}{2}\|\mathbf{y}_1\|^2} e^{-\frac{1}{2}\|\mathbf{y}_3\|^2} \begin{cases} a \|\mathbf{y}_2\|^\alpha & , \quad \|\mathbf{y}_2\|^2 < S_0 \\ b e^{-\frac{1}{2}\|\mathbf{y}_2\|^2} & , \quad \|\mathbf{y}_2\|^2 \geq S_0 \end{cases} \quad , \quad (24)$$

where  $J = \prod_{i=1}^n \sigma_i$  is the Jacobian of the coordinate transformation. The parameters and  $a, b, \alpha \geq 0$  may be tuned to improve the efficiency of the numerical integration (subject, of course, to the constraint that the PDF  $\rho$  is properly normalised).

We see that the sampling density  $\rho$  factorises into three terms which only depend on the components  $\mathbf{y}_1$ ,  $\mathbf{y}_2$  and  $\mathbf{y}_3$ , respectively. The task of generating points distributed according to  $\rho$  thus reduces to the task of generating components  $\mathbf{y}_1$ ,  $\mathbf{y}_2$  and  $\mathbf{y}_3$  distributed according to the respective factors. For  $\mathbf{y}_1$  and  $\mathbf{y}_3$  these factors are Gaussian, so generating the components  $\mathbf{y}_1$  and  $\mathbf{y}_3$  is trivial. The generation of the component  $\mathbf{y}_2$  requires special care.

Before we address this problem let us talk about the case of non-nested models. Assume that we have two theory functions  $\tilde{\mathbf{x}}_1 \equiv \tilde{\mathbf{x}}_1(\boldsymbol{\xi})$  and  $\tilde{\mathbf{x}}_2 \equiv \tilde{\mathbf{x}}_2(\boldsymbol{\eta})$  with parameter spaces  $\Omega_1$  and  $\Omega_2$ , respectively, of arbitrary and possibly different dimension. Our null hypothesis is that theory 2 is realised with parameters  $\boldsymbol{\eta}_0 \in \Omega_2$ . We therefore approximate the theory manifold  $M_2$  by its tangent hyperplane  $H_2$  at  $\tilde{\mathbf{x}}_{20} \equiv \tilde{\mathbf{x}}_2(\boldsymbol{\eta}_0)$ . Since  $M_2$  is no subset of  $M_1$  we have to approximate  $M_1$  by its tangent hyperplane at some other parameters  $\boldsymbol{\xi}_0 \in \Omega_1$ . If  $\boldsymbol{\eta}_0$  is the maximum likelihood estimate of some measured data  $\mathbf{x}_0$ , i.e.  $\boldsymbol{\eta}_0 = \hat{\boldsymbol{\eta}}(\mathbf{x}_0)$ , an obvious choice for  $\boldsymbol{\xi}_0$  would be the maximum likelihood estimate of  $\mathbf{x}_0$  in theory 1, i.e.  $\boldsymbol{\xi}_0 = \hat{\boldsymbol{\xi}}(\mathbf{x}_0)$ . In any case we define hyperplanes



**Figure 2:** Orthogonal decomposition of a three-dimensional sample space for non-nested models. The tangent hyperplane  $H'_1$  of model 1 is one-dimensional and the tangent hyperplane  $H'_2$  of model 2 is two-dimensional. The thin dashed line is the projection of  $H'_1$  onto  $H'_2$ . The blue colour indicates the probability density for the toy observable vector  $\mathbf{y}$ .

$H_1$  and  $H_2$  analogous to (21):

$$\begin{aligned}
 H_1 &= \left\{ \tilde{\mathbf{x}}_{10} + \mathbf{h} \mid \mathbf{h} \in \text{span} \left( \left. \frac{\partial \tilde{\mathbf{x}}_1(\boldsymbol{\xi})}{\partial \xi_1} \right|_{\xi=\xi_0}, \left. \frac{\partial \tilde{\mathbf{x}}_1(\boldsymbol{\xi})}{\partial \xi_2} \right|_{\xi=\xi_0}, \dots \right) \right\}, \\
 H_2 &= \left\{ \tilde{\mathbf{x}}_{20} + \mathbf{h} \mid \mathbf{h} \in \text{span} \left( \left. \frac{\partial \tilde{\mathbf{x}}_2(\boldsymbol{\eta})}{\partial \eta_1} \right|_{\eta=\eta_0}, \left. \frac{\partial \tilde{\mathbf{x}}_2(\boldsymbol{\eta})}{\partial \eta_2} \right|_{\eta=\eta_0}, \dots \right) \right\}, \quad (25)
 \end{aligned}$$

with  $\tilde{\mathbf{x}}_{10} = \tilde{\mathbf{x}}_1(\boldsymbol{\xi}_0)$  and  $\tilde{\mathbf{x}}_{20} = \tilde{\mathbf{x}}_2(\boldsymbol{\eta}_0)$ . We define the matrix  $V^{-1}$  as in (22) and construct coordinates  $\mathbf{y} \equiv \mathbf{y}(\mathbf{x})$  according to (23), but with  $\tilde{\mathbf{x}}_0$  replaced by  $\tilde{\mathbf{x}}_{20}$ . Let  $H'_1$  and  $H'_2$  be the images of  $H_1$  and  $H_2$ , respectively, under the coordinate transformation  $\mathbf{y}$  and let  $\tilde{\mathbf{y}}_{10} = \mathbf{y}(\tilde{\mathbf{x}}_{10})$ . The image of  $\tilde{\mathbf{x}}_{20}$  under  $\mathbf{y}$  is the origin. Since  $H_2$  contains  $\tilde{\mathbf{x}}_{20}$  the hyperplane  $H'_2$  contains the origin and is thus a linear subspace of  $\mathbb{R}^n$ .  $H'_1$ , on the other hand, contains  $\tilde{\mathbf{y}}_{10}$  but not necessarily the origin, so it is not a linear subspace. For  $n = 3$ , a two-dimensional  $H'_2$  and a one-dimensional  $H'_1$ , this situation is depicted in Fig. 2.

We see that for non-nested models the boundaries of the region with  $S(\mathbf{x}) < S_0$  are curved, even if a linear approximation is used for the theory manifolds. This makes it harder to construct a sampling distribution which avoids this region. We shall try it anyway: let  $H'_{21} \subset H'_2$  be the subspace obtained by shifting  $H_1$  by  $-\tilde{\mathbf{y}}_{10}$  (so that it contains the origin) and projecting it onto  $H'_2$ . Furthermore, let  $H'_{22}$  be the orthogonal complement of  $H'_{21}$  in  $H_2$ . Finally, let  $H'_\perp$  be the orthogonal complement of  $H'_2$  in  $\mathbb{R}^n$ . The projection of  $H'_1$  onto  $H'_{22}$  is then a single point  $\mathbf{c}$ , which can be obtained by projecting  $\tilde{\mathbf{y}}_{10}$  onto  $H'_{22}$ . Any vector  $\mathbf{y}$  can now be written as the sum of three orthogonal components  $\mathbf{y}_1$ ,  $\mathbf{y}_2$  and  $\mathbf{y}_3$  with  $\mathbf{y}_1 \in H'_{21}$ ,  $\mathbf{y}_2 \in H'_{22}$  and  $\mathbf{y}_3 \in H'_\perp$ . (See Fig. 2.) The distance between  $\mathbf{y}$  and  $H'_1$  is larger than  $\|\mathbf{c} - \mathbf{y}_2\|$  since projecting any vector on a

lower-dimensional subspace reduces its length and the projection of any vector pointing from  $\mathbf{y}$  to  $H'_1$  onto the subspace  $H'_{22}$  is  $\mathbf{c} - \mathbf{y}_2$ .

Now recall the test statistic  $S_2$  from (13), which we constructed to test theory 2 against the “union” of theories 1 and 2. In the approximation where the theory manifolds  $M_1$  and  $M_2$  are equal to the hyperplanes  $H_1$  and  $H_2$  the functions  $D_1^{\min}$  and  $D_2^{\min}$  satisfy

$$D_2^{\min}(\mathbf{x}) = \|\mathbf{y}_3\|^2 + n \ln(2\pi) \quad , \quad D_1^{\min}(\mathbf{x}) \geq \|\mathbf{c} - \mathbf{y}_2\|^2 + n \ln(2\pi) \quad . \quad (26)$$

Thus,  $S_2(\mathbf{x}) < S_0$  holds if

$$\|\mathbf{y}_3\|^2 < S_0 + \|\mathbf{c} - \mathbf{y}_2\|^2 \quad . \quad (27)$$

Note, however, that this condition is sufficient but not necessary.  $S_2(\mathbf{x})$  must be smaller than  $S_0$  in the region defined by (27), but it may also be smaller than  $S_0$  outside this region. Nonetheless, a good choice for the sampling density  $\rho$  will be one which avoids the region defined by (27). Analogous to (24), this density can be constructed as follows:

$$\rho(\mathbf{x}) = J e^{-\frac{1}{2}\|\mathbf{y}_1\|^2} e^{-\frac{1}{2}\|\mathbf{y}_2\|^2} \begin{cases} a \|\mathbf{y}_3\|^\alpha & , \quad \|\mathbf{y}_3\|^2 < S_0 + \|\mathbf{c} - \mathbf{y}_2\|^2 \\ b e^{-\frac{1}{2}\|\mathbf{y}_3\|^2} & , \quad \|\mathbf{y}_3\|^2 \geq S_0 + \|\mathbf{c} - \mathbf{y}_2\|^2 \end{cases} \quad , \quad (28)$$

where  $J = \prod_{i=1}^n \sigma_i$  is again the Jacobian of the coordinate transformation  $\mathbf{y}$ . After requiring that  $\rho$  is properly normalised, there are still two free parameters which can be tuned to improve the efficiency of the numerical integration. As in (24), the density  $\rho$  factorises into three terms. The first two are Gaussian and only depend on the components  $\mathbf{y}_1$  and  $\mathbf{y}_2$ , respectively. Generating components  $\mathbf{y}_1$  and  $\mathbf{y}_2$  with the correct statistical distribution is therefore trivial. A new complication in (28) is that the last factor, i.e. the distribution of  $\mathbf{y}_3$ , now depends on  $\mathbf{y}_2$ . This simply means that we have to generate a value for  $\mathbf{y}_2$  *before* we generate  $\mathbf{y}_3$ .

The remaining problem is to generate a random vector  $\mathbf{z}$  of some dimension  $m$  distributed according to the PDF

$$\rho'(\mathbf{z}) = \begin{cases} a \|\mathbf{z}\|^\alpha & , \quad \|\mathbf{z}\|^2 < \Delta^2 \\ b e^{-\frac{1}{2}\|\mathbf{z}\|^2} & , \quad \|\mathbf{z}\|^2 \geq \Delta^2 \end{cases} \quad (29)$$

with some  $\Delta > 0$ . (In (24) we have  $\mathbf{z} = \mathbf{y}_2$ ,  $m = \nu$  and  $\Delta^2 = S_0$  while in (28) we have  $\mathbf{z} = \mathbf{y}_3$ ,  $m = \dim(H'_\perp)$  and  $\Delta^2 = S_0 + \|\mathbf{c} - \mathbf{y}_2\|^2$ .) We first note that the PDF  $\rho'$  is rotationally invariant. The length  $r = \|\mathbf{z}\|$  of the vector  $\mathbf{z}$  is then distributed according to a PDF

$$\tilde{\rho}'(r) = \frac{2\pi^{m/2}}{\Gamma(m/2)} r^{m-1} \begin{cases} a r^\alpha & , \quad r^2 < \Delta^2 \\ b e^{-\frac{1}{2}r^2} & , \quad r^2 \geq \Delta^2 \end{cases} \quad . \quad (30)$$

We may write this as

$$\tilde{\rho}'(r) = f \tilde{\rho}'_{<}(r) + (1 - f) \tilde{\rho}'_{>}(r) \quad (31)$$

where  $f \in [0, 1]$  is a free parameter and

$$\tilde{\rho}'_{<}(r) = \frac{m + \alpha}{\Delta^{m+\alpha}} \theta(r) \theta(\Delta - r) \quad , \quad \tilde{\rho}'_{>}(r) = \frac{r^{m-1} e^{-\frac{1}{2}r^2} \theta(r - \Delta)}{2^{(m-2)/2} \Gamma(m/2) (1 - P_{m/2}(\frac{1}{2}\Delta^2))} \quad (32)$$

are PDFs normalised to 1. Here,  $P_{m/2}$  is the normalised lower incomplete Gamma function. Since  $\tilde{\rho}'_{<}$  and  $\tilde{\rho}'_{>}$  are normalised, the parameter  $f$  is just the fraction of sample points that will be put in the “inner region” with  $r < \Delta$ . For a given  $f$ , the corresponding values of  $a$  and  $b$  are

$$a = \frac{f \Gamma(m/2) (m + \alpha)}{2 \pi^{m/2} \Delta^{m+\alpha}} \quad , \quad b = \frac{1 - f}{(2 \pi)^{m/2} (1 - P_{m/2}(\frac{1}{2}\Delta^2))} \quad . \quad (33)$$

By integrating  $\tilde{\rho}'_{<}$ ,  $\tilde{\rho}'_{>}$  and  $\tilde{\rho}'$  from 0 to  $r$  we obtain the *cumulative distribution functions* (CDFs)

$$\begin{aligned} \text{CDF}_{\tilde{\rho}'_{<}}(r) &= \frac{r^{m+\alpha}}{\Delta^{m+\alpha}} \quad , \quad \text{CDF}_{\tilde{\rho}'_{>}}(r) = \frac{P_{m/2}(\frac{1}{2}r^2) - P_{m/2}(\frac{1}{2}\Delta^2)}{1 - P_{m/2}(\frac{1}{2}\Delta^2)} \\ \Rightarrow \text{CDF}_{\tilde{\rho}'}(r) &= \begin{cases} f \text{CDF}_{\tilde{\rho}'_{<}}(r) & , \quad r < \Delta \\ f + (1 - f) \text{CDF}_{\tilde{\rho}'_{>}}(r) & , \quad r \geq \Delta \end{cases} \quad . \end{aligned} \quad (34)$$

To generate random variables  $r$  distributed according to  $\tilde{\rho}'$  we need the inverse of  $\text{CDF}_{\tilde{\rho}'}$ :

$$\begin{aligned} \text{CDF}_{\tilde{\rho}'_{<}}^{-1}(q) &= \Delta q^{1/(m+\alpha)} \quad , \quad \text{CDF}_{\tilde{\rho}'_{>}}^{-1}(q) = \sqrt{2 P_{m/2}^{-1}(q + (1 - q) P_{m/2}(\frac{1}{2}\Delta^2))} \\ \Rightarrow \text{CDF}_{\tilde{\rho}'}^{-1}(q) &= \begin{cases} \text{CDF}_{\tilde{\rho}'_{<}}^{-1}(\frac{q}{f}) & , \quad q < f \\ \text{CDF}_{\tilde{\rho}'_{>}}^{-1}(\frac{q-f}{1-f}) & , \quad q \geq f \end{cases} \quad , \end{aligned} \quad (35)$$

where  $P_{m/2}^{-1}$  is the inverse of the normalised lower incomplete Gamma function. Random vectors  $\mathbf{z}$  distributed according to  $\rho$  may now be generated in the following way: First generate a vector  $\mathbf{z}'$  according to a  $m$ -dimensional normal distribution. Then pick a uniformly distributed random variable  $q \in [0, 1]$  and set  $r = \text{CDF}_{\tilde{\rho}'}^{-1}(q)$ . The variable  $r$  is then distributed according to  $\tilde{\rho}'$ . The vector  $\mathbf{z}$  with the correct random distribution is  $\mathbf{z} = (r/\|\mathbf{z}'\|)\mathbf{z}'$ .

## 5 Introducing *myFitter*

The ideas for the numerical computation of  $p$ -values outlined in the last section have been implemented in a publicly available code called *myFitter*. The source code is available at Hepforge [8]. Detailed documentation is included in the source distribution. Here I just want to provide a brief description of the user interface and discuss some details of the implementation.

`myFitter` is a C++ class library and makes extensive use of inheritance and polymorphism to separate the tasks of fitting a model to experimental data and computing  $p$ -values from the tasks of implementing the observables (as functions of the model’s parameters) or the input function  $D$  (as a function of the observables). The main classes the user will have to deal with are:

**Model** This is the base class for all models implemented by the user. It essentially represents the theory function  $\tilde{\mathbf{x}}$  from earlier sections, i.e. the map from the model’s parameter space to the space of observables. The base class provides functionality for storing “current” values of parameters, observables and derivatives of observables with respect to the parameters, setting ranges in which parameters are allowed to float or fixing them (so that they do not float at all). It can also randomly sample the parameter space and build up a dictionary of parameter values and the corresponding observable values. This dictionary can be used to find good starting points for numerical minimisations of the input function. To implement their own model, the user has to subclass `Model` and overload the method `calc()` which computes the observables based on the current values of the parameters. They may also overload the method `calc_deriv()`, which calculates the derivatives of all observables with respect to all parameters. The default implementation uses simple numerical differentiation.

**InputComponent** This is the base class for objects that represent terms in the input function  $D$  (see Sec. 2). Each input component represents the contribution from one or more observables  $x_i$  to the input function. To calculate the value of the input function, the contributions of all `InputComponent` objects are added up. This is done by another class, `InputFunction`, which acts as a container for `InputComponent` objects. Derived classes of `InputComponent` must overload the method `calc( $\tilde{\mathbf{x}}, \mathbf{x}$ )`, which takes two vectors as arguments (the first being the “predicted” values of the observables and the second being the “measured” ones) and returns the contribution of the term to the input function. Additionally, the methods `calc_deriv()` and `get_hessian()` must be implemented, which calculate the derivatives with respect to the  $\tilde{x}_i$  and the Hessian matrix for the minimum at  $\mathbf{x} = \tilde{\mathbf{x}}$ . Ready-to-use implementations for the most common input components are also available. These classes are: `GaussianIC` (for single observables with a Gaussian and possibly systematic errors), `AsymmetricGaussianIC` (for single observables with asymmetric Gaussian error bars and possibly systematic errors) and `CorrelatedGaussianIC` (for several observables with Gaussian errors and a correlation matrix).

**Fitter** Objects of this type are responsible for fitting the parameters of models (represented by `Model` objects) to experimental data (represented by an `InputFunction` object) and for computing  $p$ -values by numerical integration. Each `Fitter` object contains an `InputFunction` object which is accessible through

the `input_function()` method and must be “filled” with `InputComponent` objects before any fits can be done. Once the input function is initialised, fits can be performed with the `local_fit()` and `global_fit()` methods. As arguments, these methods take the `Model` object to be fitted and (optionally) a vector of central values for the observables. If no central values are given, the defaults from the `input_function()` are used. The difference between these methods is that `local_fit()` uses the current values of the model parameters as starting point for the minimisation of the input function, while `global_fit()` uses the dictionary created by a previous call to the model’s `scan()` method. The  $p$ -values for likelihood ratio tests of nested and non-nested models can be calculated with the methods `calc_nested_lrt_pvalue()` and `calc_lrt_pvalue()`, respectively. As arguments, these two methods take the models to be compared. Note that, for `calc_nested_lrt_pvalue()` to work, the second model must be a restricted version of the first, i.e. a copy of the first object with some additional parameters fixed. In addition to these methods, the `Fitter` class contains numerous options and flags that control the accuracy and various other aspects of the minimisation and integration routines. These options are described in the package documentation. Most notably, the  $p$ -value integrations can be *parallelised without additional programming efforts* by the user.

Both, for the case of nested and non-nested models, the efficiency of the integration can be improved further by *adaptive* integration techniques, where the shape of the sampling density  $\rho$  is tuned *during* the actual integration. For the adaptation, the implementation in `myFitter` uses the OmniComp/Dvegas package [9] by Nikolas Kauer, which implements the VEGAS algorithm [10] and was developed in the context of [11, 12]. Thanks to OmniComp, parallelised integration is fully supported.

To maximise the likelihood function, `myFitter` uses a custom implementation of the BFGS method for numerical optimisation [13–16]. The optimisation terminates successfully when the length of the gradient of the likelihood function drops below a certain value configurable by the user. Other optimization algorithms can be implemented by subclassing the `Minimizer` class and assigning an instance of this class to the `Fitter` object via the `Fitter::minimizer()` method. The problem of minimising a function of bounded parameters (i.e. of parameters that have an upper or lower limit) is solved in the usual way by smoothly and invertably mapping the real axis  $\mathbb{R}$  to the allowed range of the parameter. Internally, `myFitter` does this with the function

$$g : (-\infty, \infty) \rightarrow (0, \infty), \quad x \mapsto f(x) = \frac{1}{2}(x + \sqrt{x^2 + 1}) \quad . \quad (36)$$

## 6 Performance Tests

The performance of the `myFitter` method for the numerical integration of Eq. 10 was compared with more generic methods in three tests using simple toy models. All al-



ternative methods use the coordinate transformation (23) to transform the PDF  $f$  to a normal distribution. The simplest method, referred to as *no adaptation* in the following, just uses importance sampling with a normal distribution as sampling density. The other two methods map the integration volume (i.e. the  $\mathbb{R}^n$ ) to the unit hypercube  $[0, 1]^n$  by using

$$t_i = \frac{1}{2} \operatorname{Erf} \frac{y_i}{\sqrt{2}} + \frac{1}{2} \quad , \quad i = 1, \dots, n \quad (37)$$

as integration variables and use the VEGAS algorithm [10] to perform the integration over the variables  $t_i$ . The VEGAS algorithm is most efficient when the features of the integrand are aligned with the coordinate axes. In one variant, called *aligned VEGAS* in the following, we perform a rotation which aligns the tangent hyperplanes of the theory manifolds with the coordinate axes before mapping to the unit cube. This usually leads to an integrand whose features are aligned with the coordinate axes. In a second variant, which we call *misaligned VEGAS*, we choose the rotation so that the theory manifolds are *not* aligned with the coordinate axes. The misaligned VEGAS method is the best possible method when no information about the theory manifolds can be used.

In the first test we study the performance of the four integration methods in the context of a model with a curved theory manifold. To this end we consider a model with seven observables  $x_1, \dots, x_7$  and four parameters  $\xi_1, \dots, \xi_4$ . The theory function  $\tilde{\mathbf{x}}$  is given by

$$\tilde{\mathbf{x}}(\boldsymbol{\xi}) = (\xi_1, \xi_2, \xi_3, \xi_4, 0, 0, -(\xi_1^2 + \xi_2^2 + \xi_3^2 + \xi_4^2)\lambda) \quad . \quad (38)$$

where  $\lambda$  is a fixed number which controls the curvature of the theory manifold. As input function we use the expression (3) for Gaussian errors with a unit covariance matrix:

$$D(\tilde{\mathbf{x}}, \mathbf{x}) = \sum_{i=1}^7 (\tilde{x}_i - x_i)^2 + 7 \ln(2\pi) \quad . \quad (39)$$

The constrained version of this model is defined by fixing  $\xi_2$  to zero and  $\xi_1$  to some other value  $r$ . The test statistic  $S$  is then defined according to (9). We take  $\mathbf{x}_0 = (0, 0, 0, 0, 1, 1, 1)$  as the actually measured data and perform the test with  $S_0 = S(\mathbf{x}_0)$ . Different choices for  $r$  lead to different values of  $S_0$  and thus to different  $p$ -values.

For two values of  $\lambda$ , the value of  $r$  was varied to obtain  $p$ -values roughly corresponding to 2, 3, 4 and 5 standard deviations. The  $p$ -value was then computed numerically with *myFitter* and the three alternative methods to a relative precision of 1%. The number of integrand evaluations needed by each method are summarised in Tab. 1. The  $p$ -values obtained by applying Wilks theorem are also shown. We see that adaptive methods always lead to a significant speedup and that the *myFitter* method performs best in all cases. For three standard deviations or less the two VEGAS methods still compete rather well with *myFitter*. At four standard deviations *myFitter* is faster than the VEGAS methods by a factor of 3 for large curvature and a factor of 10 for small curvature. At five standard deviations only *myFitter* is able to compute the  $p$ -value with



$\lambda$	$p$ -value	$p$ -value (Wilks)	<i>myFitter</i> [ $10^3$ ]	aligned VEGAS [ $10^3$ ]	misaligned VEGAS [ $10^3$ ]	no adaptation [ $10^3$ ]
0.1	$5.1 \cdot 10^{-2}$	$5.6 \cdot 10^{-2}$	40	60	90	200
	$2.9 \cdot 10^{-3}$	$3.3 \cdot 10^{-3}$	150	240	360	–
	$6.1 \cdot 10^{-5}$	$8.1 \cdot 10^{-5}$	210	2100	3000	–
	$5.4 \cdot 10^{-7}$	$7.3 \cdot 10^{-7}$	270	–	–	–
1.0	$5.2 \cdot 10^{-2}$	$8.8 \cdot 10^{-2}$	50	60	70	200
	$3.0 \cdot 10^{-3}$	$5.8 \cdot 10^{-3}$	180	210	210	–
	$8.2 \cdot 10^{-5}$	$20.1 \cdot 10^{-5}$	700	2100	1800	–
	$8.1 \cdot 10^{-7}$	$29.7 \cdot 10^{-7}$	6000	–	–	–

**Table 1:** Results of the test with a curved theory manifold. The curvature is controlled by  $\lambda$  (see Eq. 38). The parameter  $\xi_1$  was fixed to different values in the constrained model, leading to the  $p$ -values shown in the second column (which roughly correspond to 2, 3, 4 and 5 standard deviations). The  $p$ -values obtained by applying Wilks theorem is shown in the third column. The numbers in the last four columns are the number of integrand evaluations needed by the four different integration methods to compute the  $p$ -value with a relative accuracy of 1%. In the empty cells, the integration was aborted after a number of evaluations which was a factor of 10 larger than the evaluations needed by slowest of the other methods.

a reasonable number of evaluations. The main reason for the poor performance of the VEGAS methods at small  $p$ -values is the fact that they require a large number of initial evaluations to find *any* points in the integration region which give a nonzero contribution to the integrand. The *myFitter* method converges faster because it “knows”, to a certain approximation, where the integrand is nonzero.

The second test compares the performance of the four integration methods in the case of models with bounded parameters. To this end, we use the theory function (38) with  $\lambda = 0$  and the input function (39). In the constrained version of the model, the parameters  $\xi_1$  and  $\xi_2$  are still fixed to  $r$  and 0, respectively. We assume again that  $\mathbf{x}_0 = (0, 0, 0, 0, 1, 1, 1)$  is the actually measured data, perform the fit with  $S_0 = S(\mathbf{x}_0)$  and vary  $r$  to change the  $p$ -value. However, in the full model the parameter  $\xi_2$  is now restricted to the interval  $[-0.25, 0.25]$ . Thus, the “effective” number of degrees of freedom of the LRT is somewhere between one and two. Consequently, we expect the  $p$ -value to lie somewhere between the results obtained from Wilks theorem with one and two degrees of freedom.

Different  $p$ -values roughly corresponding to 2, 3, 4 and 5 standard deviations were again computed with the four integration methods to a relative precision of 1%. The required number of integrand evaluations are shown in Tab. 2. Again the *myFitter* method performs best in all cases. The ‘misaligned VEGAS’ and ‘no adaptation’ methods are significantly slower even for large  $p$ -values. The performance of the ‘aligned

$p$ -value	$p$ -value (Wilks)	$myFitter$ [ $10^3$ ]	aligned VEGAS [ $10^3$ ]	misaligned VEGAS [ $10^3$ ]	no adaptation [ $10^3$ ]
$4.4 \cdot 10^{-2}$	$11.0 \cdot 10^{-2}$	30	30	110	240
$2.8 \cdot 10^{-3}$	$9.5 \cdot 10^{-3}$	30	50	460	4100
$6.3 \cdot 10^{-5}$	$27.4 \cdot 10^{-5}$	30	1500	5100	–
$5.4 \cdot 10^{-7}$	$29.0 \cdot 10^{-7}$	40	–	–	–

**Table 2:** Results of test with bounded parameters (see text). The parameter  $\xi_1$  was fixed to different values in the constrained model, leading to the  $p$ -values shown in the first column (which roughly correspond to 2, 3, 4 and 5 standard deviations). The  $p$ -values obtained by applying Wilks theorem (with two degrees of freedom) is shown in the second column. The numbers in the last four columns are the number of integrand evaluations needed by the four different integration methods to compute the  $p$ -value with a relative accuracy of 1%. In the empty cells, the integration was aborted after a number of evaluations which was a factor of 10 larger than the evaluations needed by slowest of the other methods.

VEGAS’ method is comparable at first, but drops significantly between 3 and 4 standard deviations. The reason is the same as in the previous test: for small  $p$ -values VEGAS needs a large number of initial evaluations in order to find enough points with a nonzero integrand value. Again, only  $myFitter$  is capable of computing  $p$ -values at the  $5\sigma$  level.

The final test is concerned with the case of non-nested models. The input function (for both models) is again given by (39). The first model has two parameters  $\xi_1, \xi_2$  and the theory function  $\tilde{\mathbf{x}}_1$  is defined by

$$\tilde{\mathbf{x}}_1(\boldsymbol{\xi}) = (\xi_1, \xi_2, 0, 0, \xi_1, \xi_2, 0) \quad . \quad (40)$$

The theory function  $\tilde{\mathbf{x}}_2$  of the second model is the same as  $\tilde{\mathbf{x}}$  from (38) with  $\lambda = 0$ . Obviously, neither theory manifold contains the other as a subset, so this is an example of non-nested models. We assume the actually measured data to be

$$\mathbf{x}_0 = (r, r, 1, 0, r, r, 1) \quad (41)$$

with  $r > 0$ . For sufficiently large values of  $r$  the maximum likelihood value of model 1 at  $\mathbf{x}_0$  is larger than that of model 2 and we have the interesting situation that the model with less parameters fits the measured data better than the model with more parameters. In this situation Wilks’ theorem is clearly not applicable, so we will concentrate on this case. We perform a LRT which compares model 2 with the ‘union’ of models 1 and 2, using the test statistic  $S_2$  from (13) and  $S_0 = S_2(\mathbf{x}_0)$ .

As before, several LRTs were performed with different values of  $r$  leading to  $p$ -values roughly corresponding to 2, 3, 4 and 5 standard deviations. The  $p$ -values were again computed with the four integration methods to a relative precision of 1%, and the

$p$ -value	$p$ -value (Wilks)	$myFitter$ [ $10^3$ ]	aligned VEGAS [ $10^3$ ]	misaligned VEGAS [ $10^3$ ]	no adaptation [ $10^3$ ]
$5.0 \cdot 10^{-2}$	$50.3 \cdot 10^{-2}$	50	70	100	200
$2.7 \cdot 10^{-3}$	$44.6 \cdot 10^{-3}$	60	210	270	4000
$6.7 \cdot 10^{-5}$	$147.0 \cdot 10^{-5}$	70	330	540	–
$5.7 \cdot 10^{-7}$	$157.3 \cdot 10^{-7}$	80	–	–	–

**Table 3:** Results of test for non-nested models. The  $p$ -values in the first column were obtained by numerical integration with different values for  $r$  (see text). The  $p$ -values in the second column were computed by applying Wilks’ theorem with two degrees of freedom. The numbers in the last four columns are the number of integrand evaluations needed by the four different integration methods to compute the  $p$ -value with a relative accuracy of 1%. In the empty cells, the integration was aborted after a number of evaluations which was a factor of 10 larger than the evaluations needed by slowest of the other methods.

required number of integrand evaluations are shown in Tab. 3. The  $p$ -values obtained by applying Wilks’ theorem with two degrees of freedom are also shown for illustration. We see that Wilks’ theorem is clearly not applicable here. As in the previous tests, the  $myFitter$  method is consistently faster than the other methods and significantly faster for small  $p$ -values.

## 7 Conclusions

Likelihood ratio tests are a popular tool in global analyses of models in particle physics. For a correct statistical interpretation of the data, reliable methods for the computation of  $p$ -values in likelihood ratio tests are needed. There are many realistic situations where Wilks’ theorem does not apply and the distribution of the test statistic is not known analytically. These include likelihood ratio tests of non-nested models or models with parameters that are only allowed to float in a finite range. Real-world examples of the former case are the global analyses [1, 3] of the Standard Model with a fourth generation of fermions where the models being compared are not nested due to the non-decoupling nature of the additional fermions. The latter case includes models where systematic errors are treated within the  $RFit$  scheme. In these situations one has to resort to numerical methods. Monte Carlo integration can be used to compute  $p$ -values numerically, but the integration usually becomes very inefficient for small  $p$ -values.

In this paper I presented an efficient approach to the numerical computation of  $p$ -values which is based on importance sampling and applies to a broad class of statistical models. In global analyses in particle physics, the predictions of a theoretical model can be described by a manifold in the space of observables. The PDF of the statistical

model is then obtained by “smearing out” the theory manifold in a way determined by the experimental uncertainties. The proposed methods use geometric information about the theory manifolds to construct suitable sampling densities for the Monte Carlo integration and substantially improve the performance of the numerical integration for small  $p$ -values. These methods are implemented in a publicly available C++ framework for likelihood ratio tests called *myFitter*.

## Acknowledgements

I would like to thank Jérôme Charles for fruitful discussions about likelihood ratio tests for non-nested models. I also thank Otto Eberhardt for checking fit results with CKMfitter and Ulrich Nierste for thorough proof reading.

## References

- [1] O. Eberhardt, G. Herbert, H. Lacker, A. Lenz, A. Menzel, U. Nierste, and M. Wiebusch, *Phys.Rev.Lett.* **109**, 241802 (2012), [arXiv:1209.1101 \[hep-ph\]](#).
- [2] O. Eberhardt, G. Herbert, H. Lacker, A. Lenz, A. Menzel, U. Nierste, and M. Wiebusch, *Phys.Rev.* **D86**, 013011 (2012), [arXiv:1204.3872 \[hep-ph\]](#).
- [3] O. Eberhardt, A. Lenz, A. Menzel, U. Nierste, and M. Wiebusch, *Phys.Rev.* **D86**, 074014 (2012), [arXiv:1207.0438 \[hep-ph\]](#).
- [4] G. Casella and R. L. Berger, *Statistical Inference*, 2nd ed. (Duxbury Press, 2001).
- [5] K. Nakamura *et al.* (Particle Data Group), *J.Phys.G* **G37**, 075021 (2010).
- [6] S. S. Wilks, *Ann. Math. Statist.* **9**, 60 (1938).
- [7] A. Höcker, H. Lacker, S. Laplace, and F. Le Diberder, *Eur.Phys.J.* **C21**, 225 (2001), [arXiv:hep-ph/0104062 \[hep-ph\]](#).
- [8] <http://myfitter.hepforge.org>.
- [9] <http://dvegas.hepforge.org>.
- [10] G. P. Lepage, *Journal of Computational Physics* **27**, 192 (1978), ISSN 0021-9991.
- [11] N. Kauer and D. Zeppenfeld, *Phys.Rev.* **D65**, 014021 (2002), [arXiv:hep-ph/0107181 \[hep-ph\]](#).
- [12] N. Kauer, *Phys.Rev.* **D67**, 054013 (2003), [arXiv:hep-ph/0212091 \[hep-ph\]](#).

- [13] C. G. Broyden, *IMA J. Appl. Math.* **6**, 76 (1970).
- [14] R. Fletcher, *Comput. J.* **13**, 317 (1970).
- [15] D. Goldfarb, *Math. Comput.* **24**, 23 (1970).
- [16] D. F. Shanno, *Math. Comput.* **24**, 647 (1970).